
PyTorch Lightning Documentation

Release 1.0.2

William Falcon et al.

Oct 15, 2020

GETTING STARTED

1	Lightning in 2 steps	1
2	How to organize PyTorch into Lightning	15
3	Rapid prototyping templates	19
4	Style guide	21
5	Fast performance tips	27
6	LightningModule	31
7	Trainer	77
8	Callback	105
9	LightningDataModule	125
10	Logging	133
11	Metrics	153
12	Step-by-step walk-through	183
13	Bolts	211
14	Community Examples	213
15	AWS/GCP training	215
16	16-bit training	217
17	Computing cluster (SLURM)	219
18	Child Modules	223
19	Debugging	225
20	Loggers	229
21	Early stopping	233
22	Fast Training	235

23	Hyperparameters	237
24	Learning Rate Finder	243
25	Multi-GPU training	247
26	Multiple Datasets	259
27	Saving and loading weights	261
28	Optimization	265
29	Performance and Bottleneck Profiler	271
30	Single GPU Training	277
31	Sequential Data	279
32	Training Tricks	281
33	Transfer Learning	285
34	TPU support	289
35	Test set	293
36	Inference in Production	295
37	ASR & TTS	297
38	Contributor Covenant Code of Conduct	311
39	Contributing	313
40	How to become a core contributor	321
41	PyTorch Lightning Governance Persons of interest	323
42	Changelog	325
43	Indices and tables	359
	Index	361

LIGHTNING IN 2 STEPS

In this guide we'll show you how to organize your PyTorch code into Lightning in 2 steps.

Organizing your code with PyTorch Lightning makes your code:

- Keep all the flexibility (this is all pure PyTorch), but removes a ton of boilerplate
- More readable by decoupling the research code from the engineering
- Easier to reproduce
- Less error prone by automating most of the training loop and tricky engineering
- Scalable to any hardware without changing your model

Here's a 3 minute conversion guide for PyTorch projects:

1.1 Step 0: Install PyTorch Lightning

You can install using `pip`

```
pip install pytorch-lightning
```

Or with `conda` (see how to install conda [here](#)):

```
conda install pytorch-lightning -c conda-forge
```

You could also use conda environments

```
conda activate my_env
pip install pytorch-lightning
```

Import the following:

```
import os
import torch
from torch import nn
import torch.nn.functional as F
from torchvision.datasets import MNIST
from torchvision import transforms
```

(continues on next page)

(continued from previous page)

```
from torch.utils.data import DataLoader
import pytorch_lightning as pl
from torch.utils.data import random_split
```

1.2 Step 1: Define LightningModule

```
class LitAutoEncoder(pl.LightningModule):

    def __init__(self):
        super().__init__()
        self.encoder = nn.Sequential(
            nn.Linear(28*28, 64),
            nn.ReLU(),
            nn.Linear(64, 3)
        )
        self.decoder = nn.Sequential(
            nn.Linear(3, 64),
            nn.ReLU(),
            nn.Linear(64, 28*28)
        )

    def forward(self, x):
        # in lightning, forward defines the prediction/inference actions
        embedding = self.encoder(x)
        return embedding

    def training_step(self, batch, batch_idx):
        # training_step defined the train loop.
        # It is independent of forward
        x, y = batch
        x = x.view(x.size(0), -1)
        z = self.encoder(x)
        x_hat = self.decoder(z)
        loss = F.mse_loss(x_hat, x)
        # Logging to TensorBoard by default
        self.log('train_loss', loss)
        return loss

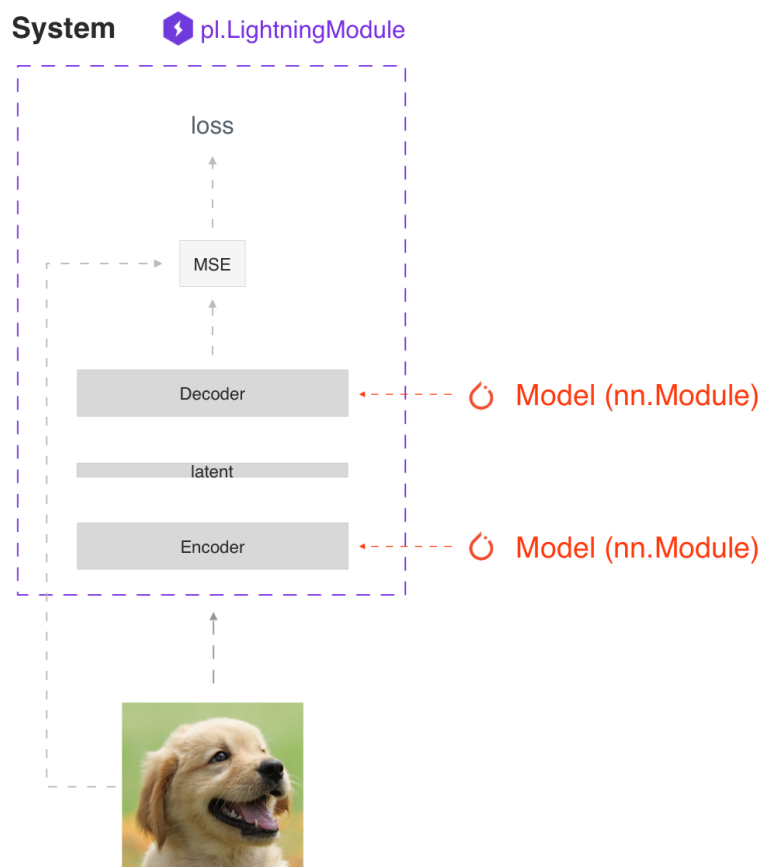
    def configure_optimizers(self):
        optimizer = torch.optim.Adam(self.parameters(), lr=1e-3)
        return optimizer
```

SYSTEM VS MODEL

A *LightningModule* defines a *system* not a model.

Examples of systems are:

- Autoencoder
- BERT
- DQN
- GAN
- Image classifier



- Seq2seq
- [SimCLR](#)
- [VAE](#)

Under the hood a `LightningModule` is still just a `torch.nn.Module` that groups all research code into a single file to make it self-contained:

- The Train loop
- The Validation loop
- The Test loop
- The Model or system of Models
- The Optimizer

You can customize any part of training (such as the backward pass) by overriding any of the 20+ hooks found in [Available Callback hooks](#)

```
class LitAutoEncoder(pl.LightningModule):  
  
    def backward(self, loss, optimizer, optimizer_idx):  
        loss.backward()
```

FORWARD vs TRAINING_STEP

In Lightning we separate training from inference. The `training_step` defines the full training loop. We encourage users to use the forward to define inference actions.

For example, in this case we could define the autoencoder to act as an embedding extractor:

```
def forward(self, x):  
    embeddings = self.encoder(x)  
    return embeddings
```

Of course, nothing is stopping you from using forward from within the `training_step`

```
def training_step(self, batch, batch_idx):  
    ...  
    z = self(x)
```

It really comes down to your application. We do however, recommend that you keep both intents separate.

- Use forward for inference (predicting).
- Use `training_step` for training.

More details in [LightningModule](#) docs.

1.3 Step 2: Fit with Lightning Trainer

First, define the data however you want. Lightning just needs a `DataLoader` for the train/val/test splits.

```
dataset = MNIST(os.getcwd(), download=True, transform=transforms.ToTensor())
train_loader = DataLoader(dataset)
```

Next, init the *LightningModule* and the PyTorch Lightning Trainer, then call fit with both the data and model.

```
# init model
autoencoder = LitAutoEncoder()

# most basic trainer, uses good defaults (auto-tensorboard, checkpoints, logs, and_
↪ more)
# trainer = pl.Trainer(gpus=8) (if you have GPUs)
trainer = pl.Trainer()
trainer.fit(autoencoder, train_loader)
```

The Trainer automates:

- Epoch and batch iteration
- Calling of `optimizer.step()`, `backward`, `zero_grad()`
- Calling of `.eval()`, enabling/disabling grads
- *Saving and loading weights*
- Tensorboard (see *Loggers* options)
- *Multi-GPU training* support
- *TPU support*
- *16-bit training* support

Tip: If you prefer to manually manage optimizers you can use the *Manual optimization* mode (ie: RL, GANs, etc...).

That's it!

These are the main 2 concepts you need to know in Lightning. All the other features of lightning are either features of the Trainer or LightningModule.

1.4 Basic features

1.4.1 Manual vs automatic optimization

Automatic optimization

With Lightning you don't need to worry about when to enable/disable grads, do a backward pass, or update optimizers as long as you return a loss with an attached graph from the *training_step*, Lightning will automate the optimization.

```
def training_step(self, batch, batch_idx):
    loss = self.encoder(batch[0])
    return loss
```

Manual optimization

However, for certain research like GANs, reinforcement learning or something with multiple optimizers or an inner loop, you can turn off automatic optimization and fully control the training loop yourself.

First, turn off automatic optimization:

```
trainer = Trainer(automatic_optimization=False)
```

Now you own the train loop!

```
def training_step(self, batch, batch_idx, opt_idx):
    (opt_a, opt_b, opt_c) = self.optimizers()

    loss_a = self.generator(batch[0])

    # use this instead of loss.backward so we can automate half precision, etc...
    self.manual_backward(loss_a, opt_a, retain_graph=True)
    self.manual_backward(loss_a, opt_a)
    opt_a.step()
    opt_a.zero_grad()

    loss_b = self.discriminator(batch[0])
    self.manual_backward(loss_b, opt_b)
    ...
```

1.4.2 Predict or Deploy

When you're done training, you have 3 options to use your LightningModule for predictions.

Option 1: Sub-models

Pull out any model inside your system for predictions.

```
# -----
# to use as embedding extractor
# -----
autoencoder = LitAutoEncoder.load_from_checkpoint('path/to/checkpoint_file.ckpt')
encoder_model = autoencoder.encoder
encoder_model.eval()

# -----
# to use as image generator
# -----
decoder_model = autoencoder.decoder
decoder_model.eval()
```

Option 2: Forward

You can also add a forward method to do predictions however you want.

```
# -----
# using the AE to extract embeddings
# -----
class LitAutoEncoder(pl.LightningModule):
    def forward(self, x):
        embedding = self.encoder(x)
        return embedding

autoencoder = LitAutoEncoder()
autoencoder = autoencoder(torch.rand(1, 28 * 28))
```

```
# -----
# or using the AE to generate images
# -----
class LitAutoEncoder(pl.LightningModule):
    def forward(self):
        z = torch.rand(1, 3)
        image = self.decoder(z)
        image = image.view(1, 1, 28, 28)
        return image

autoencoder = LitAutoEncoder()
image_sample = autoencoder()
```

Option 3: Production

For production systems onnx or torchscript are much faster. Make sure you have added a forward method or trace only the sub-models you need.

```
# -----
# torchscript
# -----
autoencoder = LitAutoEncoder()
torch.jit.save(autoencoder.to_torchscript(), "model.pt")
os.path.isfile("model.pt")
```

```
# -----
# onnx
# -----
with tempfile.NamedTemporaryFile(suffix='.onnx', delete=False) as tmpfile:
    autoencoder = LitAutoEncoder()
    input_sample = torch.randn((1, 28 * 28))
    autoencoder.to_onnx(tmpfile.name, input_sample, export_params=True)
    os.path.isfile(tmpfile.name)
```

1.4.3 Using CPUs/GPUs/TPUs

It's trivial to use CPUs, GPUs or TPUs in Lightning. There's **NO NEED** to change your code, simply change the Trainer options.

```
# train on CPU
trainer = pl.Trainer()
```

```
# train on 8 CPUs
trainer = pl.Trainer(num_processes=8)
```

```
# train on 1024 CPUs across 128 machines
trainer = pl.Trainer(
    num_processes=8,
    num_nodes=128
)
```

```
# train on 1 GPU
trainer = pl.Trainer(gpus=1)
```

```
# train on multiple GPUs across nodes (32 gpus here)
trainer = pl.Trainer(
    gpus=4,
    num_nodes=8
)
```

```
# train on gpu 1, 3, 5 (3 gpus total)
trainer = pl.Trainer(gpus=[1, 3, 5])
```

```
# Multi GPU with mixed precision
trainer = pl.Trainer(gpus=2, precision=16)
```

```
# Train on TPUs
trainer = pl.Trainer(tpu_cores=8)
```

Without changing a SINGLE line of your code, you can now do the following with the above code:

```
# train on TPUs using 16 bit precision
# using only half the training data and checking validation every quarter of a
↪training epoch
trainer = pl.Trainer(
    tpu_cores=8,
    precision=16,
    limit_train_batches=0.5,
    val_check_interval=0.25
)
```

1.4.4 Checkpoints

Lightning automatically saves your model. Once you've trained, you can load the checkpoints as follows:

```
model = LitModel.load_from_checkpoint(path)
```

The above checkpoint contains all the arguments needed to init the model and set the state dict. If you prefer to do it manually, here's the equivalent

```
# load the ckpt
ckpt = torch.load('path/to/checkpoint.ckpt')

# equivalent to the above
model = LitModel()
model.load_state_dict(ckpt['state_dict'])
```

1.4.5 Data flow

Each loop (training, validation, test) has three hooks you can implement:

- `x_step`
- `x_step_end`
- `x_epoch_end`

To illustrate how data flows, we'll use the training loop (ie: `x=training`)

```
outs = []
for batch in data:
    out = training_step(batch)
    outs.append(out)
training_epoch_end(outs)
```

The equivalent in Lightning is:

```
def training_step(self, batch, batch_idx):
    prediction = ...
    return prediction

def training_epoch_end(self, training_step_outputs):
    for prediction in predictions:
        # do something with these
```

In the event that you use DP or DDP2 distributed modes (ie: split a batch across GPUs), use the `x_step_end` to manually aggregate (or don't implement it to let lightning auto-aggregate for you).

```
for batch in data:
    model_copies = copy_model_per_gpu(model, num_gpus)
    batch_split = split_batch_per_gpu(batch, num_gpus)

    gpu_outs = []
    for model, batch_part in zip(model_copies, batch_split):
        # LightningModule hook
        gpu_out = model.training_step(batch_part)
```

(continues on next page)

(continued from previous page)

```
gpu_outs.append(gpu_out)

# LightningModule hook
out = training_step_end(gpu_outs)
```

The lightning equivalent is:

```
def training_step(self, batch, batch_idx):
    loss = ...
    return loss

def training_step_end(self, losses):
    gpu_0_loss = losses[0]
    gpu_1_loss = losses[1]
    return (gpu_0_loss + gpu_1_loss) * 1/2
```

Tip: The validation and test loops have the same structure.

1.4.6 Logging

To log to Tensorboard, your favorite logger, and/or the progress bar, use the `log()` method which can be called from any method in the `LightningModule`.

```
def training_step(self, batch, batch_idx):
    self.log('my_metric', x)
```

The `log()` method has a few options:

- `on_step` (logs the metric at that step in training)
- `on_epoch` (automatically accumulates and logs at the end of the epoch)
- `prog_bar` (logs to the progress bar)
- `logger` (logs to the logger like Tensorboard)

Depending on where `log` is called from, Lightning auto-determines the correct mode for you. But of course you can override the default behavior by manually setting the flags

Note: Setting `on_epoch=True` will accumulate your logged values over the full training epoch.

```
def training_step(self, batch, batch_idx):
    self.log('my_loss', loss, on_step=True, on_epoch=True, prog_bar=True, logger=True)
```

Note: The loss value shown in the progress bar is smoothed (averaged) over the last values, so it differs from the actual loss returned in train/validation step.

You can also use any method of your logger directly:

```
def training_step(self, batch, batch_idx):
    tensorboard = self.logger.experiment
    tensorboard.any_summary_writer_method_you_want()
```

Once your training starts, you can view the logs by using your favorite logger or booting up the Tensorboard logs:

```
tensorboard --logdir ./lightning_logs
```

Note: Lightning automatically shows the loss value returned from `training_step` in the progress bar. So, no need to explicitly log like this `self.log('loss', loss, prog_bar=True)`.

Read more about [Loggers](#).

1.4.7 Optional extensions

Callbacks

A callback is an arbitrary self-contained program that can be executed at arbitrary parts of the training loop.

Here's an example adding a not-so-fancy learning rate decay rule:

```
class DecayLearningRate(pl.Callback)

    def __init__(self):
        self.old_lrs = []

    def on_train_start(self, trainer, pl_module):
        # track the initial learning rates
        for opt_idx in optimizer in enumerate(trainer.optimizers):
            group = []
            for param_group in optimizer.param_groups:
                group.append(param_group['lr'])
            self.old_lrs.append(group)

    def on_train_epoch_end(self, trainer, pl_module, outputs):
        for opt_idx in optimizer in enumerate(trainer.optimizers):
            old_lr_group = self.old_lrs[opt_idx]
            new_lr_group = []
            for p_idx, param_group in enumerate(optimizer.param_groups):
                old_lr = old_lr_group[p_idx]
                new_lr = old_lr * 0.98
                new_lr_group.append(new_lr)
                param_group['lr'] = new_lr
            self.old_lrs[opt_idx] = new_lr_group
```

Things you can do with a callback:

- Send emails at some point in training
- Grow the model
- Update learning rates
- Visualize gradients

- ...
- You are only limited by your imagination

Learn more about custom callbacks.

LightningDataModules

DataLoaders and data processing code tends to end up scattered around. Make your data code reusable by organizing it into a LightningDataModule.

```
class MNISTDataModule(pl.LightningDataModule):

    def __init__(self, batch_size=32):
        super().__init__()
        self.batch_size = batch_size

    # When doing distributed training, Datamodules have two optional arguments for
    # granular control over download/prepare/splitting data:

    # OPTIONAL, called only on 1 GPU/machine
    def prepare_data(self):
        MNIST(os.getcwd(), train=True, download=True)
        MNIST(os.getcwd(), train=False, download=True)

    # OPTIONAL, called for every GPU/machine (assigning state is OK)
    def setup(self, stage):
        # transforms
        transform=transforms.Compose([
            transforms.ToTensor(),
            transforms.Normalize((0.1307,), (0.3081,))
        ])
        # split dataset
        if stage == 'fit':
            mnist_train = MNIST(os.getcwd(), train=True, transform=transform)
            self.mnist_train, self.mnist_val = random_split(mnist_train, [55000, ↪
5000])
        if stage == 'test':
            self.mnist_test = MNIST(os.getcwd(), train=False, transform=transform)

    # return the dataloader for each split
    def train_dataloader(self):
        mnist_train = DataLoader(self.mnist_train, batch_size=self.batch_size)
        return mnist_train

    def val_dataloader(self):
        mnist_val = DataLoader(self.mnist_val, batch_size=self.batch_size)
        return mnist_val

    def test_dataloader(self):
        mnist_test = DataLoader(self.mnist_test, batch_size=self.batch_size)
        return mnist_test
```

LightningDataModule is designed to enable sharing and reusing data splits and transforms across different projects. It encapsulates all the steps needed to process data: downloading, tokenizing, processing etc.

Now you can simply pass your LightningDataModule to the Trainer:


```
# init model
model = LitModel()

# init data
dm = MNISTDataModule()

# train
trainer = pl.Trainer()
trainer.fit(model, dm)

# test
trainer.test(datamodule=dm)
```

DataModules are specifically useful for building models based on data. Read more on [LightningDataModule](#).

1.4.8 Debugging

Lightning has many tools for debugging. Here is an example of just a few of them:

```
# use only 10 train batches and 3 val batches
trainer = pl.Trainer(limit_train_batches=10, limit_val_batches=3)
```

```
# Automatically overfit the same batch of your model for a sanity test
trainer = pl.Trainer(overfit_batches=1)
```

```
# unit test all the code- hits every line of your code once to see if you have bugs,
# instead of waiting hours to crash on validation
trainer = pl.Trainer(fast_dev_run=True)
```

```
# train only 20% of an epoch
trainer = pl.Trainer(limit_train_batches=0.2)
```

```
# run validation every 25% of a training epoch
trainer = pl.Trainer(val_check_interval=0.25)
```

```
# Profile your code to find speed/memory bottlenecks
Trainer(profile=True)
```

1.5 Other cool features

Once you define and train your first Lightning model, you might want to try other cool features like

- *Automatic early stopping*
- *Automatic truncated-back-propagation-through-time*
- *Automatically scale your batch size*
- *Automatically find a good learning rate*

- *Load checkpoints directly from S3*
- *Scale to massive compute clusters*
- *Use multiple dataloaders per train/val/test loop*
- *Use multiple optimizers to do reinforcement learning or even GANs*

Or read our *Step-by-step walk-through* to learn more!

1.5.1 Grid AI

Grid AI is our native solution for large scale training and tuning on the cloud provider of your choice.

[Click here](#) to request early-access.

1.6 Community

Our community of core maintainers and thousands of expert researchers is active on our [Slack](#) and [Forum](#). Drop by to hang out, ask Lightning questions or even discuss research!

1.6.1 Masterclass

We also offer a Masterclass to teach you the advanced uses of Lightning.



HOW TO ORGANIZE PYTORCH INTO LIGHTNING

To enable your code to work with Lightning, here's how to organize PyTorch into Lightning

2.1 1. Move your computational code

Move the model architecture and forward pass to your *LightningModule*.

```
class LitModel(LightningModule):  
  
    def __init__(self):  
        super().__init__()   
        self.layer_1 = torch.nn.Linear(28 * 28, 128)  
        self.layer_2 = torch.nn.Linear(128, 10)  
  
    def forward(self, x):  
        x = x.view(x.size(0), -1)  
        x = self.layer_1(x)  
        x = F.relu(x)  
        x = self.layer_2(x)  
        return x
```

2.2 2. Move the optimizer(s) and schedulers

Move your optimizers to the `configure_optimizers()` hook.

```
class LitModel(LightningModule):  
  
    def configure_optimizers(self):  
        optimizer = torch.optim.Adam(self.parameters(), lr=1e-3)  
        return optimizer
```

2.3 3. Find the train loop “meat”

Lightning automates most of the training for you, the epoch and batch iterations, all you need to keep is the training step logic. This should go into the `training_step()` hook (make sure to use the hook parameters, `batch` and `batch_idx` in this case):

```
class LitModel(LightningModule):  
  
    def training_step(self, batch, batch_idx):  
        x, y = batch  
        y_hat = self(x)  
        loss = F.cross_entropy(y_hat, y)  
        return loss
```

2.4 4. Find the val loop “meat”

To add an (optional) validation loop add logic to the `validation_step()` hook (make sure to use the hook parameters, `batch` and `batch_idx` in this case).

```
class LitModel(LightningModule):  
  
    def validation_step(self, batch, batch_idx):  
        x, y = batch  
        y_hat = self(x)  
        val_loss = F.cross_entropy(y_hat, y)  
        return val_loss
```

Note: `model.eval()` and `torch.no_grad()` are called automatically for validation

2.5 5. Find the test loop “meat”

To add an (optional) test loop add logic to the `test_step()` hook (make sure to use the hook parameters, `batch` and `batch_idx` in this case).

```
class LitModel(LightningModule):  
  
    def test_step(self, batch, batch_idx):  
        x, y = batch  
        y_hat = self(x)  
        loss = F.cross_entropy(y_hat, y)  
        return loss
```

Note: `model.eval()` and `torch.no_grad()` are called automatically for testing.

The test loop will not be used until you call.

```
trainer.test()
```

Tip: `.test()` loads the best checkpoint automatically

2.6 6. Remove any `.cuda()` or `to.device()` calls

Your *LightningModule* can automatically run on any hardware!

RAPID PROTOTYPING TEMPLATES

Use these templates for rapid prototyping

3.1 General Use

Use case	Description	link
Scratch model	To prototype quickly / debug with random data	
Scratch model with manual optimization	To prototype quickly / debug with random data	

STYLE GUIDE

A main goal of Lightning is to improve readability and reproducibility. Imagine looking into any GitHub repo, finding a lightning module and knowing exactly where to look to find the things you care about.

The goal of this style guide is to encourage Lightning code to be structured similarly.

4.1 LightningModule

These are best practices about structuring your LightningModule

4.1.1 Systems vs models

The main principle behind a LightningModule is that a full system should be self-contained. In Lightning we differentiate between a system and a model.

A model is something like a resnet18, RNN, etc.

A system defines how a collection of models interact with each other. Examples of this are:

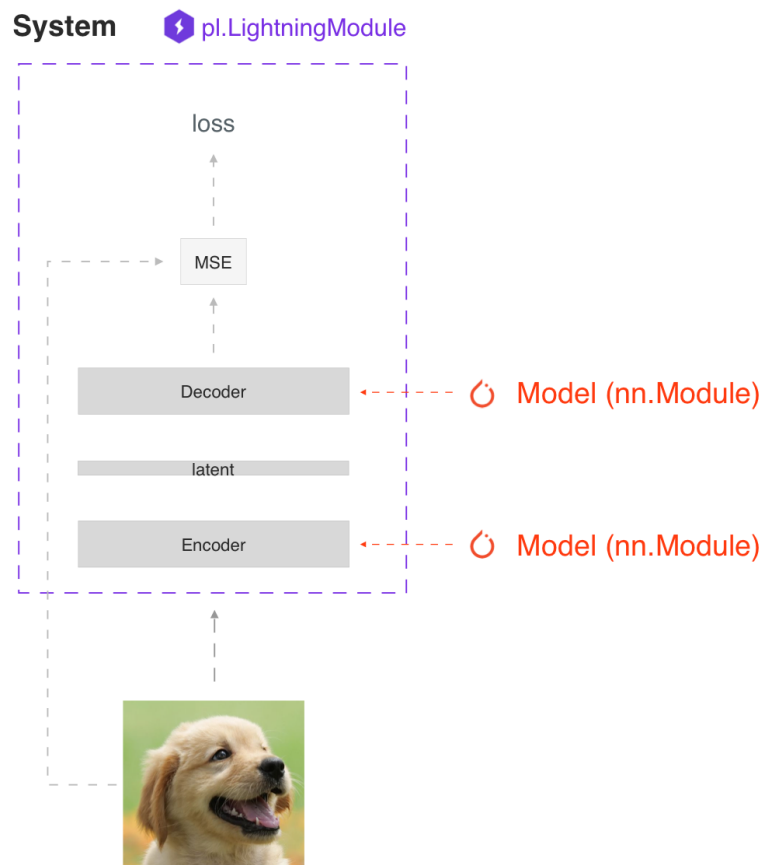
- GANs
- Seq2Seq
- BERT
- etc

A LightningModule can define both a system and a model.

Here's a LightningModule that defines a model:

```
class LitModel(pl.LightningModule):
    def __init__(self, num_layers: int = 3)
        super().__init__()
        self.layer_1 = nn.Linear(...)
        self.layer_2 = nn.Linear(...)
        self.layer_3 = nn.Linear(...)
```

Here's a lightningModule that defines a system:



```
class LitModel(pl.LightningModule):
    def __init__(self, encoder: nn.Module = None, decoder: nn.Module = None):
        super().__init__()
        self.encoder = encoder
        self.decoder = decoder
```

For fast prototyping it's often useful to define all the computations in a `LightningModule`. For reusability and scalability it might be better to pass in the relevant backbones.

4.1.2 Self-contained

A Lightning module should be self-contained. A good test to see how self-contained your model is, is to ask yourself this question:

“Can someone drop this file into a Trainer without knowing anything about the internals?”

For example, we couple the optimizer with a model because the majority of models require a specific optimizer with a specific learning rate scheduler to work well.

4.1.3 Init

The first place where `LightningModules` tend to stop being self-contained is in the `init`. Try to define all the relevant sensible defaults in the `init` so that the user doesn't have to guess.

Here's an example where a user will have to go hunt through files to figure out how to `init` this `LightningModule`.

```
class LitModel(pl.LightningModule):
    def __init__(self, params):
        self.lr = params.lr
        self.coef_x = params.coef_x
```

Models defined as such leave you with many questions; what is `coef_x`? is it a string? a float? what is the range? etc...

Instead, be explicit in your `init`

```
class LitModel(pl.LightningModule):
    def __init__(self, encoder: nn.Module, coeff_x: float = 0.2, lr: float = 1e-3)
```

Now the user doesn't have to guess. Instead they know the value type and the model has a sensible default where the user can see the value immediately.

4.1.4 Method order

The only required methods in the `LightningModule` are:

- `init`
- `training_step`
- `configure_optimizers`

However, if you decide to implement the rest of the optional methods, the recommended order is:

- model/system definition (`init`)
- if doing inference, define `forward`

- training hooks
- validation hooks
- test hooks
- `configure_optimizers`
- any other hooks

In practice, this code looks like:

```
class LitModel(pl.LightningModule):  
  
    def __init__(...):  
  
    def forward(...):  
  
    def training_step(...)  
  
    def training_step_end(...)  
  
    def training_epoch_end(...)  
  
    def validation_step(...)  
  
    def validation_step_end(...)  
  
    def validation_epoch_end(...)  
  
    def test_step(...)  
  
    def test_step_end(...)  
  
    def test_epoch_end(...)  
  
    def configure_optimizers(...)  
  
    def any_extra_hook(...)
```

4.1.5 Forward vs `training_step`

We recommend using `forward` for inference/predictions and keeping `training_step` independent

```
def forward(...):  
    embeddings = self.encoder(x)  
  
def training_step(...):  
    x, y = ...  
    z = self.encoder(x)  
    pred = self.decoder(z)  
    ...
```

However, when using `DataParallel`, you will need to call `forward` manually

```
def training_step(...):  
    x, y = ...  
    z = self(x)  # < ----- instead of self.encoder(x)
```

(continues on next page)

(continued from previous page)

```
pred = self.decoder(z)
...
```

4.2 Data

These are best practices for handling data.

4.2.1 Dataloaders

Lightning uses dataloaders to handle all the data flow through the system. Whenever you structure dataloaders, make sure to tune the number of workers for maximum efficiency.

Warning: Make sure not to use `ddp_spawn` with `num_workers > 0` or you will bottleneck your code.

4.2.2 DataModules

Lightning introduced datamodules. The problem with dataloaders is that sharing full datasets is often still challenging because all these questions need to be answered:

- What splits were used?
- How many samples does this dataset have?
- What transforms were used?
- etc...

It's for this reason that we recommend you use datamodules. This is specially important when collaborating because it will save your team a lot of time as well.

All they need to do is drop a datamodule into a lightning trainer and not worry about what was done to the data.

This is true for both academic and corporate settings where data cleaning and ad-hoc instructions slow down the progress of iterating through ideas.

FAST PERFORMANCE TIPS

Lightning builds in all the micro-optimizations we can find to increase your performance. But we can only automate so much.

Here are some additional things you can do to increase your performance.

5.1 Dataloaders

When building your `DataLoader` set `num_workers > 0` and `pin_memory=True` (only for GPUs).

```
Dataloader(dataset, num_workers=8, pin_memory=True)
```

5.1.1 num_workers

The question of how many `num_workers` is tricky. Here's a summary of some references, [1], and our suggestions.

1. `num_workers=0` means ONLY the main process will load batches (that can be a bottleneck).
2. `num_workers=1` means ONLY one worker (just not the main process) will load data but it will still be slow.
3. The `num_workers` depends on the batch size and your machine.
4. A general place to start is to set `num_workers` equal to the number of CPUs on that machine.

Warning: Increasing `num_workers` will ALSO increase your CPU memory consumption.

The best thing to do is to increase the `num_workers` slowly and stop once you see no more improvement in your training speed.

5.1.2 Spawn

When using `distributed_backend=ddp_spawn` (the `ddp` default) or TPU training, the way multiple GPUs/TPU cores are used is by calling `.spawn()` under the hood. The problem is that PyTorch has issues with `num_workers > 0` when using `.spawn()`. For this reason we recommend you use `distributed_backend=ddp` so you can increase the `num_workers`, however your script has to be callable like so:

```
python my_program.py --gpus X
```

5.2 .item(), .numpy(), .cpu()

Don't call `.item()` anywhere in your code. Use `.detach()` instead to remove the connected graph calls. Lightning takes a great deal of care to be optimized for this.

5.3 empty_cache()

Don't call this unnecessarily! Every time you call this ALL your GPUs have to wait to sync.

5.4 Construct tensors directly on the device

LightningModules know what device they are on! Construct tensors on the device directly to avoid CPU->Device transfer.

```
# bad
t = torch.rand(2, 2).cuda()

# good (self is LightningModule)
t = torch.rand(2, 2, device=self.device)
```

For tensors that need to be model attributes, it is best practice to register them as buffers in the modules's `__init__` method:

```
# bad
self.t = torch.rand(2, 2, device=self.device)

# good
self.register_buffer("t", torch.rand(2, 2))
```

5.5 Use DDP not DP

DP performs three GPU transfers for EVERY batch:

1. Copy model to device.
2. Copy data to device.
3. Copy outputs of each device back to master.

Whereas DDP only performs 1 transfer to sync gradients. Because of this, DDP is MUCH faster than DP.

5.6 16-bit precision

Use 16-bit to decrease the memory consumption (and thus increase your batch size). On certain GPUs (V100s, 2080tis), 16-bit calculations are also faster. However, know that 16-bit and multi-processing (any DDP) can have issues. Here are some common problems.

1. **CUDA error: an illegal memory access was encountered.** The solution is likely setting a specific CUDA, CUDNN, PyTorch version combination.
2. CUDA error: device-side assert triggered. This is a general catch-all error. To see the actual error run your script like so:

```
# won't see what the error is
python main.py

# will see what the error is
CUDA_LAUNCH_BLOCKING=1 python main.py
```

Tip: We also recommend using 16-bit native found in PyTorch 1.6. Just install this version and Lightning will automatically use it.

LIGHTNINGMODULE

A `LightningModule` organizes your PyTorch code into 5 sections

- Computations (`init`).
- Train loop (`training_step`)
- Validation loop (`validation_step`)
- Test loop (`test_step`)
- Optimizers (`configure_optimizers`)

Notice a few things.

1. It's the SAME code.
2. The PyTorch code IS NOT abstracted - just organized.
3. All the other code that's not in the `LightningModule` has been automated for you by the trainer.

```
net = Net()
trainer = Trainer()
trainer.fit(net)
```

4. There are no `.cuda()` or `.to()` calls... Lightning does these for you.

```
# don't do in lightning
x = torch.Tensor(2, 3)
x = x.cuda()
x = x.to(device)
```

(continues on next page)

(continued from previous page)

```
# do this instead
x = x # leave it alone!

# or to init a new tensor
new_x = torch.Tensor(2, 3)
new_x = new_x.type_as(x)
```

5. There are no samplers for distributed, Lightning also does this for you.

```
# Don't do in Lightning...
data = MNIST(...)
sampler = DistributedSampler(data)
DataLoader(data, sampler=sampler)

# do this instead
data = MNIST(...)
DataLoader(data)
```

6. A `LightningModule` is a `torch.nn.Module` but with added functionality. Use it as such!

```
net = Net.load_from_checkpoint(PATH)
net.freeze()
out = net(x)
```

Thus, to use Lightning, you just need to organize your code which takes about 30 minutes, (and let's be real, you probably should do anyhow).

6.1 Minimal Example

Here are the only required methods.

```
>>> import pytorch_lightning as pl
>>> class LitModel(pl.LightningModule):
...     def __init__(self):
...         super().__init__()
...         self.l1 = torch.nn.Linear(28 * 28, 10)
...
...     def forward(self, x):
...         return torch.relu(self.l1(x.view(x.size(0), -1)))
...
...     def training_step(self, batch, batch_idx):
...         x, y = batch
```

(continues on next page)

(continued from previous page)

```

...         y_hat = self(x)
...         loss = F.cross_entropy(y_hat, y)
...         return loss
...
...     def configure_optimizers(self):
...         return torch.optim.Adam(self.parameters(), lr=0.02)

```

Which you can train by doing:

```

train_loader = DataLoader(MNIST(os.getcwd(), download=True, transform=transforms.
    ↳ToTensor()))
trainer = pl.Trainer()
model = LitModel()

trainer.fit(model, train_loader)

```

The LightningModule has many convenience methods, but the core ones you need to know about are:

Name	Description
init	Define computations here
forward	Use for inference only (separate from training_step)
training_step	the full training loop
validation_step	the full validation loop
test_step	the full test loop
configure_optimizers	define optimizers and LR schedulers

6.2 Training

6.2.1 Training loop

To add a training loop use the *training_step* method

```

class LitClassifier(pl.LightningModule):

    def __init__(self, model):
        super().__init__()
        self.model = model

    def training_step(self, batch, batch_idx):
        x, y = batch
        y_hat = self.model(x)
        loss = F.cross_entropy(y_hat, y)
        return loss

```

Under the hood, Lightning does the following (pseudocode):

```

# put model in train mode
model.train()
torch.set_grad_enabled(True)

```

(continues on next page)

(continued from previous page)

```
outs = []
for batch in train_dataloader:
    # forward
    out = training_step(val_batch)

    # backward
    loss.backward()

    # apply and clear grads
    optimizer.step()
    optimizer.zero_grad()
```

Training epoch-level metrics

If you want to calculate epoch-level metrics and log them, use the `.log` method

```
def training_step(self, batch, batch_idx):
    x, y = batch
    y_hat = self.model(x)
    loss = F.cross_entropy(y_hat, y)

    # logs metrics for each training_step,
    # and the average across the epoch, to the progress bar and logger
    self.log('train_loss', loss, on_step=True, on_epoch=True, prog_bar=True,
    ↪ logger=True)
    return loss
```

The `.log` object automatically reduces the requested metrics across the full epoch. Here's the pseudocode of what it does under the hood:

```
outs = []
for batch in train_dataloader:
    # forward
    out = training_step(val_batch)

    # backward
    loss.backward()

    # apply and clear grads
    optimizer.step()
    optimizer.zero_grad()

epoch_metric = torch.mean(torch.stack([x['train_loss'] for x in outs]))
```

Train epoch-level operations

If you need to do something with all the outputs of each *training_step*, override *training_epoch_end* yourself.

```
def training_step(self, batch, batch_idx):
    x, y = batch
    y_hat = self.model(x)
    loss = F.cross_entropy(y_hat, y)
    preds = ...
    return {'loss': loss, 'other_stuff': preds}

def training_epoch_end(self, training_step_outputs):
    for pred in training_step_outputs:
        # do something
```

The matching pseudocode is:

```
outs = []
for batch in train_dataloader:
    # forward
    out = training_step(val_batch)

    # backward
    loss.backward()

    # apply and clear grads
    optimizer.step()
    optimizer.zero_grad()

training_epoch_end(outs)
```

Training with DataParallel

When training using a *distributed_backend* that splits data from each batch across GPUs, sometimes you might need to aggregate them on the master GPU for processing (dp, or ddp2).

In this case, implement the *training_step_end* method

```
def training_step(self, batch, batch_idx):
    x, y = batch
    y_hat = self.model(x)
    loss = F.cross_entropy(y_hat, y)
    pred = ...
    return {'loss': loss, 'pred': pred}

def training_step_end(self, batch_parts):
    gpu_0_prediction = batch_parts.pred[0]['pred']
    gpu_1_prediction = batch_parts.pred[1]['pred']

    # do something with both outputs
    return (batch_parts[0]['loss'] + batch_parts[1]['loss']) / 2

def training_epoch_end(self, training_step_outputs):
    for out in training_step_outputs:
        # do something with preds
```

The full pseudocode that lightning does under the hood is:

```
outs = []
for train_batch in train_dataloader:
    batches = split_batch(train_batch)
    dp_outs = []
    for sub_batch in batches:
        # 1
        dp_out = training_step(sub_batch)
        dp_outs.append(dp_out)

    # 2
    out = training_step_end(dp_outs)
    outs.append(out)

# do something with the outputs for all batches
# 3
training_epoch_end(outs)
```

6.2.2 Validation loop

To add a validation loop, override the `validation_step` method of the `LightningModule`:

```
class LitModel(pl.LightningModule):
    def validation_step(self, batch, batch_idx):
        x, y = batch
        y_hat = self.model(x)
        loss = F.cross_entropy(y_hat, y)
        self.log('val_loss', loss)
```

Under the hood, Lightning does the following:

```
# ...
for batch in train_dataloader:
    loss = model.training_step()
    loss.backward()
    # ...

    if validate_at_some_point:
        # disable grads + batchnorm + dropout
        torch.set_grad_enabled(False)
        model.eval()

        # ----- VAL LOOP -----
        for val_batch in model.val_dataloader:
            val_out = model.validation_step(val_batch)
        # ----- VAL LOOP -----

        # enable grads + batchnorm + dropout
        torch.set_grad_enabled(True)
        model.train()
```


Validation epoch-level metrics

If you need to do something with all the outputs of each *validation_step*, override *validation_epoch_end*.

```
def validation_step(self, batch, batch_idx):
    x, y = batch
    y_hat = self.model(x)
    loss = F.cross_entropy(y_hat, y)
    pred = ...
    return pred

def validation_epoch_end(self, validation_step_outputs):
    for pred in validation_step_outputs:
        # do something with a pred
```

Validating with DataParallel

When training using a *distributed_backend* that splits data from each batch across GPUs, sometimes you might need to aggregate them on the master GPU for processing (dp, or ddp2).

In this case, implement the *validation_step_end* method

```
def validation_step(self, batch, batch_idx):
    x, y = batch
    y_hat = self.model(x)
    loss = F.cross_entropy(y_hat, y)
    pred = ...
    return {'loss': loss, 'pred': pred}

def validation_step_end(self, batch_parts):
    gpu_0_prediction = batch_parts.pred[0]['pred']
    gpu_1_prediction = batch_parts.pred[1]['pred']

    # do something with both outputs
    return (batch_parts[0]['loss'] + batch_parts[1]['loss']) / 2

def validation_epoch_end(self, validation_step_outputs):
    for out in validation_step_outputs:
        # do something with preds
```

The full pseudocode that lightning does under the hood is:

```
outs = []
for batch in dataloader:
    batches = split_batch(batch)
    dp_outs = []
    for sub_batch in batches:
        # 1
        dp_out = validation_step(sub_batch)
        dp_outs.append(dp_out)

    # 2
    out = validation_step_end(dp_outs)
    outs.append(out)

# do something with the outputs for all batches
```

(continues on next page)

(continued from previous page)

```
# 3
validation_epoch_end(outs)
```

6.2.3 Test loop

The process for adding a test loop is the same as the process for adding a validation loop. Please refer to the section above for details.

The only difference is that the test loop is only called when `.test()` is used:

```
model = Model()
trainer = Trainer()
trainer.fit()

# automatically loads the best weights for you
trainer.test(model)
```

There are two ways to call `test()`:

```
# call after training
trainer = Trainer()
trainer.fit(model)

# automatically auto-loads the best weights
trainer.test(test_dataloaders=test_dataloader)

# or call with pretrained model
model = MyLightningModule.load_from_checkpoint(PATH)
trainer = Trainer()
trainer.test(model, test_dataloaders=test_dataloader)
```

6.3 Inference

For research, LightningModules are best structured as systems.

```
import pytorch_lightning as pl
import torch
from torch import nn

class Autoencoder(pl.LightningModule):

    def __init__(self, latent_dim=2):
        super().__init__()
        self.encoder = nn.Sequential(nn.Linear(28 * 28, 256), nn.ReLU(), nn.
↳Linear(256, latent_dim))
        self.decoder = nn.Sequential(nn.Linear(latent_dim, 256), nn.ReLU(), nn.
↳Linear(256, 28 * 28))

    def training_step(self, batch, batch_idx):
```

(continues on next page)

(continued from previous page)

```

x, _ = batch

# encode
x = x.view(x.size(0), -1)
z = self.encoder(x)

# decode
recons = self.decoder(z)

# reconstruction
reconstruction_loss = nn.functional.mse_loss(recons, x)
return reconstruction_loss

def validation_step(self, batch, batch_idx):
    x, _ = batch
    x = x.view(x.size(0), -1)
    z = self.encoder(x)
    recons = self.decoder(z)
    reconstruction_loss = nn.functional.mse_loss(recons, x)
    self.log('val_reconstruction', reconstruction_loss)

def configure_optimizers(self):
    return torch.optim.Adam(self.parameters(), lr=0.0002)

```

Which can be trained like this:

```

autoencoder = Autoencoder()
trainer = pl.Trainer(gpus=1)
trainer.fit(autoencoder, train_dataloader, val_dataloader)

```

This simple model generates examples that look like this (the encoders and decoders are too weak)



The methods above are part of the lightning interface:

- training_step
- validation_step
- test_step
- configure_optimizers

Note that in this case, the train loop and val loop are exactly the same. We can of course reuse this code.

```

class Autoencoder(pl.LightningModule):

    def __init__(self, latent_dim=2):
        super().__init__()
        self.encoder = nn.Sequential(nn.Linear(28 * 28, 256), nn.ReLU(), nn.
↳ Linear(256, latent_dim))

```

(continues on next page)

(continued from previous page)

```

        self.decoder = nn.Sequential(nn.Linear(latent_dim, 256), nn.ReLU(), nn.
↪Linear(256, 28 * 28))

    def training_step(self, batch, batch_idx):
        loss = self.shared_step(batch)

        return loss

    def validation_step(self, batch, batch_idx):
        loss = self.shared_step(batch)
        self.log('val_loss', loss)

    def shared_step(self, batch):
        x, _ = batch

        # encode
        x = x.view(x.size(0), -1)
        z = self.encoder(x)

        # decode
        recons = self.decoder(z)

        # loss
        return nn.functional.mse_loss(recons, x)

    def configure_optimizers(self):
        return torch.optim.Adam(self.parameters(), lr=0.0002)

```

We create a new method called *shared_step* that all loops can use. This method name is arbitrary and NOT reserved.

6.3.1 Inference in research

In the case where we want to perform inference with the system we can add a *forward* method to the LightningModule.

```

class Autoencoder(pl.LightningModule):
    def forward(self, x):
        return self.decoder(x)

```

The advantage of adding a forward is that in complex systems, you can do a much more involved inference procedure, such as text generation:

```

class Seq2Seq(pl.LightningModule):

    def forward(self, x):
        embeddings = self(x)
        hidden_states = self.encoder(embeddings)
        for h in hidden_states:
            # decode
            ...
        return decoded

```

6.3.2 Inference in production

For cases like production, you might want to iterate different models inside a LightningModule.

```
import pytorch_lightning as pl
from pytorch_lightning.metrics import functional as FM

class ClassificationTask(pl.LightningModule):

    def __init__(self, model):
        super().__init__()
        self.model = model

    def training_step(self, batch, batch_idx):
        x, y = batch
        y_hat = self.model(x)
        loss = F.cross_entropy(y_hat, y)
        return loss

    def validation_step(self, batch, batch_idx):
        x, y = batch
        y_hat = self.model(x)
        loss = F.cross_entropy(y_hat, y)
        acc = FM.accuracy(y_hat, y)

        # loss is tensor. The Checkpoint Callback is monitoring 'checkpoint_on'
        metrics = {'val_acc': acc, 'val_loss': loss}
        self.log_dict(metrics)
        return metrics

    def test_step(self, batch, batch_idx):
        metrics = self.validation_step(batch, batch_idx)
        metrics = {'test_acc': metrics['val_acc'], 'test_loss': metrics['val_loss']}
        self.log_dict(metrics)

    def configure_optimizers(self):
        return torch.optim.Adam(self.model.parameters(), lr=0.02)
```

Then pass in any arbitrary model to be fit with this task

```
for model in [resnet50(), vgg16(), BidirectionalRNN()]:
    task = ClassificationTask(model)

    trainer = Trainer(gpus=2)
    trainer.fit(task, train_dataloader, val_dataloader)
```

Tasks can be arbitrarily complex such as implementing GAN training, self-supervised or even RL.

```
class GANTask(pl.LightningModule):

    def __init__(self, generator, discriminator):
        super().__init__()
        self.generator = generator
        self.discriminator = discriminator
        ...
```

When used like this, the model can be separated from the Task and thus used in production without needing to keep it in a *LightningModule*.

- You can export to onnx.
- Or trace using Jit.
- or run in the python runtime.

```
task = ClassificationTask(model)

trainer = Trainer(gpus=2)
trainer.fit(task, train_dataloader, val_dataloader)

# use model after training or load weights and drop into the production system
model.eval()
y_hat = model(x)
```

6.4 LightningModule API

6.4.1 Methods

configure_optimizers

`LightningModule.configure_optimizers()`

Choose what optimizers and learning-rate schedulers to use in your optimization. Normally you'd need one. But in the case of GANs or similar you might have multiple.

Returns

Any of these 6 options.

- Single optimizer.
- List or Tuple - List of optimizers.
- Two lists - The first list has multiple optimizers, the second a list of LR schedulers (or lr_dict).
- Dictionary, with an 'optimizer' key, and (optionally) a 'lr_scheduler' key which value is a single LR scheduler or lr_dict.
- Tuple of dictionaries as described, with an optional 'frequency' key.
- None - Fit will run without any optimizer.

Note: The 'frequency' value is an int corresponding to the number of sequential batches optimized with the specific optimizer. It should be given to none or to all of the optimizers. There is a difference between passing multiple optimizers in a list, and passing multiple optimizers in dictionaries with a frequency of 1: In the former case, all optimizers will operate on the given batch in each optimization step. In the latter, only one optimizer will operate on the given batch at every step.

The lr_dict is a dictionary which contains scheduler and its associated configuration. It has five keys. The default configuration is shown below.

```
{
    'scheduler': lr_scheduler, # The LR scheduler
    'interval': 'epoch', # The unit of the scheduler's step size
```

(continues on next page)

(continued from previous page)

```

'frequency': 1, # The frequency of the scheduler
'reduce_on_plateau': False, # For ReduceLROnPlateau scheduler
'monitor': 'val_loss' # Metric to monitor
}

```

If user only provides LR schedulers, then their configuration will set to default as shown above.

Examples

```

# most cases
def configure_optimizers(self):
    opt = Adam(self.parameters(), lr=1e-3)
    return opt

# multiple optimizer case (e.g.: GAN)
def configure_optimizers(self):
    generator_opt = Adam(self.model_gen.parameters(), lr=0.01)
    discriminator_opt = Adam(self.model_disc.parameters(), lr=0.02)
    return generator_opt, discriminator_opt

# example with learning rate schedulers
def configure_optimizers(self):
    generator_opt = Adam(self.model_gen.parameters(), lr=0.01)
    discriminator_opt = Adam(self.model_disc.parameters(), lr=0.02)
    discriminator_sched = CosineAnnealing(discriminator_opt, T_max=10)
    return [generator_opt, discriminator_opt], [discriminator_sched]

# example with step-based learning rate schedulers
def configure_optimizers(self):
    gen_opt = Adam(self.model_gen.parameters(), lr=0.01)
    dis_opt = Adam(self.model_disc.parameters(), lr=0.02)
    gen_sched = {'scheduler': ExponentialLR(gen_opt, 0.99),
                 'interval': 'step'} # called after each training step
    dis_sched = CosineAnnealing(discriminator_opt, T_max=10) # called every epoch
    return [gen_opt, dis_opt], [gen_sched, dis_sched]

# example with optimizer frequencies
# see training procedure in `Improved Training of Wasserstein GANs`, Algorithm 1
# https://arxiv.org/abs/1704.00028
def configure_optimizers(self):
    gen_opt = Adam(self.model_gen.parameters(), lr=0.01)
    dis_opt = Adam(self.model_disc.parameters(), lr=0.02)
    n_critic = 5
    return (
        {'optimizer': dis_opt, 'frequency': n_critic},
        {'optimizer': gen_opt, 'frequency': 1}
    )

```

Note: Some things to know:

- Lightning calls `.backward()` and `.step()` on each optimizer and learning rate scheduler as needed.
- If you use 16-bit precision (`precision=16`), Lightning will automatically handle the optimizers for you.

- If you use multiple optimizers, `training_step()` will have an additional `optimizer_idx` parameter.
- If you use LBFGS Lightning handles the closure function automatically for you.
- If you use multiple optimizers, gradients will be calculated only for the parameters of current optimizer at each training step.
- If you need to control how often those optimizers step or override the default `.step()` schedule, override the `optimizer_step()` hook.
- If you only want to call a learning rate scheduler every `x` step or epoch, or want to monitor a custom metric, you can specify these in a `lr_dict`:

```
{
    'scheduler': lr_scheduler,
    'interval': 'step', # or 'epoch'
    'monitor': 'val_f1',
    'frequency': x,
}
```

forward

`LightningModule.forward(*args, **kwargs)`

Same as `torch.nn.Module.forward()`, however in Lightning you want this to define the operations you want to use for prediction (i.e.: on a server or as a feature extractor).

Normally you'd call `self()` from your `training_step()` method. This makes it easy to write a complex system for training with the outputs you'd want in a prediction setting.

You may also find the `auto_move_data()` decorator useful when using the module outside Lightning in a production setting.

Parameters

- **`*args`** – Whatever you decide to pass into the forward method.
- **`**kwargs`** – Keyword arguments are also possible.

Returns Predicted output

Examples

```
# example if we were using this model as a feature extractor
def forward(self, x):
    feature_maps = self.convnet(x)
    return feature_maps

def training_step(self, batch, batch_idx):
    x, y = batch
    feature_maps = self(x)
    logits = self.classifier(feature_maps)

    # ...
    return loss

# splitting it this way allows model to be used a feature extractor
```

(continues on next page)

(continued from previous page)

```

model = MyModelAbove()

inputs = server.get_request()
results = model(inputs)
server.write_results(results)

# -----
# This is in stark contrast to torch.nn.Module where normally you would have this:
def forward(self, batch):
    x, y = batch
    feature_maps = self.convnet(x)
    logits = self.classifier(feature_maps)
    return logits

```

freeze

`LightningModule.freeze()`
Freeze all params for inference.

Example

```

model = MyLightningModule(...)
model.freeze()

```

Return type `None`

log

`LightningModule.log(name, value, prog_bar=False, logger=True, on_step=None, on_epoch=None, reduce_fx=torch.mean, tbptt_reduce_fx=torch.mean, tbptt_pad_token=0, enable_graph=False, sync_dist=False, sync_dist_op='mean', sync_dist_group=None)`

Log a key, value

Example:

```
self.log('train_loss', loss)
```

The default behavior per hook is as follows

Table 1: * also applies to the test loop

LightningModule Hook	on_step	on_epoch	prog_bar	logger
training_step	T	F	F	T
training_step_end	T	F	F	T
training_epoch_end	F	T	F	T
validation_step*	F	T	F	T
validation_step_end*	F	T	F	T
validation_epoch_end*	F	T	F	T

Parameters

- **name** (str) – key name
- **value** (Any) – value name
- **prog_bar** (bool) – if True logs to the progress bar
- **logger** (bool) – if True logs to the logger
- **on_step** (Optional[bool]) – if True logs at this step. None auto-logs at the training_step but not validation/test_step
- **on_epoch** (Optional[bool]) – if True logs epoch accumulated metrics. None auto-logs at the val/test step but not training_step
- **reduce_fx** (Callable) – Torch.mean by default
- **tbptt_reduce_fx** (Callable) – function to reduce on truncated back prop
- **tbptt_pad_token** (int) – token to use for padding
- **enable_graph** (bool) – if True, will not auto detach the graph
- **sync_dist** (bool) – if True, reduces the metric across GPUs/TPUs
- **sync_dist_op** (Union[Any, str]) – the op to sync across
- **sync_dist_group** (Optional[Any]) – the ddp group

log_dict

LightningModule.**log_dict** (dictionary, prog_bar=False, logger=True, on_step=None, on_epoch=None, reduce_fx=torch.mean, tbptt_reduce_fx=torch.mean, tbptt_pad_token=0, enable_graph=False, sync_dist=False, sync_dist_op='mean', sync_dist_group=None)

Log a dictionary of values at once

Example:

```
values = {'loss': loss, 'acc': acc, ..., 'metric_n': metric_n}
self.log_dict(values)
```

Parameters

- **dictionary** (dict) – key value pairs (str, tensors)
- **prog_bar** (bool) – if True logs to the progress base
- **logger** (bool) – if True logs to the logger
- **on_step** (Optional[bool]) – if True logs at this step. None auto-logs for training_step but not validation/test_step
- **on_epoch** (Optional[bool]) – if True logs epoch accumulated metrics. None auto-logs for val/test step but not training_step
- **reduce_fx** (Callable) – Torch.mean by default
- **tbptt_reduce_fx** (Callable) – function to reduce on truncated back prop
- **tbptt_pad_token** (int) – token to use for padding
- **enable_graph** (bool) – if True, will not auto detach the graph
- **sync_dist** (bool) – if True, reduces the metric across GPUs/TPUs

- `sync_dist_op` (Union[Any, str]) – the op to sync across
- `sync_dist_group` (Optional[Any]) – the ddp group:

print

LightningModule.**print** (*args, **kwargs)

Prints only from process 0. Use this in any distributed mode to log only once.

Parameters

- `*args` – The thing to print. Will be passed to Python’s built-in print function.
- `**kwargs` – Will be passed to Python’s built-in print function.

Example

```
def forward(self, x):
    self.print(x, 'in forward')
```

Return type None

save_hyperparameters

LightningModule.**save_hyperparameters** (*args, frame=None)

Save all model arguments.

Parameters `args` – single object of *dict*, *Namespace* or *OmegaConf* or string names or argumentst from class `__init__`

```
>>> from collections import OrderedDict
>>> class ManuallyArgsModel(LightningModule):
...     def __init__(self, arg1, arg2, arg3):
...         super().__init__()
...         # manually assign arguments
...         self.save_hyperparameters('arg1', 'arg3')
...     def forward(self, *args, **kwargs):
...         ...
>>> model = ManuallyArgsModel(1, 'abc', 3.14)
>>> model.hparams
"arg1": 1
"arg3": 3.14
```

```
>>> class AutomaticArgsModel(LightningModule):
...     def __init__(self, arg1, arg2, arg3):
...         super().__init__()
...         # equivalent automatic
...         self.save_hyperparameters()
...     def forward(self, *args, **kwargs):
...         ...
>>> model = AutomaticArgsModel(1, 'abc', 3.14)
>>> model.hparams
"arg1": 1
"arg2": abc
"arg3": 3.14
```

```
>>> class SingleArgModel(LightningModule):
...     def __init__(self, params):
...         super().__init__()
...         # manually assign single argument
...         self.save_hyperparameters(params)
...     def forward(self, *args, **kwargs):
...         ...
>>> model = SingleArgModel(Namespace(p1=1, p2='abc', p3=3.14))
>>> model.hparams
"p1": 1
"p2": abc
"p3": 3.14
```

Return type `None`

test_step

`LightningModule.test_step(*args, **kwargs)`

Operates on a single batch of data from the test set. In this step you'd normally generate examples or calculate anything of interest such as accuracy.

```
# the pseudocode for these calls
test_outs = []
for test_batch in test_data:
    out = test_step(test_batch)
    test_outs.append(out)
test_epoch_end(test_outs)
```

Parameters

- **batch** `(Tensor | (Tensor, ...) | [Tensor, ...])` – The output of your `DataLoader`. A tensor, tuple or list.
- **batch_idx** `(int)` – The index of this batch.
- **dataloader_idx** `(int)` – The index of the dataloader that produced this batch (only if multiple test datasets used).

Returns

Any of.

- Any object or value
- *None* - Testing will skip to the next batch

```
# if you have one test dataloader:
def test_step(self, batch, batch_idx)

# if you have multiple test dataloaders:
def test_step(self, batch, batch_idx, dataloader_idx)
```

Examples

```
# CASE 1: A single test dataset
def test_step(self, batch, batch_idx):
    x, y = batch

    # implement your own
    out = self(x)
    loss = self.loss(out, y)

    # log 6 example images
    # or generated text... or whatever
    sample_imgs = x[:6]
    grid = torchvision.utils.make_grid(sample_imgs)
    self.logger.experiment.add_image('example_images', grid, 0)

    # calculate acc
    labels_hat = torch.argmax(out, dim=1)
    test_acc = torch.sum(y == labels_hat).item() / (len(y) * 1.0)

    # log the outputs!
    self.log_dict({'test_loss': loss, 'test_acc': test_acc})
```

If you pass in multiple validation datasets, `test_step()` will have an additional argument.

```
# CASE 2: multiple test datasets
def test_step(self, batch, batch_idx, dataloader_idx):
    # dataloader_idx tells you which dataset this is.
```

Note: If you don't need to validate you don't need to implement this method.

Note: When the `test_step()` is called, the model has been put in eval mode and PyTorch gradients have been disabled. At the end of the test epoch, the model goes back to training mode and gradients are enabled.

test_step_end

`LightningModule.test_step_end(*args, **kwargs)`

Use this when testing with dp or ddp2 because `test_step()` will operate on only part of the batch. However, this is still optional and only needed for things like softmax or NCE loss.

Note: If you later switch to ddp or some other mode, this will still be called so that you don't have to change your code.

```
# pseudocode
sub_batches = split_batches_for_dp(batch)
batch_parts_outputs = [test_step(sub_batch) for sub_batch in sub_batches]
test_step_end(batch_parts_outputs)
```

Parameters `batch_parts_outputs` – What you return in `test_step()` for each batch part.

Returns None or anything

```
# WITHOUT test_step_end
# if used in DP or DDP2, this batch is 1/num_gpus large
def test_step(self, batch, batch_idx):
    # batch is 1/num_gpus big
    x, y = batch

    out = self(x)
    loss = self.softmax(out)
    self.log('test_loss', loss)

# -----
# with test_step_end to do softmax over the full batch
def test_step(self, batch, batch_idx):
    # batch is 1/num_gpus big
    x, y = batch

    out = self.encoder(x)
    return out

def test_epoch_end(self, output_results):
    # this out is now the full size of the batch
    all_test_step_outs = output_results.out
    loss = nce_loss(all_test_step_outs)
    self.log('test_loss', loss)
```

See also:

See the [Multi-GPU training](#) guide for more details.

test_epoch_end

`LightningModule.test_epoch_end(outputs)`

Called at the end of a test epoch with the output of all test steps.

```
# the pseudocode for these calls
test_outs = []
for test_batch in test_data:
    out = test_step(test_batch)
    test_outs.append(out)
test_epoch_end(test_outs)
```

Parameters `outputs` `List[Any]` – List of outputs you defined in `test_step_end()`, or if there are multiple dataloaders, a list containing a list of outputs for each dataloader

Return type `None`

Returns `None`

Note: If you didn't define a `test_step()`, this won't be called.

Examples

With a single dataloader:

```
def test_epoch_end(self, outputs):
    # do something with the outputs of all test batches
    all_test_preds = test_step_outputs.predictions

    some_result = calc_all_results(all_test_preds)
    self.log(some_result)
```

With multiple dataloaders, *outputs* will be a list of lists. The outer list contains one entry per dataloader, while the inner list contains the individual outputs of each test step for that dataloader.

```
def test_epoch_end(self, outputs):
    final_value = 0
    for dataloader_outputs in outputs:
        for test_step_out in dataloader_outputs:
            # do something
            final_value += test_step_out

    self.log('final_metric', final_value)
```

to_onnx

`LightningModule.to_onnx` (*file_path*, *input_sample=None*, ***kwargs*)

Saves the model in ONNX format

Parameters

- **file_path** (str) – The path of the file the model should be saved to.
- **input_sample** (Optional[`Tensor`]) – A sample of an input tensor for tracing.
- ****kwargs** – Will be passed to `torch.onnx.export` function.

Example

```
>>> class SimpleModel(LightningModule):
...     def __init__(self):
...         super().__init__()
...         self.l1 = torch.nn.Linear(in_features=64, out_features=4)
...
...     def forward(self, x):
...         return torch.relu(self.l1(x.view(x.size(0), -1)))

>>> with tempfile.NamedTemporaryFile(suffix='.onnx', delete=False) as tmpfile:
...     model = SimpleModel()
...     input_sample = torch.randn((1, 64))
...     model.to_onnx(tmpfile.name, input_sample, export_params=True)
...     os.path.isfile(tmpfile.name)
True
```

to_torchscript

`LightningModule.to_torchscript` (*file_path=None*, *method='script'*, *example_inputs=None*, ***kwargs*)

By default compiles the whole model to a `ScriptModule`. If you want to use tracing, please provided the argument `method='trace'` and make sure that either the `example_inputs` argument is provided, or the model has `self.example_input_array` set. If you would like to customize the modules that are scripted you should override this method. In case you want to return multiple modules, we recommend using a dictionary.

Parameters

- **file_path** (Optional[str]) – Path where to save the torchscript. Default: None (no file saved).
- **method** (Optional[str]) – Whether to use TorchScript's script or trace method. Default: 'script'
- **example_inputs** (Optional[Tensor]) – Tensor to be used to do tracing when method is set to 'trace'. Default: None (Use `self.example_input_array`)
- ****kwargs** – Additional arguments that will be passed to the `torch.jit.script()` or `torch.jit.trace()` function.

Note:

- Requires the implementation of the `forward()` method.
 - The exported script will be set to evaluation mode.
 - It is recommended that you install the latest supported version of PyTorch to use this feature without limitations. See also the `torch.jit` documentation for supported features.
-

Example

```
>>> class SimpleModel(LightningModule):
...     def __init__(self):
...         super().__init__()
...         self.l1 = torch.nn.Linear(in_features=64, out_features=4)
...
...     def forward(self, x):
...         return torch.relu(self.l1(x.view(x.size(0), -1)))
...
>>> model = SimpleModel()
>>> torch.jit.save(model.to_torchscript(), "model.pt")
>>> os.path.isfile("model.pt")
True
```

Return type `Union[ScriptModule, Dict[str, ScriptModule]]`

Returns This `LightningModule` as a torchscript, regardless of whether `file_path` is defined or not.

training_step

`LightningModule.training_step(*args, **kwargs)`

Here you compute and return the training loss and some additional metrics for e.g. the progress bar or logger.

Parameters

- `batch` (Tensor | (Tensor, ...) | [Tensor, ...]) – The output of your `DataLoader`. A tensor, tuple or list.
- `batch_idx` (int) – Integer displaying index of this batch
- `optimizer_idx` (int) – When using multiple optimizers, this argument will also be present.
- `hiddens` (Tensor) – Passed in if `truncated_bptt_steps > 0`.

Returns

Any of.

- `Tensor` - The loss tensor
- `dict` - A dictionary. Can include any keys, but must include the key 'loss'
- `None` - Training will skip to the next batch

In this step you'd normally do the forward pass and calculate the loss for a batch. You can also do fancier things like multiple forward passes or something model specific.

Example:

```
def training_step(self, batch, batch_idx):
    x, y, z = batch
    out = self.encoder(x)
    loss = self.loss(out, x)
    return loss
```

If you define multiple optimizers, this step will be called with an additional `optimizer_idx` parameter.

```
# Multiple optimizers (e.g.: GANs)
def training_step(self, batch, batch_idx, optimizer_idx):
    if optimizer_idx == 0:
        # do training_step with encoder
    if optimizer_idx == 1:
        # do training_step with decoder
```

If you add truncated back propagation through time you will also get an additional argument with the hidden states of the previous step.

```
# Truncated back-propagation through time
def training_step(self, batch, batch_idx, hiddens):
    # hiddens are the hidden states from the previous truncated backprop step
    ...
    out, hiddens = self.lstm(data, hiddens)
    ...
    return {'loss': loss, 'hiddens': hiddens}
```

Note: The loss value shown in the progress bar is smoothed (averaged) over the last values, so it differs from the actual loss returned in train/validation step.

training_step_end

LightningModule.**training_step_end**(*args, **kwargs)

Use this when training with dp or ddp2 because `training_step()` will operate on only part of the batch. However, this is still optional and only needed for things like softmax or NCE loss.

Note: If you later switch to ddp or some other mode, this will still be called so that you don't have to change your code

```
# pseudocode
sub_batches = split_batches_for_dp(batch)
batch_parts_outputs = [training_step(sub_batch) for sub_batch in sub_batches]
training_step_end(batch_parts_outputs)
```

Parameters `batch_parts_outputs` – What you return in *training_step* for each batch part.

Returns Anything

When using dp/ddp2 distributed backends, only a portion of the batch is inside the *training_step*:

```
def training_step(self, batch, batch_idx):
    # batch is 1/num_gpus big
    x, y = batch

    out = self(x)

    # softmax uses only a portion of the batch in the denominator
    loss = self.softmax(out)
    loss = nce_loss(loss)
    return loss
```

If you wish to do something with all the parts of the batch, then use this method to do it:

```
def training_step(self, batch, batch_idx):
    # batch is 1/num_gpus big
    x, y = batch

    out = self.encoder(x)
    return {'pred': out}

def training_step_end(self, training_step_outputs):
    gpu_0_pred = training_step_outputs[0]['pred']
    gpu_1_pred = training_step_outputs[1]['pred']
    gpu_n_pred = training_step_outputs[n]['pred']

    # this softmax now uses the full batch
    loss = nce_loss([gpu_0_pred, gpu_1_pred, gpu_n_pred])
    return loss
```

See also:

See the *Multi-GPU training* guide for more details.

training_epoch_end

LightningModule.**training_epoch_end**(*outputs*)

Called at the end of the training epoch with the outputs of all training steps. Use this in case you need to do something with all the outputs for every training_step.

```
# the pseudocode for these calls
train_outs = []
for train_batch in train_data:
    out = training_step(train_batch)
    train_outs.append(out)
training_epoch_end(train_outs)
```

Parameters **outputs** (`List[Any]`) – List of outputs you defined in `training_step()`, or if there are multiple dataloaders, a list containing a list of outputs for each dataloader.

Return type `None`

Returns `None`

Note: If this method is not overridden, this won't be called.

Example:

```
def training_epoch_end(self, training_step_outputs):
    # do something with all training_step outputs
    return result
```

With multiple dataloaders, `outputs` will be a list of lists. The outer list contains one entry per dataloader, while the inner list contains the individual outputs of each training step for that dataloader.

```
def training_epoch_end(self, training_step_outputs):
    for out in training_step_outputs:
        # do something here
```

unfreeze

LightningModule.**unfreeze**()

Unfreeze all parameters for training.

```
model = MyLightningModule(...)
model.unfreeze()
```

Return type `None`

validation_step

LightningModule.**validation_step**(*args, **kwargs)

Operates on a single batch of data from the validation set. In this step you'd might generate examples or calculate anything of interest like accuracy.

```
# the pseudocode for these calls
val_outs = []
for val_batch in val_data:
    out = validation_step(train_batch)
    val_outs.append(out)
validation_epoch_end(val_outs)
```

Parameters

- **batch** `(Tensor | (Tensor, ...) | [Tensor, ...])` – The output of your `DataLoader`. A tensor, tuple or list.
- **batch_idx** `(int)` – The index of this batch
- **dataloader_idx** `(int)` – The index of the dataloader that produced this batch (only if multiple val datasets used)

Returns

Any of.

- Any object or value
- *None* - Validation will skip to the next batch

```
# pseudocode of order
out = validation_step()
if defined('validation_step_end'):
    out = validation_step_end(out)
out = validation_epoch_end(out)
```

```
# if you have one val dataloader:
def validation_step(self, batch, batch_idx)

# if you have multiple val dataloaders:
def validation_step(self, batch, batch_idx, dataloader_idx)
```

Examples

```
# CASE 1: A single validation dataset
def validation_step(self, batch, batch_idx):
    x, y = batch

    # implement your own
    out = self(x)
    loss = self.loss(out, y)

    # log 6 example images
    # or generated text... or whatever
    sample_imgs = x[:6]
```

(continues on next page)

(continued from previous page)

```

grid = torchvision.utils.make_grid(sample_imgs)
self.logger.experiment.add_image('example_images', grid, 0)

# calculate acc
labels_hat = torch.argmax(out, dim=1)
val_acc = torch.sum(y == labels_hat).item() / (len(y) * 1.0)

# log the outputs!
self.log_dict({'val_loss': loss, 'val_acc': val_acc})

```

If you pass in multiple val datasets, `validation_step` will have an additional argument.

```

# CASE 2: multiple validation datasets
def validation_step(self, batch, batch_idx, dataloader_idx):
    # dataloader_idx tells you which dataset this is.

```

Note: If you don't need to validate you don't need to implement this method.

Note: When the `validation_step()` is called, the model has been put in eval mode and PyTorch gradients have been disabled. At the end of validation, the model goes back to training mode and gradients are enabled.

validation_step_end

`LightningModule.validation_step_end(*args, **kwargs)`

Use this when validating with dp or ddp2 because `validation_step()` will operate on only part of the batch. However, this is still optional and only needed for things like softmax or NCE loss.

Note: If you later switch to ddp or some other mode, this will still be called so that you don't have to change your code.

```

# pseudocode
sub_batches = split_batches_for_dp(batch)
batch_parts_outputs = [validation_step(sub_batch) for sub_batch in sub_batches]
validation_step_end(batch_parts_outputs)

```

Parameters `batch_parts_outputs` – What you return in `validation_step()` for each batch part.

Returns None or anything

```

# WITHOUT validation_step_end
# if used in DP or DDP2, this batch is 1/num_gpus large
def validation_step(self, batch, batch_idx):
    # batch is 1/num_gpus big
    x, y = batch

    out = self.encoder(x)
    loss = self.softmax(out)

```

(continues on next page)

(continued from previous page)

```

    loss = nce_loss(loss)
    self.log('val_loss', loss)

# -----
# with validation_step_end to do softmax over the full batch
def validation_step(self, batch, batch_idx):
    # batch is 1/num_gpus big
    x, y = batch

    out = self(x)
    return out

def validation_epoch_end(self, val_step_outputs):
    for out in val_step_outputs:
        # do something with these

```

See also:

See the *Multi-GPU training* guide for more details.

validation_epoch_end

LightningModule.**validation_epoch_end**(*outputs*)

Called at the end of the validation epoch with the outputs of all validation steps.

```

# the pseudocode for these calls
val_outs = []
for val_batch in val_data:
    out = validation_step(val_batch)
    val_outs.append(out)
validation_epoch_end(val_outs)

```

Parameters **outputs** (List[Any]) – List of outputs you defined in `validation_step()`, or if there are multiple dataloaders, a list containing a list of outputs for each dataloader.

Return type None

Returns None

Note: If you didn't define a `validation_step()`, this won't be called.

Examples

With a single dataloader:

```

def validation_epoch_end(self, val_step_outputs):
    for out in val_step_outputs:
        # do something

```

With multiple dataloaders, *outputs* will be a list of lists. The outer list contains one entry per dataloader, while the inner list contains the individual outputs of each validation step for that dataloader.

```
def validation_epoch_end(self, outputs):
    for dataloader_output_result in outputs:
        dataloader_outs = dataloader_output_result.dataloader_i_outputs

    self.log('final_metric', final_value)
```

6.4.2 Properties

These are properties available in a LightningModule.

current_epoch

The current epoch

```
def training_step(...):
    if self.current_epoch == 0:
```

device

The device the module is on. Use it to keep your code device agnostic

```
def training_step(...):
    z = torch.rand(2, 3, device=self.device)
```

global_rank

The global_rank of this LightningModule. Lightning saves logs, weights etc only from global_rank = 0. You normally do not need to use this property

Global rank refers to the index of that GPU across ALL GPUs. For example, if using 10 machines, each with 4 GPUs, the 4th GPU on the 10th machine has global_rank = 39

global_step

The current step (does not reset each epoch)

```
def training_step(...):
    self.logger.experiment.log_image(..., step=self.global_step)
```

hparams

After calling `save_hyperparameters` anything passed to `init()` is available via `hparams`.

```
def __init__(self, learning_rate):
    self.save_hyperparameters()

def configure_optimizers(self):
    return Adam(self.parameters(), lr=self.hparams.learning_rate)
```

logger

The current logger being used (tensorboard or other supported logger)

```
def training_step(...):
    # the generic logger (same no matter if tensorboard or other supported logger)
    self.logger

    # the particular logger
    tensorboard_logger = self.logger.experiment
```

local_rank

The `local_rank` of this `LightningModule`. Lightning saves logs, weights etc only from `global_rank = 0`. You normally do not need to use this property

Local rank refers to the rank on that machine. For example, if using 10 machines, the GPU at index 0 on each machine has `local_rank = 0`.

precision

The type of precision used:

```
def training_step(...):
    if self.precision == 16:
```

trainer

Pointer to the trainer

```
def training_step(...):
    max_steps = self.trainer.max_steps
    any_flag = self.trainer.any_flag
```

use_amp

True if using Automatic Mixed Precision (AMP)

use_ddp

True if using ddp

use_ddp2

True if using ddp2

use_dp

True if using dp

use_tpu

True if using TPUs

6.4.3 Hooks

This is the pseudocode to describe how all the hooks are called during a call to `.fit()`

```
def fit(...):
    on_fit_start()

    if global_rank == 0:
        # prepare data is called on GLOBAL_ZERO only
        prepare_data()

    for gpu/tpu in gpu/tpus:
        train_on_device(model.copy())

    on_fit_end()

def train_on_device(model):
    # setup is called PER DEVICE
    setup()
    configure_optimizers()
    on_pretrain_routine_start()
```

(continues on next page)

(continued from previous page)

```
for epoch in epochs:
    train_loop()

teardown()

def train_loop():
    on_train_epoch_start()
    train_outs = []
    for train_batch in train_dataloader():
        on_train_batch_start()

        # ----- train_step methods -----
        out = training_step(batch)
        train_outs.append(out)

        loss = out.loss

        backward()
        on_after_backward()
        optimizer_step()
        on_before_zero_grad()
        optimizer_zero_grad()

        on_train_batch_end(out)

        if should_check_val:
            val_loop()

    # end training epoch
    logs = training_epoch_end(outs)

def val_loop():
    model.eval()
    torch.set_grad_enabled(False)

    on_validation_epoch_start()
    val_outs = []
    for val_batch in val_dataloader():
        on_validation_batch_start()

        # ----- val_step methods -----
        out = validation_step(val_batch)
        val_outs.append(out)

        on_validation_batch_end(out)

    validation_epoch_end(val_outs)
    on_validation_epoch_end()

    # set up for train
    model.train()
    torch.set_grad_enabled(True)
```

backward

`LightningModule.backward(loss, optimizer, optimizer_idx)`

Override backward with your own implementation if you need to.

Parameters

- **loss** (Tensor) – Loss is already scaled by accumulated grads
- **optimizer** (Optimizer) – Current optimizer being used
- **optimizer_idx** (int) – Index of the current optimizer being used

Called to perform backward step. Feel free to override as needed. The loss passed in has already been scaled for accumulated gradients if requested.

Example:

```
def backward(self, loss, optimizer, optimizer_idx):
    loss.backward()
```

Return type `None`

get_progress_bar_dict

`LightningModule.get_progress_bar_dict()`

Implement this to override the default items displayed in the progress bar. By default it includes the average loss value, split index of BPTT (if used) and the version of the experiment when using a logger.

```
Epoch 1:   4%|          | 40/1095 [00:03<01:37, 10.84it/s, loss=4.501, v_num=10]
```

Here is an example how to override the defaults:

```
def get_progress_bar_dict(self):
    # don't show the version number
    items = super().get_progress_bar_dict()
    items.pop("v_num", None)
    return items
```

Return type `Dict[str, Union[int, str]]`

Returns Dictionary with the items to be displayed in the progress bar.

manual_backward

`LightningModule.manual_backward(loss, optimizer, *args, **kwargs)`

Call this directly from your `training_step` when doing optimizations manually. By using this we can ensure that all the proper scaling when using 16-bit etc has been done for you

This function forwards all args to the `.backward()` call as well.

Tip: In manual mode we still automatically clip grads if `Trainer(gradient_clip_val=x)` is set

Example:

```
def training_step(...):
    (opt_a, opt_b) = self.optimizers()
    loss = ...
    # automatically applies scaling, etc...
    self.manual_backward(loss, opt_a)
```

Return type `None`

`on_after_backward`

`LightningModule.on_after_backward()`

Called in the training loop after `loss.backward()` and before optimizers do anything. This is the ideal place to inspect or log gradient information.

Example:

```
def on_after_backward(self):
    # example to inspect gradient information in tensorboard
    if self.trainer.global_step % 25 == 0: # don't make the tf file huge
        params = self.state_dict()
        for k, v in params.items():
            grads = v
            name = k
            self.logger.experiment.add_histogram(tag=name, values=grads,
                                                  global_step=self.trainer.global_
↪step)
```

Return type `None`

`on_before_zero_grad`

`LightningModule.on_before_zero_grad(optimizer)`

Called after `optimizer.step()` and before `optimizer.zero_grad()`.

Called in the training loop after taking an optimizer step and before zeroing grads. Good place to inspect weight information with weights updated.

This is where it is called:

```
for optimizer in optimizers:
    optimizer.step()
    model.on_before_zero_grad(optimizer) # < ---- called here
    optimizer.zero_grad()
```

Parameters `optimizer` (Optimizer) – The optimizer for which grads should be zeroed.

Return type `None`

on_fit_start

`ModelHooks.on_fit_start()`

Called at the very beginning of fit. If on DDP it is called on every process

on_fit_end

`ModelHooks.on_fit_end()`

Called at the very end of fit. If on DDP it is called on every process

on_load_checkpoint

`LightningModule.on_load_checkpoint(checkpoint)`

Do something with the checkpoint. Gives model a chance to load something before `state_dict` is restored.

Parameters `checkpoint` `(Dict[str, Any])` – A dictionary with variables from the checkpoint.

Return type `None`

on_save_checkpoint

`LightningModule.on_save_checkpoint(checkpoint)`

Give the model a chance to add something to the checkpoint. `state_dict` is already there.

Parameters `checkpoint` `(Dict[str, Any])` – A dictionary in which you can save variables to save in a checkpoint. Contents need to be pickleable.

Return type `None`

on_pretrain_routine_start

`ModelHooks.on_pretrain_routine_start()`

Called at the beginning of the pretrain routine (between fit and train start).

- fit
- pretrain_routine start
- pretrain_routine end
- training_start

Return type `None`

on_pretrain_routine_end

`ModelHooks.on_pretrain_routine_end()`

Called at the end of the pretrain routine (between fit and train start).

- fit
- pretrain_routine start
- pretrain_routine end
- training_start

Return type `None`

on_test_batch_start

`ModelHooks.on_test_batch_start(batch, batch_idx, dataloader_idx)`

Called in the test loop before anything happens for that batch.

Parameters

- **batch** `(Any)` – The batched data as it is returned by the training `DataLoader`.
- **batch_idx** `(int)` – the index of the batch
- **dataloader_idx** `(int)` – the index of the dataloader

Return type `None`

on_test_batch_end

`ModelHooks.on_test_batch_end(outputs, batch, batch_idx, dataloader_idx)`

Called in the test loop after the batch.

Parameters

- **outputs** `(Any)` – The outputs of `test_step_end(test_step(x))`
- **batch** `(Any)` – The batched data as it is returned by the training `DataLoader`.
- **batch_idx** `(int)` – the index of the batch
- **dataloader_idx** `(int)` – the index of the dataloader

Return type `None`

on_test_epoch_start

`ModelHooks.on_test_epoch_start()`

Called in the test loop at the very beginning of the epoch.

Return type `None`

on_test_epoch_end

`ModelHooks.on_test_epoch_end()`

Called in the test loop at the very end of the epoch.

Return type `None`

on_train_batch_start

`ModelHooks.on_train_batch_start(batch, batch_idx, dataloader_idx)`

Called in the training loop before anything happens for that batch.

If you return -1 here, you will skip training for the rest of the current epoch.

Parameters

- `batch` (Any) – The batched data as it is returned by the training DataLoader.
- `batch_idx` (int) – the index of the batch
- `dataloader_idx` (int) – the index of the dataloader

Return type `None`

on_train_batch_end

`ModelHooks.on_train_batch_end(outputs, batch, batch_idx, dataloader_idx)`

Called in the training loop after the batch.

Parameters

- `outputs` (Any) – The outputs of `validation_step_end(validation_step(x))`
- `batch` (Any) – The batched data as it is returned by the training DataLoader.
- `batch_idx` (int) – the index of the batch
- `dataloader_idx` (int) – the index of the dataloader

Return type `None`

on_train_epoch_start

`ModelHooks.on_train_epoch_start()`

Called in the training loop at the very beginning of the epoch.

Return type `None`

on_train_epoch_end

`ModelHooks.on_train_epoch_end(outputs)`

Called in the training loop at the very end of the epoch.

Return type `None`

on_validation_batch_start

`ModelHooks.on_validation_batch_start(batch, batch_idx, dataloader_idx)`

Called in the validation loop before anything happens for that batch.

Parameters

- `batch` (Any) – The batched data as it is returned by the training DataLoader.
- `batch_idx` (int) – the index of the batch
- `dataloader_idx` (int) – the index of the dataloader

Return type `None`

`on_validation_batch_end`

`ModelHooks.on_validation_batch_end(outputs, batch, batch_idx, dataloader_idx)`

Called in the validation loop after the batch.

Parameters

- `outputs` [\[Any\]](#) – The outputs of `validation_step_end(validation_step(x))`
- `batch` [\[Any\]](#) – The batched data as it is returned by the training `DataLoader`.
- `batch_idx` [\[int\]](#) – the index of the batch
- `dataloader_idx` [\[int\]](#) – the index of the dataloader

Return type `None`

`on_validation_epoch_start`

`ModelHooks.on_validation_epoch_start()`

Called in the validation loop at the very beginning of the epoch.

Return type `None`

`on_validation_epoch_end`

`ModelHooks.on_validation_epoch_end()`

Called in the validation loop at the very end of the epoch.

Return type `None`

`optimizer_step`

`LightningModule.optimizer_step(epoch, batch_idx, optimizer, optimizer_idx, second_order_closure=None, on_tpu=False, using_native_amp=False, using_lbfgs=False)`

Override this method to adjust the default way the `Trainer` calls each optimizer. By default, Lightning calls `step()` and `zero_grad()` as shown in the example once per optimizer.

Parameters

- `epoch` [\[int\]](#) – Current epoch
- `batch_idx` [\[int\]](#) – Index of current batch
- `optimizer` [\[Optimizer\]](#) – A PyTorch optimizer
- `optimizer_idx` [\[int\]](#) – If you used multiple optimizers this indexes into that list.
- `second_order_closure` [\[Optional\[Callable\]\]](#) – closure for second order methods
- `on_tpu` [\[bool\]](#) – true if TPU backward is required
- `using_native_amp` [\[bool\]](#) – True if using native amp
- `using_lbfgs` [\[bool\]](#) – True if the matching optimizer is lbfgs

Examples

```
# DEFAULT
def optimizer_step(self, current_epoch, batch_idx, optimizer, optimizer_idx,
                  second_order_closure, on_tpu, using_native_amp, using_lbfgs):
    optimizer.step()

# Alternating schedule for optimizer steps (i.e.: GANs)
def optimizer_step(self, current_epoch, batch_idx, optimizer, optimizer_idx,
                  second_order_closure, on_tpu, using_native_amp, using_lbfgs):
    # update generator opt every 2 steps
    if optimizer_idx == 0:
        if batch_idx % 2 == 0 :
            optimizer.step()
            optimizer.zero_grad()

    # update discriminator opt every 4 steps
    if optimizer_idx == 1:
        if batch_idx % 4 == 0 :
            optimizer.step()
            optimizer.zero_grad()

    # ...
    # add as many optimizers as you want
```

Here's another example showing how to use this for more advanced things such as learning rate warm-up:

```
# learning rate warm-up
def optimizer_step(self, current_epoch, batch_idx, optimizer,
                  optimizer_idx, second_order_closure, on_tpu, using_native_amp,
→ using_lbfgs):
    # warm up lr
    if self.trainer.global_step < 500:
        lr_scale = min(1., float(self.trainer.global_step + 1) / 500.)
        for pg in optimizer.param_groups:
            pg['lr'] = lr_scale * self.learning_rate

    # update params
    optimizer.step()
    optimizer.zero_grad()
```

Note: If you also override the `on_before_zero_grad()` model hook don't forget to add the call to it before `optimizer.zero_grad()` yourself.

Return type `None`

optimizer_zero_grad

LightningModule.**optimizer_zero_grad**(*epoch, batch_idx, optimizer, optimizer_idx*)

prepare_data

LightningModule.**prepare_data**()
Use this to download and prepare data.

Warning: DO NOT set state to the model (use *setup* instead) since this is NOT called on every GPU in DDP/TPU

Example:

```
def prepare_data(self):  
    # good  
    download_data()  
    tokenize()  
    etc()  
  
    # bad  
    self.split = data_split  
    self.some_state = some_other_state()
```

In DDP `prepare_data` can be called in two ways (using `Trainer(prepare_data_per_node)`):

1. Once per node. This is the default and is only called on `LOCAL_RANK=0`.
2. Once in total. Only called on `GLOBAL_RANK=0`.

Example:

```
# DEFAULT  
# called once per node on LOCAL_RANK=0 of that node  
Trainer(prepare_data_per_node=True)  
  
# call on GLOBAL_RANK=0 (great for shared file systems)  
Trainer(prepare_data_per_node=False)
```

This is called before requesting the dataloaders:

```
model.prepare_data()  
    if ddp/tpu: init()  
model.setup(stage)  
model.train_dataloader()  
model.val_dataloader()  
model.test_dataloader()
```

Return type `None`

setup

`ModelHooks.setup(stage)`

Called at the beginning of fit and test. This is a good hook when you need to build models dynamically or adjust something about them. This hook is called on every process when using DDP.

Parameters `stage` (`str`) – either ‘fit’ or ‘test’

Example:

```
class LitModel(...):
    def __init__(self):
        self.ll = None

    def prepare_data(self):
        download_data()
        tokenize()

        # don't do this
        self.something = else

    def setup(stage):
        data = Load_data(...)
        self.ll = nn.Linear(28, data.num_classes)
```

tbptt_split_batch

`LightningModule.tbptt_split_batch(batch, split_size)`

When using truncated backpropagation through time, each batch must be split along the time dimension. Lightning handles this by default, but for custom behavior override this function.

Parameters

- `batch` (`Tensor`) – Current batch
- `split_size` (`int`) – The size of the split

Return type `list`

Returns List of batch splits. Each split will be passed to `training_step()` to enable truncated back propagation through time. The default implementation splits root level Tensors and Sequences at `dim=1` (i.e. time dim). It assumes that each time dim is the same length.

Examples

```
def tbptt_split_batch(self, batch, split_size):
    splits = []
    for t in range(0, time_dims[0], split_size):
        batch_split = []
        for i, x in enumerate(batch):
            if isinstance(x, torch.Tensor):
                split_x = x[:, t:t + split_size]
            elif isinstance(x, collections.Sequence):
                split_x = [None] * len(x)
                for batch_idx in range(len(x)):
                    split_x[batch_idx] = x[batch_idx][t:t + split_size]
```

(continues on next page)

(continued from previous page)

```
        batch_split.append(split_x)

    splits.append(batch_split)

    return splits
```

Note: Called in the training loop after `on_batch_start()` if `truncated_bptt_steps > 0`. Each returned batch split is passed separately to `training_step()`.

teardown

`ModelHooks.teardown(stage)`

Called at the end of fit and test.

Parameters `stage` (`str`) – either ‘fit’ or ‘test’

train_dataloader

`LightningModule.train_dataloader()`

Implement a PyTorch `DataLoader` for training.

Return type `DataLoader`

Returns Single PyTorch `DataLoader`.

The `dataloader` you return will not be called every epoch unless you set `reload_dataloaders_every_epoch` to `True`.

For data processing use the following pattern:

- download in `prepare_data()`
- process and split in `setup()`

However, the above are only necessary for distributed processing.

Warning: do not assign state in `prepare_data`

- `fit()`
- ...
- `prepare_data()`
- `setup()`
- `train_dataloader()`

Note: Lightning adds the correct sampler for distributed and arbitrary hardware. There is no need to set it yourself.

Example

```
def train_dataloader(self):
    transform = transforms.Compose([transforms.ToTensor(),
                                    transforms.Normalize((0.5,), (1.0,))])
    dataset = MNIST(root='/path/to/mnist/', train=True, transform=transform,
                    download=True)
    loader = torch.utils.data.DataLoader(
        dataset=dataset,
        batch_size=self.batch_size,
        shuffle=True
    )
    return loader
```

val_dataloader

LightningModule.**val_dataloader**()

Implement one or multiple PyTorch DataLoaders for validation.

The dataloader you return will not be called every epoch unless you set `reload_dataloaders_every_epoch` to `True`.

It's recommended that all data downloads and preparation happen in `prepare_data()`.

- `fit()`
- ...
- `prepare_data()`
- `train_dataloader()`
- `val_dataloader()`
- `test_dataloader()`

Note: Lightning adds the correct sampler for distributed and arbitrary hardware. There is no need to set it yourself.

Return type `Union[DataLoader, List[DataLoader]]`

Returns Single or multiple PyTorch DataLoaders.

Examples

```
def val_dataloader(self):
    transform = transforms.Compose([transforms.ToTensor(),
                                    transforms.Normalize((0.5,), (1.0,))])
    dataset = MNIST(root='/path/to/mnist/', train=False,
                    transform=transform, download=True)
    loader = torch.utils.data.DataLoader(
        dataset=dataset,
        batch_size=self.batch_size,
        shuffle=False
    )
```

(continues on next page)

(continued from previous page)

```
    return loader

# can also return multiple dataloaders
def val_dataloader(self):
    return [loader_a, loader_b, ..., loader_n]
```

Note: If you don't need a validation dataset and a `validation_step()`, you don't need to implement this method.

Note: In the case where you return multiple validation dataloaders, the `validation_step()` will have an argument `dataloader_idx` which matches the order here.

test_dataloader

`LightningModule.test_dataloader()`

Implement one or multiple PyTorch DataLoaders for testing.

The dataloader you return will not be called every epoch unless you set `reload_dataloaders_every_epoch` to `True`.

For data processing use the following pattern:

- download in `prepare_data()`
- process and split in `setup()`

However, the above are only necessary for distributed processing.

Warning: do not assign state in `prepare_data`

- `fit()`
- ...
- `prepare_data()`
- `setup()`
- `train_dataloader()`
- `val_dataloader()`
- `test_dataloader()`

Note: Lightning adds the correct sampler for distributed and arbitrary hardware. There is no need to set it yourself.

Return type `Union[DataLoader, List[DataLoader]]`

Returns Single or multiple PyTorch DataLoaders.

Example

```
def test_dataloader(self):
    transform = transforms.Compose([transforms.ToTensor(),
                                    transforms.Normalize((0.5,), (1.0,))])
    dataset = MNIST(root='/path/to/mnist/', train=False, transform=transform,
                    download=True)
    loader = torch.utils.data.DataLoader(
        dataset=dataset,
        batch_size=self.batch_size,
        shuffle=False
    )

    return loader

# can also return multiple dataloaders
def test_dataloader(self):
    return [loader_a, loader_b, ..., loader_n]
```

Note: If you don't need a test dataset and a `test_step()`, you don't need to implement this method.

Note: In the case where you return multiple test dataloaders, the `test_step()` will have an argument `dataloader_idx` which matches the order here.

transfer_batch_to_device

`DataHooks.transfer_batch_to_device(batch, device)`

Override this hook if your `DataLoader` returns tensors wrapped in a custom data structure.

The data types listed below (and any arbitrary nesting of them) are supported out of the box:

- `torch.Tensor` or anything that implements `.to(...)`
- `list`
- `dict`
- `tuple`
- `torchtext.data.batch.Batch`

For anything else, you need to define how the data is moved to the target device (CPU, GPU, TPU, ...).

Example:

```
def transfer_batch_to_device(self, batch, device):
    if isinstance(batch, CustomBatch):
        # move all tensors in your custom data structure to the device
        batch.samples = batch.samples.to(device)
        batch.targets = batch.targets.to(device)
    else:
        batch = super().transfer_batch_to_device(data, device)
    return batch
```

Parameters

- **batch** (Any) – A batch of data that needs to be transferred to a new device.
- **device** (device) – The target device as defined in PyTorch.

Return type Any

Returns A reference to the data on the new device.

Note: This hook should only transfer the data and not modify it, nor should it move the data to any other device than the one passed in as argument (unless you know what you are doing).

Note: This hook only runs on single GPU training (no data-parallel). If you need multi-GPU support for your custom batch objects, you need to define your custom `DistributedDataParallel` or `LightningDistributedDataParallel` and override `configure_ddp()`.

See also:

- `move_data_to_device()`
- `apply_to_collection()`

TRAINER

Once you’ve organized your PyTorch code into a `LightningModule`, the Trainer automates everything else.

This abstraction achieves the following:

1. You maintain control over all aspects via PyTorch code without an added abstraction.
2. The trainer uses best practices embedded by contributors and users from top AI labs such as Facebook AI Research, NYU, MIT, Stanford, etc...
3. The trainer allows overriding any key part that you don’t want automated.

7.1 Basic use

This is the basic use of the trainer:

```
model = MyLightningModule()

trainer = Trainer()
trainer.fit(model, train_dataloader, val_dataloader)
```

7.2 Trainer in Python scripts

In Python scripts, it's recommended you use a main function to call the Trainer.

```
from argparse import ArgumentParser

def main(hparams):
    model = LightningModule()
    trainer = Trainer(gpus=hparams.gpus)
    trainer.fit(model)

if __name__ == '__main__':
    parser = ArgumentParser()
    parser.add_argument('--gpus', default=None)
    args = parser.parse_args()

    main(args)
```

So you can run it like so:

```
python main.py --gpus 2
```

Note: Pro-tip: You don't need to define all flags manually. Lightning can add them automatically

```
from argparse import ArgumentParser

def main(args):
    model = LightningModule()
    trainer = Trainer.from_argparse_args(args)
    trainer.fit(model)

if __name__ == '__main__':
    parser = ArgumentParser()
    parser = Trainer.add_argparse_args(parser)
    args = parser.parse_args()

    main(args)
```

So you can run it like so:

```
python main.py --gpus 2 --max_steps 10 --limit_train_batches 10 --any_trainer_arg x
```

Note: If you want to stop a training run early, you can press “Ctrl + C” on your keyboard. The trainer will catch the *KeyboardInterrupt* and attempt a graceful shutdown, including running callbacks such as *on_train_end*. The trainer object will also set an attribute *interrupted* to *True* in such cases. If you have a callback which shuts down compute resources, for example, you can conditionally run the shutdown logic for only uninterrupted runs.

7.3 Testing

Once you're done training, feel free to run the test set! (Only right before publishing your paper or pushing to production)

```
trainer.test(test_dataloader=test_dataloader)
```

7.4 Deployment / prediction

You just trained a LightningModule which is also just a torch.nn.Module. Use it to do whatever!

```
# load model
pretrained_model = LightningModule.load_from_checkpoint(PATH)
pretrained_model.freeze()

# use it for finetuning
def forward(self, x):
    features = pretrained_model(x)
    classes = classifier(features)

# or for prediction
out = pretrained_model(x)
api_write({'response': out})
```

You may wish to run the model on a variety of devices. Instead of moving the data manually to the correct device, decorate the forward method (or any other method you use for inference) with `auto_move_data()` and Lightning will take care of the rest.

7.5 Reproducibility

To ensure full reproducibility from run to run you need to set seeds for pseudo-random generators, and set deterministic flag in Trainer.

Example:

```
from pytorch_lightning import Trainer, seed_everything

seed_everything(42)
# sets seeds for numpy, torch, python.random and PYTHONHASHSEED.
model = Model()
trainer = Trainer(deterministic=True)
```

7.6 Trainer flags

7.6.1 accelerator

The accelerator backend to use (previously known as `distributed_backend`).

- (``dp``) is `DataParallel` (split batch among GPUs of same machine)
- (``ddp``) is `DistributedDataParallel` (each gpu on each node trains, and syncs grads)
- (``ddp_cpu``) is `DistributedDataParallel` on CPU (same as `ddp`, but does not use GPUs. Useful for multi-node CPU training or single-node debugging. Note that this will **not** give a speedup on a single node, since Torch already makes efficient use of multiple CPUs on a single machine.)
- (``ddp2``) **dp on node, ddp across nodes. Useful for things like increasing** the number of negative samples

```
# default used by the Trainer
trainer = Trainer(distributed_backend=None)
```

Example:

```
# dp = DataParallel
trainer = Trainer(gpus=2, distributed_backend='dp')

# ddp = DistributedDataParallel
trainer = Trainer(gpus=2, num_nodes=2, distributed_backend='ddp')

# ddp2 = DistributedDataParallel + dp
trainer = Trainer(gpus=2, num_nodes=2, distributed_backend='ddp2')
```

Note: this option does not apply to TPU. TPUs use ``ddp`` by default (over each core)

You can also modify hardware behavior by subclassing an existing accelerator to adjust for your needs.

Example:

```
class MyOwnDDP(DDPAccelerator):
    ...

Trainer(accelerator=MyOwnDDP())
```

Warning: Passing in custom accelerators is experimental but work is in progress to enable full compatibility.

7.6.2 accumulate_grad_batches

Accumulates grads every k batches or as set up in the dict. Trainer also calls `optimizer.step()` for the last indivisible step number.

```
# default used by the Trainer (no accumulation)
trainer = Trainer(accumulate_grad_batches=1)
```

Example:

```
# accumulate every 4 batches (effective batch size is batch*4)
trainer = Trainer(accumulate_grad_batches=4)

# no accumulation for epochs 1-4. accumulate 3 for epochs 5-10. accumulate 20 after_
↳that
trainer = Trainer(accumulate_grad_batches={5: 3, 10: 20})
```

7.6.3 amp_backend

Use PyTorch AMP ('native') (available PyTorch 1.6+), or NVIDIA apex ('apex').

```
# using PyTorch built-in AMP, default used by the Trainer
trainer = Trainer(amp_backend='native')

# using NVIDIA Apex
trainer = Trainer(amp_backend='apex')
```

7.6.4 amp_level

The optimization level to use (O1, O2, etc. ...) for 16-bit GPU precision (using NVIDIA apex under the hood).

Check [NVIDIA apex docs](#) for level

Example:

```
# default used by the Trainer
trainer = Trainer(amp_level='O2')
```

7.6.5 automatic_optimization

When set to False, Lightning does not automate the optimization process. This means you are responsible for your own optimizer behavior

Example:

```
def training_step(self, batch, batch_idx):
    opt = self.optimizers()

    loss = ...
    self.manual_backward(loss, opt)
    opt.step()
    opt.zero_grad()
```

This is not recommended when using a single optimizer, instead it's recommended when using 2+ optimizers AND you are an expert user. Most useful for research like RL, sparse coding and GAN research.

In the multi-optimizer case, ignore the optimizer_idx flag and use the optimizers directly

Example:

```
def training_step(self, batch, batch_idx, optimizer_idx):
    (opt_a, opt_b) = self.optimizers()

    gen_loss = ...
    self.manual_backward(gen_loss, opt_a)
    opt_a.step()
    opt_a.zero_grad()

    disc_loss = ...
    self.manual_backward(disc_loss, opt_b)
    opt_b.step()
    opt_b.zero_grad()
```

7.6.6 auto_scale_batch_size

Automatically tries to find the largest batch size that fits into memory, before any training.

```
# default used by the Trainer (no scaling of batch size)
trainer = Trainer(auto_scale_batch_size=None)

# run batch size scaling, result overrides hparams.batch_size
trainer = Trainer(auto_scale_batch_size='binsearch')

# call tune to find the batch size
trainer.tune(model)
```

7.6.7 auto_select_gpus

If enabled and *gpus* is an integer, pick available gpus automatically. This is especially useful when GPUs are configured to be in “exclusive mode”, such that only one process at a time can access them.

Example:

```
# no auto selection (picks first 2 gpus on system, may fail if other process is_
↳occupying)
trainer = Trainer(gpus=2, auto_select_gpus=False)

# enable auto selection (will find two available gpus on system)
trainer = Trainer(gpus=2, auto_select_gpus=True)
```

7.6.8 auto_lr_find

Runs a learning rate finder algorithm (see this [paper](#)) when calling `trainer.tune()`, to find optimal initial learning rate.

```
# default used by the Trainer (no learning rate finder)
trainer = Trainer(auto_lr_find=False)
```

Example:

```
# run learning rate finder, results override hparams.learning_rate
trainer = Trainer(auto_lr_find=True)

# call tune to find the lr
trainer.tune(model)
```

Example:

```
# run learning rate finder, results override hparams.my_lr_arg
trainer = Trainer(auto_lr_find='my_lr_arg')

# call tune to find the lr
trainer.tune(model)
```

Note: See the *learning rate finder guide*.

7.6.9 benchmark

If true enables `cuda.benchmark`. This flag is likely to increase the speed of your system if your input sizes don't change. However, if it does, then it will likely make your system slower.

The speedup comes from allowing the `cuda` auto-tuner to find the best algorithm for the hardware [\[see discussion here\]](#).

Example:

```
# default used by the Trainer
trainer = Trainer(benchmark=False)
```

7.6.10 deterministic

If true enables `cuda.deterministic`. Might make your system slower, but ensures reproducibility. Also sets `$HOROVOD_FUSION_THRESHOLD=0`.

For more info check [\[pytorch docs\]](#).

Example:

```
# default used by the Trainer
trainer = Trainer(deterministic=False)
```

7.6.11 callbacks

Add a list of *Callback*. These callbacks DO NOT replace the explicit callbacks (loggers or *ModelCheckpoint*).

Note: Only user defined callbacks (ie: Not *ModelCheckpoint*)

```
# a list of callbacks
callbacks = [PrintCallback()]
trainer = Trainer(callbacks=callbacks)
```

Example:

```
from pytorch_lightning.callbacks import Callback

class PrintCallback(Callback):
    def on_train_start(self, trainer, pl_module):
        print("Training is started!")
    def on_train_end(self, trainer, pl_module):
        print("Training is done.")
```


7.6.12 check_val_every_n_epoch

Check val every n train epochs.

Example:

```
# default used by the Trainer
trainer = Trainer(check_val_every_n_epoch=1)

# run val loop every 10 training epochs
trainer = Trainer(check_val_every_n_epoch=10)
```

7.6.13 checkpoint_callback

Pass in a callback for checkpointing. Checkpoints capture the exact value of all parameters used by a model. By default Lightning saves a checkpoint for you in your current working directory, with the state of your last training epoch, but you can override the default behavior by Initializing the *ModelCheckpoint* callback, and passing it to Trainer *checkpoint_callback* flag.

```
from pytorch_lightning.callbacks import ModelCheckpoint

# default used by the Trainer
checkpoint_callback = ModelCheckpoint(
    filepath=os.getcwd(),
    save_top_k=True,
    verbose=True,
    monitor='checkpoint_on',
    mode='min',
    prefix=''
)

trainer = Trainer(checkpoint_callback=checkpoint_callback)
```

To disable automatic checkpointing, set this to *False*.

```
trainer = Trainer(checkpoint_callback=False)
```

See also *Saving and Loading Weights*.

7.6.14 default_root_dir

Default path for logs and weights when no logger or `pytorch_lightning.callbacks.ModelCheckpoint` callback passed. On certain clusters you might want to separate where logs and checkpoints are stored. If you don't then use this argument for convenience. Paths can be local paths or remote paths such as `s3://bucket/path` or `hdfs://path/`. Credentials will need to be set up to use remote filepaths.

Example:

```
# default used by the Trainer
trainer = Trainer(default_root_path=os.getcwd())
```

7.6.15 distributed_backend

This has been renamed “accelerator”.

7.6.16 fast_dev_run

Runs 1 batch of train, test and val to find any bugs (ie: a sort of unit test).

Under the hood the pseudocode looks like this:

```
# loading
__init__()
prepare_data

# test training step
training_batch = next(train_dataloader)
training_step(training_batch)

# test val step
val_batch = next(val_dataloader)
out = validation_step(val_batch)
validation_epoch_end([out])
```

```
# default used by the Trainer
trainer = Trainer(fast_dev_run=False)

# runs 1 train, val, test batch and program ends
trainer = Trainer(fast_dev_run=True)
```

7.6.17 gpus

- Number of GPUs to train on (int)
- or which GPUs to train on (list)
- can handle strings

```
# default used by the Trainer (ie: train on CPU)
trainer = Trainer(gpus=None)

# equivalent
trainer = Trainer(gpus=0)
```

Example:

```
# int: train on 2 gpus
trainer = Trainer(gpus=2)

# list: train on GPUs 1, 4 (by bus ordering)
trainer = Trainer(gpus=[1, 4])
trainer = Trainer(gpus='1, 4') # equivalent

# -1: train on all gpus
trainer = Trainer(gpus=-1)
trainer = Trainer(gpus='-1') # equivalent

# combine with num_nodes to train on multiple GPUs across nodes
# uses 8 gpus in total
trainer = Trainer(gpus=2, num_nodes=4)

# train only on GPUs 1 and 4 across nodes
trainer = Trainer(gpus=[1, 4], num_nodes=4)
```

See also:

- [Multi-GPU training guide](#).

7.6.18 gradient_clip_val

Gradient clipping value

- 0 means don't clip.

```
# default used by the Trainer
trainer = Trainer(gradient_clip_val=0.0)
```

7.6.19 limit_test_batches

How much of test dataset to check.

```
# default used by the Trainer
trainer = Trainer(limit_test_batches=1.0)

# run through only 25% of the test set each epoch
trainer = Trainer(limit_test_batches=0.25)

# run for only 10 batches
trainer = Trainer(limit_test_batches=10)
```

In the case of multiple test dataloaders, the limit applies to each dataloader individually.

7.6.20 limit_val_batches

How much of validation dataset to check. Useful when debugging or testing something that happens at the end of an epoch.

```
# default used by the Trainer
trainer = Trainer(limit_val_batches=1.0)

# run through only 25% of the validation set each epoch
trainer = Trainer(limit_val_batches=0.25)

# run for only 10 batches
trainer = Trainer(limit_val_batches=10)
```

In the case of multiple validation dataloaders, the limit applies to each dataloader individually.

7.6.21 log_gpu_memory

Options:

- None
- 'min_max'
- 'all'

```
# default used by the Trainer
trainer = Trainer(log_gpu_memory=None)

# log all the GPUs (on master node only)
trainer = Trainer(log_gpu_memory='all')
```

(continues on next page)

(continued from previous page)

```
# log only the min and max memory on the master node
trainer = Trainer(log_gpu_memory='min_max')
```

Note: Might slow performance because it uses the output of nvidia-smi.

7.6.22 flush_logs_every_n_steps

Writes logs to disk this often.

```
# default used by the Trainer
trainer = Trainer(flush_logs_every_n_steps=100)
```

See also:

- [Logging](#)

7.6.23 logger

Logger (or iterable collection of loggers) for experiment tracking.

```
from pytorch_lightning.loggers import TensorBoardLogger

# default logger used by trainer
logger = TensorBoardLogger(
    save_dir=os.getcwd(),
    version=1,
    name='lightning_logs'
)
Trainer(logger=logger)
```

7.6.24 max_epochs

Stop training once this number of epochs is reached

```
# default used by the Trainer
trainer = Trainer(max_epochs=1000)
```

7.6.25 min_epochs

Force training for at least these many epochs

```
# default used by the Trainer
trainer = Trainer(min_epochs=1)
```

7.6.26 max_steps

Stop training after this number of steps Training will stop if max_steps or max_epochs have reached (earliest).

```
# Default (disabled)
trainer = Trainer(max_steps=None)

# Stop after 100 steps
trainer = Trainer(max_steps=100)
```

7.6.27 min_steps

Force training for at least these number of steps. Trainer will train model for at least min_steps or min_epochs (latest).

```
# Default (disabled)
trainer = Trainer(min_steps=None)

# Run at least for 100 steps (disable min_epochs)
trainer = Trainer(min_steps=100, min_epochs=0)
```

7.6.28 num_nodes

Number of GPU nodes for distributed training.

```
# default used by the Trainer
trainer = Trainer(num_nodes=1)

# to train on 8 nodes
trainer = Trainer(num_nodes=8)
```

7.6.29 num_processes

Number of processes to train with. Automatically set to the number of GPUs when using `distributed_backend="ddp"`. Set to a number greater than 1 when using `distributed_backend="ddp_cpu"` to mimic distributed training on a machine without GPUs. This is useful for debugging, but **will not** provide any speedup, since single-process Torch already makes efficient use of multiple CPUs.

```
# Simulate DDP for debugging on your GPU-less laptop
trainer = Trainer(distributed_backend="ddp_cpu", num_processes=2)
```

7.6.30 num_sanity_val_steps

Sanity check runs `n` batches of val before starting the training routine. This catches any bugs in your validation without having to wait for the first validation check. The Trainer uses 2 steps by default. Turn it off or modify it here.

```
# default used by the Trainer
trainer = Trainer(num_sanity_val_steps=2)

# turn it off
trainer = Trainer(num_sanity_val_steps=0)

# check all validation data
trainer = Trainer(num_sanity_val_steps=-1)
```

Example:

```
python -m torch_xla.distributed.xla_dist
--tpu=$TPU_POD_NAME
--conda-env=torch-xla-nightly
--env=XLA_USE_BF16=1
-- python your_trainer_file.py
```

7.6.31 plugins

Plugins allow you to connect arbitrary backends, precision libraries, SLURM, etc... For example:

- DDP
- SLURM
- TorchElastic
- Apex

To define your own behavior, subclass the relevant class and pass it in. Here's an example linking up your own cluster.

```
from pytorch_lightning.cluster_environments import cluster_environment

class MyCluster(ClusterEnvironment):

    def master_address(self):
        return your_master_address

    def master_port(self):
        return your_master_port

    def world_size(self):
        return the_world_size

trainer = Trainer(cluster_environment=cluster_environment())
```

7.6.32 prepare_data_per_node

If True will call `prepare_data()` on `LOCAL_RANK=0` for every node. If False will only call from `NODE_RANK=0`, `LOCAL_RANK=0`

```
# default
Trainer(prepare_data_per_node=True)

# use only NODE_RANK=0, LOCAL_RANK=0
Trainer(prepare_data_per_node=False)
```

7.6.33 tpu_cores

- How many TPU cores to train on (1 or 8).
- Which TPU core to train on [1-8]

A single TPU v2 or v3 has 8 cores. A TPU pod has up to 2048 cores. A slice of a POD means you get as many cores as you request.

Your effective batch size is `batch_size * total tpu cores`.

Note: No need to add a `DistributedDataSampler`, Lightning automatically does it for you.

This parameter can be either 1 or 8.

```
# your_trainer_file.py

# default used by the Trainer (ie: train on CPU)
trainer = Trainer(tpu_cores=None)

# int: train on a single core
```

(continues on next page)

(continued from previous page)

```

trainer = Trainer(tpu_cores=1)

# list: train on a single selected core
trainer = Trainer(tpu_cores=[2])

# int: train on all cores few cores
trainer = Trainer(tpu_cores=8)

# for 8+ cores must submit via xla script with
# a max of 8 cores specified. The XLA script
# will duplicate script onto each TPU in the POD
trainer = Trainer(tpu_cores=8)

```

To train on more than 8 cores (ie: a POD), submit this script using the xla_dist script.

Example:

```

python -m torch_xla.distributed.xla_dist
--tpu=$TPU_POD_NAME
--conda-env=torch-xla-nightly
--env=XLA_USE_BF16=1
-- python your_trainer_file.py

```

7.6.34 overfit_batches

Uses this much data of the training set. If nonzero, will use the same training set for validation and testing. If the training dataloaders have *shuffle=True*, Lightning will automatically disable it.

Useful for quickly debugging or trying to overfit on purpose.

```

# default used by the Trainer
trainer = Trainer(overfit_batches=0.0)

# use only 1% of the train set (and use the train set for val and test)
trainer = Trainer(overfit_batches=0.01)

# overfit on 10 of the same batches
trainer = Trainer(overfit_batches=10)

```

7.6.35 precision

Full precision (32), half precision (16). Can be used on CPU, GPU or TPUs.

If used on TPU will use torch.bfloat16 but tensor printing will still show torch.float32.

```

# default used by the Trainer
trainer = Trainer(precision=32)

```

(continues on next page)

(continued from previous page)

```
# 16-bit precision
trainer = Trainer(precision=16)
```

Example:

```
# one day
trainer = Trainer(precision=8|4|2)
```

7.6.36 process_position

Orders the progress bar. Useful when running multiple trainers on the same node.

```
# default used by the Trainer
trainer = Trainer(process_position=0)
```

Note: This argument is ignored if a custom callback is passed to `callbacks`.

7.6.37 profiler

To profile individual steps during training and assist in identifying bottlenecks.

See the *profiler documentation*. for more details.

```
from pytorch_lightning.profiler import SimpleProfiler, AdvancedProfiler

# default used by the Trainer
trainer = Trainer(profiler=None)

# to profile standard training events
trainer = Trainer(profiler=True)

# equivalent to profiler=True
trainer = Trainer(profiler=SimpleProfiler())

# advanced profiler for function-level stats
trainer = Trainer(profiler=AdvancedProfiler())
```

7.6.38 progress_bar_refresh_rate

How often to refresh progress bar (in steps). In notebooks, faster refresh rates (lower number) is known to crash them because of their screen refresh rates, so raise it to 50 or more.

```
# default used by the Trainer
trainer = Trainer(progress_bar_refresh_rate=1)

# disable progress bar
trainer = Trainer(progress_bar_refresh_rate=0)
```

Note: This argument is ignored if a custom callback is passed to `callbacks`.

7.6.39 reload_dataloaders_every_epoch

Set to `True` to reload dataloaders every epoch.

```
# if False (default)
train_loader = model.train_dataloader()
for epoch in epochs:
    for batch in train_loader:
        ...

# if True
for epoch in epochs:
    train_loader = model.train_dataloader()
    for batch in train_loader:
```

7.6.40 replace_sampler_ddp

Enables auto adding of distributed sampler. By default it will add `shuffle=True` for train sampler and `shuffle=False` for val/test sampler. If you want to customize it, you can set `replace_sampler_ddp=False` and add your own distributed sampler.

```
# default used by the Trainer
trainer = Trainer(replace_sampler_ddp=True)
```

By setting to `False`, you have to add your own distributed sampler:

```
# default used by the Trainer
sampler = torch.utils.data.distributed.DistributedSampler(dataset, shuffle=True)
dataloader = DataLoader(dataset, batch_size=32, sampler=sampler)
```

7.6.41 resume_from_checkpoint

To resume training from a specific checkpoint pass in the path here.

```
# default used by the Trainer
trainer = Trainer(resume_from_checkpoint=None)

# resume from a specific checkpoint
trainer = Trainer(resume_from_checkpoint='some/path/to/my_checkpoint.ckpt')
```

7.6.42 log_every_n_steps

How often to add logging rows (does not write to disk)

```
# default used by the Trainer
trainer = Trainer(log_every_n_steps=50)
```

See also:

- [Logging](#)

7.6.43 sync_batchnorm

Enable synchronization between batchnorm layers across all GPUs.

```
trainer = Trainer(sync_batchnorm=True)
```

7.6.44 track_grad_norm

- no tracking (-1)
- Otherwise tracks that norm (2 for 2-norm)

```
# default used by the Trainer
trainer = Trainer(track_grad_norm=-1)

# track the 2-norm
trainer = Trainer(track_grad_norm=2)
```

7.6.45 limit_train_batches

How much of training dataset to check. Useful when debugging or testing something that happens at the end of an epoch.

```
# default used by the Trainer
trainer = Trainer(limit_train_batches=1.0)
```

Example:

```
# default used by the Trainer
trainer = Trainer(limit_train_batches=1.0)

# run through only 25% of the training set each epoch
trainer = Trainer(limit_train_batches=0.25)

# run through only 10 batches of the training set each epoch
trainer = Trainer(limit_train_batches=10)
```

7.6.46 truncated_bptt_steps

Truncated back prop breaks performs backprop every k steps of a much longer sequence.

If this is enabled, your batches will automatically get truncated and the trainer will apply Truncated Backprop to it.

(Williams et al. “An efficient gradient-based algorithm for on-line training of recurrent network trajectories.”)

```
# default used by the Trainer (ie: disabled)
trainer = Trainer(truncated_bptt_steps=None)

# backprop every 5 steps in a batch
trainer = Trainer(truncated_bptt_steps=5)
```

Note: Make sure your batches have a sequence dimension.

Lightning takes care to split your batch along the time-dimension.

```
# we use the second as the time dimension
# (batch, time, ...)
sub_batch = batch[0, 0:t, ...]
```

Using this feature requires updating your LightningModule’s `pytorch_lightning.core.LightningModule.training_step()` to include a *hiddens* arg with the hidden

```
# Truncated back-propagation through time
def training_step(self, batch, batch_idx, hiddens):
    # hiddens are the hiddens from the previous truncated backprop step
    out, hiddens = self.lstm(data, hiddens)
```

(continues on next page)

(continued from previous page)

```
return {
    "loss": ...,
    "hiddens": hiddens  # remember to detach() this
}
```

To modify how the batch is split, override `pytorch_lightning.core.LightningModule.tbptt_split_batch()`:

```
class LitMNIST(LightningModule):
    def tbptt_split_batch(self, batch, split_size):
        # do your own splitting on the batch
        return splits
```

7.6.47 val_check_interval

How often within one training epoch to check the validation set. Can specify as float or int.

- use (float) to check within a training epoch
- use (int) to check every n steps (batches)

```
# default used by the Trainer
trainer = Trainer(val_check_interval=1.0)

# check validation set 4 times during a training epoch
trainer = Trainer(val_check_interval=0.25)

# check validation set every 1000 training batches
# use this when using IterableDataset and your dataset has no length
# (ie: production cases with streaming data)
trainer = Trainer(val_check_interval=1000)
```

7.6.48 weights_save_path

Directory of where to save weights if specified.

```
# default used by the Trainer
trainer = Trainer(weights_save_path=os.getcwd())

# save to your custom path
trainer = Trainer(weights_save_path='my/path')
```

Example:

```
# if checkpoint callback used, then overrides the weights path
# **NOTE: this saves weights to some/path NOT my/path
checkpoint = ModelCheckpoint(filepath='some/path')
trainer = Trainer(
    checkpoint_callback=checkpoint,
    weights_save_path='my/path'
)
```

7.6.49 weights_summary

Prints a summary of the weights when training begins. Options: 'full', 'top', None.

```
# default used by the Trainer (ie: print summary of top level modules)
trainer = Trainer(weights_summary='top')

# print full summary of all modules and submodules
trainer = Trainer(weights_summary='full')

# don't print a summary
trainer = Trainer(weights_summary=None)
```

7.7 Trainer class API

```
class pytorch_lightning.trainer.Trainer (logger=True, checkpoint_callback=True,
                                         callbacks=None, default_root_dir=None,
                                         gradient_clip_val=0, process_position=0,
                                         num_nodes=1, num_processes=1,
                                         gpus=None, auto_select_gpus=False,
                                         tpu_cores=None, log_gpu_memory=None,
                                         progress_bar_refresh_rate=1, over-
                                         fit_batches=0.0, track_grad_norm=-
                                         1, check_val_every_n_epoch=1,
                                         fast_dev_run=False, accumulate_grad_batches=1,
                                         max_epochs=1000, min_epochs=1,
                                         max_steps=None, min_steps=None,
                                         limit_train_batches=1.0, limit_val_batches=1.0,
                                         limit_test_batches=1.0, val_check_interval=1.0,
                                         flush_logs_every_n_steps=100,
                                         log_every_n_steps=50, accelera-
                                         tor=None, sync_batchnorm=False, pre-
                                         cision=32, weights_summary='top',
                                         weights_save_path=None,
                                         num_sanity_val_steps=2, trun-
                                         cated_bptt_steps=None, re-
                                         sume_from_checkpoint=None, profiler=None,
                                         benchmark=False, deterministic=False,
                                         reload_dataloaders_every_epoch=False,
                                         auto_lr_find=False, replace_sampler_ddp=True,
                                         terminate_on_nan=False,
                                         auto_scale_batch_size=False, pre-
                                         pare_data_per_node=True, plugins=None,
                                         amp_backend='native', amp_level='O2',
                                         distributed_backend=None, auto-
                                         matic_optimization=True)
```

Bases: `pytorch_lightning.trainer.properties.TrainerProperties`,
`pytorch_lightning.trainer.callback_hook.TrainerCallbackHookMixin`,
`pytorch_lightning.trainer.model_hooks.TrainerModelHooksMixin`,
`pytorch_lightning.trainer.optimizers.TrainerOptimizersMixin`,
`pytorch_lightning.trainer.logging.TrainerLoggingMixin`, `pytorch_lightning.`
`trainer.training_tricks.TrainerTrainingTricksMixin`, `pytorch_lightning.`
`trainer.data_loading.TrainerDataLoadingMixin`

Customize every aspect of training via flags

Parameters

- **accelerator** `//` (`Union[str, Accelerator, None]`) – Previously known as `distributed_backend` (`dp`, `ddp`, `ddp2`, etc...). Can also take in an accelerator object for custom hardware.
- **accumulate_grad_batches** `//` (`Union[int, Dict[int, int], List[list]]`) – Accumulates grads every k batches or as set up in the dict.
- **amp_backend** `//` (`str`) – The mixed precision backend to use (“native” or “apex”)
- **amp_level** `//` (`str`) – The optimization level to use (O1, O2, etc...).
- **auto_lr_find** `//` (`Union[bool, str]`) – If set to `True`, will make `trainer.tune()` run a

learning rate finder, trying to optimize initial learning for faster convergence. `trainer.tune()` method will set the suggested learning rate in `self.lr` or `self.learning_rate` in the LightningModule. To use a different key set a string instead of `True` with the key name.

- **`auto_scale_batch_size`** (Union[str, bool]) – If set to `True`, will *initially* run a batch size finder trying to find the largest batch size that fits into memory. The result will be stored in `self.batch_size` in the LightningModule. Additionally, can be set to either *power* that estimates the batch size through a power search or *binsearch* that estimates the batch size through a binary search.
- **`auto_select_gpus`** (bool) – If enabled and *gpus* is an integer, pick available gpus automatically. This is especially useful when GPUs are configured to be in “exclusive mode”, such that only one process at a time can access them.
- **`benchmark`** (bool) – If true enables `cuda.cudnn.benchmark`.
- **`callbacks`** (Optional[List[Callback]]) – Add a list of callbacks.
- **`checkpoint_callback`** (Union[ModelCheckpoint, bool]) – Callback for checkpointing.
- **`check_val_every_n_epoch`** (int) – Check val every *n* train epochs.
- **`default_root_dir`** (Optional[str]) – Default path for logs and weights when no logger/ckpt_callback passed. Default: `os.getcwd()`. Can be remote file paths such as `s3://mybucket/path` or `'hdfs://path'`
- **`deterministic`** (bool) – If true enables `cuda.cudnn.deterministic`.
- **`distributed_backend`** (Optional[str]) – deprecated. Please use ‘accelerator’
- **`fast_dev_run`** (bool) – runs 1 batch of train, test and val to find any bugs (ie: a sort of unit test).
- **`flush_logs_every_n_steps`** (int) – How often to flush logs to disk (defaults to every 100 steps).
- **`gpus`** (Union[int, str, List[int], None]) – number of gpus to train on (int) or which GPUs to train on (list or str) applied per node
- **`gradient_clip_val`** (float) – 0 means don’t clip.
- **`limit_train_batches`** (Union[int, float]) – How much of training dataset to check (floats = percent, int = num_batches)
- **`limit_val_batches`** (Union[int, float]) – How much of validation dataset to check (floats = percent, int = num_batches)
- **`limit_test_batches`** (Union[int, float]) – How much of test dataset to check (floats = percent, int = num_batches)
- **`logger`** (Union[LightningLoggerBase, Iterable[LightningLoggerBase], bool]) – Logger (or iterable collection of loggers) for experiment tracking.
- **`log_gpu_memory`** (Optional[str]) – None, ‘min_max’, ‘all’. Might slow performance
- **`log_every_n_steps`** (int) – How often to log within steps (defaults to every 50 steps).
- **`automatic_optimization`** (bool) – If False you are responsible for calling `.backward`, `.step`, `zero_grad`. Meant to be used with multiple optimizers by advanced users.

- **prepare_data_per_node** (bool) – If True, each LOCAL_RANK=0 will call prepare data. Otherwise only NODE_RANK=0, LOCAL_RANK=0 will prepare data
- **process_position** (int) – orders the progress bar when running multiple models on same machine.
- **progress_bar_refresh_rate** (int) – How often to refresh progress bar (in steps). Value 0 disables progress bar. Ignored when a custom callback is passed to `callbacks`.
- **profiler** (Union[BaseProfiler, bool, None]) – To profile individual steps during training and assist in identifying bottlenecks.
- **overfit_batches** (Union[int, float]) – Overfit a percent of training data (float) or a set number of batches (int). Default: 0.0
- **plugins** (Optional[list]) – Plugins allow modification of core behavior like ddp and amp.
- **precision** (int) – Full precision (32), half precision (16). Can be used on CPU, GPU or TPUs.
- **max_epochs** (int) – Stop training once this number of epochs is reached.
- **min_epochs** (int) – Force training for at least these many epochs
- **max_steps** (Optional[int]) – Stop training after this number of steps. Disabled by default (None).
- **min_steps** (Optional[int]) – Force training for at least these number of steps. Disabled by default (None).
- **num_nodes** (int) – number of GPU nodes for distributed training.
- **num_sanity_val_steps** (int) – Sanity check runs n validation batches before starting the training routine. Set it to -1 to run all batches in all validation dataloaders. Default: 2
- **reload_dataloaders_every_epoch** (bool) – Set to True to reload dataloaders every epoch.
- **replace_sampler_ddp** (bool) – Explicitly enables or disables sampler replacement. If not specified this will toggled automatically when DDP is used. By default it will add `shuffle=True` for train sampler and `shuffle=False` for val/test sampler. If you want to customize it, you can set `replace_sampler_ddp=False` and add your own distributed sampler.
- **resume_from_checkpoint** (Optional[str]) – To resume training from a specific checkpoint pass in the path here. This can be a URL.
- **sync_batchnorm** (bool) – Synchronize batch norm layers between process groups/whole world.
- **terminate_on_nan** (bool) – If set to True, will terminate training (by raising a *ValueError*) at the end of each training batch, if any of the parameters or the loss are NaN or +/-inf.
- **tpu_cores** (Union[int, str, List[int], None]) – How many TPU cores to train on (1 or 8) / Single TPU to train on [1]
- **track_grad_norm** (Union[int, float, str]) – -1 no tracking. Otherwise tracks that p-norm. May be set to 'inf' infinity-norm.
- **truncated_bptt_steps** (Optional[int]) – Truncated back prop breaks performs backprop every k steps of much longer sequence.

- **val_check_interval** (Union[int, float]) – How often to check the validation set. Use float to check within a training epoch, use int to check every n steps (batches).
- **weights_summary** (Optional[str]) – Prints a summary of the weights when training begins.
- **weights_save_path** (Optional[str]) – Where to save weights if specified. Will override default_root_dir for checkpoints only. Use this if for whatever reason you need the checkpoints stored in a different place than the logs written in default_root_dir. Can be remote file paths such as s3://mybucket/path or 'hdfs://path/' Defaults to default_root_dir.

fit (model, train_dataloader=None, val_dataloaders=None, datamodule=None)

Runs the full optimization routine.

Parameters

- **datamodule** (Optional[LightningDataModule]) – A instance of LightningDataModule.
- **model** (LightningModule) – Model to fit.
- **train_dataloader** (Optional[DataLoader]) – A Pytorch DataLoader with training samples. If the model has a predefined train_dataloader method this will be skipped.
- **val_dataloaders** (Union[DataLoader, List[DataLoader], None]) – Either a single Pytorch DataLoader or a list of them, specifying validation samples. If the model has a predefined val_dataloaders method this will be skipped

test (model=None, test_dataloaders=None, ckpt_path='best', verbose=True, datamodule=None)

Separates from fit to make sure you never run on your test set until you want to.

Parameters

- **ckpt_path** (Optional[str]) – Either best or path to the checkpoint you wish to test. If None, use the weights from the last epoch to test. Default to best.
- **datamodule** (Optional[LightningDataModule]) – A instance of LightningDataModule.
- **model** (Optional[LightningModule]) – The model to test.
- **test_dataloaders** (Union[DataLoader, List[DataLoader], None]) – Either a single Pytorch DataLoader or a list of them, specifying validation samples.
- **verbose** (bool) – If True, prints the test results

Returns The final test result dictionary. If no test_epoch_end is defined returns a list of dictionaries

tune (model, train_dataloader=None, val_dataloaders=None, datamodule=None)

Runs routines to tune hyperparameters before training.

Parameters

- **datamodule** (Optional[LightningDataModule]) – A instance of LightningDataModule.
- **model** (LightningModule) – Model to tune.
- **train_dataloader** (Optional[DataLoader]) – A Pytorch DataLoader with training samples. If the model has a predefined train_dataloader method this will be skipped.

- **val_dataloaders** `//` (`Union[DataLoader, List[DataLoader], None]`) – Either a single Pytorch Dataloader or a list of them, specifying validation samples. If the model has a predefined `val_dataloaders` method this will be skipped

`pytorch_lightning.trainer.seed_everything(seed=None)`

Function that sets seed for pseudo-random number generators in: `pytorch`, `numpy`, `python.random` In addition, sets the env variable `PL_GLOBAL_SEED` which will be passed to spawned subprocesses (e.g. `ddp_spawn` backend).

Parameters `seed` `//` (`Optional[int]`) – the integer value seed for global random state in Lightning. If `None`, will read seed from `PL_GLOBAL_SEED` env variable or select it randomly.

Return type `int`

CALLBACK

A callback is a self-contained program that can be reused across projects.

Lightning has a callback system to execute callbacks when needed. Callbacks should capture NON-ESSENTIAL logic that is NOT required for your *LightningModule* to run.

Here's the flow of how the callback hooks are executed:

An overall Lightning system should have:

1. Trainer for all engineering
2. LightningModule for all research code.
3. Callbacks for non-essential code.

Example:

```
from pytorch_lightning.callbacks import Callback

class MyPrintingCallback(Callback):

    def on_init_start(self, trainer):
        print('Starting to init trainer!')

    def on_init_end(self, trainer):
        print('trainer is init now')

    def on_train_end(self, trainer, pl_module):
        print('do something when training ends')

trainer = Trainer(callbacks=[MyPrintingCallback()])
```

```
Starting to init trainer!
trainer is init now
```

We successfully extended functionality without polluting our super clean *LightningModule* research code.

8.1 Examples

You can do pretty much anything with callbacks.

- Add a MLP to fine-tune self-supervised networks.
 - Find how to modify an image input to trick the classification result.
 - Interpolate the latent space of any variational model.
 - Log images to Tensorboard for any model.
-

8.2 Built-in Callbacks

Lightning has a few built-in callbacks.

Note: For a richer collection of callbacks, check out our [bolts library](#).

<i>Callback</i>	Abstract base class used to build new callbacks.
<i>EarlyStopping</i>	Monitor a validation metric and stop training when it stops improving.
<i>GPUStatsMonitor</i>	Automatically monitors and logs GPU stats during training stage.
<i>GradientAccumulationScheduler</i>	Change gradient accumulation factor according to scheduling.
<i>LearningRateMonitor</i>	Automatically monitor and logs learning rate for learning rate schedulers during training.
<i>ModelCheckpoint</i>	Save the model after every epoch by monitoring a quantity.
<i>ProgressBar</i>	This is the default progress bar used by Lightning.
<i>ProgressBarBase</i>	The base class for progress bars in Lightning.

8.2.1 Callback

class pytorch_lightning.callbacks.Callback

Bases: `abc.ABC`

Abstract base class used to build new callbacks.

Subclass this class and override any of the relevant hooks

on_batch_end (*trainer*, *pl_module*)

Called when the training batch ends.

on_batch_start (*trainer*, *pl_module*)

Called when the training batch begins.

on_epoch_end (*trainer*, *pl_module*)

Called when the epoch ends.

on_epoch_start (*trainer*, *pl_module*)

Called when the epoch begins.

on_fit_end (*trainer, pl_module*)
Called when fit ends

on_fit_start (*trainer, pl_module*)
Called when fit begins

on_init_end (*trainer*)
Called when the trainer initialization ends, model has not yet been set.

on_init_start (*trainer*)
Called when the trainer initialization begins, model has not yet been set.

on_keyboard_interrupt (*trainer, pl_module*)
Called when the training is interrupted by KeyboardInterrupt.

on_load_checkpoint (*checkpointed_state*)
Called when loading a model checkpoint, use to reload state.

on_pretrain_routine_end (*trainer, pl_module*)
Called when the pretrain routine ends.

on_pretrain_routine_start (*trainer, pl_module*)
Called when the pretrain routine begins.

on_sanity_check_end (*trainer, pl_module*)
Called when the validation sanity check ends.

on_sanity_check_start (*trainer, pl_module*)
Called when the validation sanity check starts.

on_save_checkpoint (*trainer, pl_module*)
Called when saving a model checkpoint, use to persist state.

on_test_batch_end (*trainer, pl_module, outputs, batch, batch_idx, dataloader_idx*)
Called when the test batch ends.

on_test_batch_start (*trainer, pl_module, batch, batch_idx, dataloader_idx*)
Called when the test batch begins.

on_test_end (*trainer, pl_module*)
Called when the test ends.

on_test_epoch_end (*trainer, pl_module*)
Called when the test epoch ends.

on_test_epoch_start (*trainer, pl_module*)
Called when the test epoch begins.

on_test_start (*trainer, pl_module*)
Called when the test begins.

on_train_batch_end (*trainer, pl_module, outputs, batch, batch_idx, dataloader_idx*)
Called when the train batch ends.

on_train_batch_start (*trainer, pl_module, batch, batch_idx, dataloader_idx*)
Called when the train batch begins.

on_train_end (*trainer, pl_module*)
Called when the train ends.

on_train_epoch_end (*trainer, pl_module, outputs*)
Called when the train epoch ends.

on_train_epoch_start (*trainer, pl_module*)
Called when the train epoch begins.

on_train_start (*trainer, pl_module*)
Called when the train begins.

on_validation_batch_end (*trainer, pl_module, outputs, batch, batch_idx, dataloader_idx*)
Called when the validation batch ends.

on_validation_batch_start (*trainer, pl_module, batch, batch_idx, dataloader_idx*)
Called when the validation batch begins.

on_validation_end (*trainer, pl_module*)
Called when the validation loop ends.

on_validation_epoch_end (*trainer, pl_module*)
Called when the val epoch ends.

on_validation_epoch_start (*trainer, pl_module*)
Called when the val epoch begins.

on_validation_start (*trainer, pl_module*)
Called when the validation loop begins.

setup (*trainer, pl_module, stage*)
Called when fit or test begins

teardown (*trainer, pl_module, stage*)
Called when fit or test ends

8.2.2 EarlyStopping

```
class pytorch_lightning.callbacks.EarlyStopping (monitor='early_stop_on',  
                                                min_delta=0.0,           patience=3,  
                                                verbose=False,         mode='auto',  
                                                strict=True)
```

Bases: `pytorch_lightning.callbacks.base.Callback`

Monitor a validation metric and stop training when it stops improving.

Parameters

- **monitor** `(str)` – quantity to be monitored. Default: `'early_stop_on'`.
- **min_delta** `(float)` – minimum change in the monitored quantity to qualify as an improvement, i.e. an absolute change of less than `min_delta`, will count as no improvement. Default: `0.0`.
- **patience** `(int)` – number of validation epochs with no improvement after which training will be stopped. Default: `3`.
- **verbose** `(bool)` – verbosity mode. Default: `False`.
- **mode** `(str)` – one of `{auto, min, max}`. In `min` mode, training will stop when the quantity monitored has stopped decreasing; in `max` mode it will stop when the quantity monitored has stopped increasing; in `auto` mode, the direction is automatically inferred from the name of the monitored quantity. Default: `'auto'`.
- **strict** `(bool)` – whether to crash the training if `monitor` is not found in the validation metrics. Default: `True`.

Example:


```
>>> from pytorch_lightning import Trainer
>>> from pytorch_lightning.callbacks import EarlyStopping
>>> early_stopping = EarlyStopping('val_loss')
>>> trainer = Trainer(callbacks=[early_stopping])
```

on_load_checkpoint (*checkpointed_state*)

Called when loading a model checkpoint, use to reload state.

on_save_checkpoint (*trainer, pl_module*)

Called when saving a model checkpoint, use to persist state.

on_train_epoch_end (*trainer, pl_module, outputs*)

Called when the train epoch ends.

on_validation_end (*trainer, pl_module*)

Called when the validation loop ends.

on_validation_epoch_end (*trainer, pl_module*)

Called when the val epoch ends.

8.2.3 GPUStatsMonitor

```
class pytorch_lightning.callbacks.GPUStatsMonitor (memory_utilization=True,
                                                    gpu_utilization=True,          in-
                                                    tra_step_time=False,             in-
                                                    ter_step_time=False,
                                                    fan_speed=False,                 tempera-
                                                    ture=False)
```

Bases: `pytorch_lightning.callbacks.base.Callback`

Automatically monitors and logs GPU stats during training stage. `GPUStatsMonitor` is a callback and in order to use it you need to assign a logger in the `Trainer`.

Parameters

- **memory_utilization** `(bool)` – Set to `True` to monitor used, free and percentage of memory utilization at the start and end of each step. Default: `True`.
- **gpu_utilization** `(bool)` – Set to `True` to monitor percentage of GPU utilization at the start and end of each step. Default: `True`.
- **intra_step_time** `(bool)` – Set to `True` to monitor the time of each step. Default: `False`.
- **inter_step_time** `(bool)` – Set to `True` to monitor the time between the end of one step and the start of the next step. Default: `False`.
- **fan_speed** `(bool)` – Set to `True` to monitor percentage of fan speed. Default: `False`.
- **temperature** `(bool)` – Set to `True` to monitor the memory and gpu temperature in degree Celsius. Default: `False`.

Example:

```
>>> from pytorch_lightning import Trainer
>>> from pytorch_lightning.callbacks import GPUStatsMonitor
>>> gpu_stats = GPUStatsMonitor()
>>> trainer = Trainer(callbacks=[gpu_stats])
```

GPU stats are mainly based on `nvidia-smi --query-gpu` command. The description of the queries is as follows:

- **fan.speed** – The fan speed value is the percent of maximum speed that the device’s fan is currently intended to run at. It ranges from 0 to 100 %. Note: The reported speed is the intended fan speed. If the fan is physically blocked and unable to spin, this output will not match the actual fan speed. Many parts do not report fan speeds because they rely on cooling via fans in the surrounding enclosure.
- **memory.used** – Total memory allocated by active contexts.
- **memory.free** – Total free memory.
- **utilization.gpu** – Percent of time over the past sample period during which one or more kernels was executing on the GPU. The sample period may be between 1 second and 1/6 second depending on the product.
- **utilization.memory** – Percent of time over the past sample period during which global (device) memory was being read or written. The sample period may be between 1 second and 1/6 second depending on the product.
- **temperature.gpu** – Core GPU temperature, in degrees C.
- **temperature.memory** – HBM memory temperature, in degrees C.

on_train_batch_end (*trainer, pl_module, outputs, batch, batch_idx, dataloader_idx*)

Called when the train batch ends.

on_train_batch_start (*trainer, pl_module, batch, batch_idx, dataloader_idx*)

Called when the train batch begins.

on_train_epoch_start (*trainer, pl_module*)

Called when the train epoch begins.

on_train_start (*trainer, pl_module*)

Called when the train begins.

8.2.4 GradientAccumulationScheduler

class `pytorch_lightning.callbacks.GradientAccumulationScheduler` (*scheduling*)

Bases: `pytorch_lightning.callbacks.base.Callback`

Change gradient accumulation factor according to scheduling.

Parameters *scheduling* (`Dict[int, int]`) – scheduling in format {epoch: accumulation_factor}

Example:

```
>>> from pytorch_lightning import Trainer
>>> from pytorch_lightning.callbacks import GradientAccumulationScheduler

# at epoch 5 start accumulating every 2 batches
>>> accumulator = GradientAccumulationScheduler(scheduling={5: 2})
>>> trainer = Trainer(callbacks=[accumulator])

# alternatively, pass the scheduling dict directly to the Trainer
>>> trainer = Trainer(accumulate_grad_batches={5: 2})
```

on_epoch_start (*trainer, pl_module*)

Called when the epoch begins.

8.2.5 LearningRateMonitor

class pytorch_lightning.callbacks.**LearningRateMonitor** (*logging_interval=None*)
 Bases: pytorch_lightning.callbacks.base.Callback

Automatically monitor and logs learning rate for learning rate schedulers during training.

Parameters **logging_interval** (Optional[str]) – set to *epoch* or *step* to log *lr* of all optimizers at the same interval, set to *None* to log at individual interval according to the *interval* key of each scheduler. Defaults to *None*.

Example:

```
>>> from pytorch_lightning import Trainer
>>> from pytorch_lightning.callbacks import LearningRateMonitor
>>> lr_monitor = LearningRateMonitor(logging_interval='step')
>>> trainer = Trainer(callbacks=[lr_monitor])
```

Logging names are automatically determined based on optimizer class name. In case of multiple optimizers of same type, they will be named *Adam*, *Adam-1* etc. If a optimizer has multiple parameter groups they will be named *Adam/pg1*, *Adam/pg2* etc. To control naming, pass in a *name* keyword in the construction of the learning rate schedulers

Example:

```
def configure_optimizer(self):
    optimizer = torch.optim.Adam(...)
    lr_scheduler = {'scheduler': torch.optim.lr_schedulers.LambdaLR(optimizer, ...
→)
                    'name': 'my_logging_name'}
    return [optimizer], [lr_scheduler]
```

on_batch_start (*trainer, pl_module*)
 Called when the training batch begins.

on_epoch_start (*trainer, pl_module*)
 Called when the epoch begins.

on_train_start (*trainer, pl_module*)
 Called before training, determines unique names for all *lr* schedulers in the case of multiple of the same type or in the case of multiple parameter groups

8.2.6 ModelCheckpoint

class pytorch_lightning.callbacks.**ModelCheckpoint** (*filepath=None, monitor=None, verbose=False, save_last=None, save_top_k=None, save_weights_only=False, mode='auto', period=1, prefix=""*)
 Bases: pytorch_lightning.callbacks.base.Callback

Save the model after every epoch by monitoring a quantity.

After training finishes, use *best_model_path* to retrieve the path to the best checkpoint file and *best_model_score* to retrieve its score.

Parameters

- **filepath** (Optional[str]) – path to save the model file. Can contain named formatting options to be auto-filled.

Example:

```
# custom path
# saves a file like: my/path/epoch=0.ckpt
>>> checkpoint_callback = ModelCheckpoint('my/path/')

# save any arbitrary metrics like `val_loss`, etc. in name
# saves a file like: my/path/epoch=2-val_loss=0.02-other_metric=0.
# 03.ckpt
>>> checkpoint_callback = ModelCheckpoint(
...     filepath='my/path/{epoch}-{val_loss:.2f}-{other_metric:.2f}'
... )
```

By default, filepath is *None* and will be set at runtime to the location specified by Trainer's default_root_dir or weights_save_path arguments, and if the Trainer uses a logger, the path will also contain logger name and version.

- **monitor** (Optional[str]) – quantity to monitor. By default it is *None* which saves a checkpoint only for the last epoch
- **verbose** (bool) – verbosity mode. Default: *False*.
- **save_last** (Optional[bool]) – When *True*, always saves the model at the end of the epoch to a file *last.ckpt*. Default: *None*.
- **save_top_k** (Optional[int]) – if save_top_k == k, the best k models according to the quantity monitored will be saved. if save_top_k == 0, no models are saved. if save_top_k == -1, all models are saved. Please note that the monitors are checked every *period* epochs. if save_top_k >= 2 and the callback is called multiple times inside an epoch, the name of the saved file will be appended with a version count starting with v0.
- **mode** (str) – one of {auto, min, max}. If save_top_k != 0, the decision to overwrite the current save file is made based on either the maximization or the minimization of the monitored quantity. For *val_acc*, this should be *max*, for *val_loss* this should be *min*, etc. In *auto* mode, the direction is automatically inferred from the name of the monitored quantity.
- **save_weights_only** (bool) – if *True*, then only the model's weights will be saved (model.save_weights(filepath)), else the full model is saved (model.save(filepath)).
- **period** (int) – Interval (number of epochs) between checkpoints.

Example:

```
>>> from pytorch_lightning import Trainer
>>> from pytorch_lightning.callbacks import ModelCheckpoint

# saves checkpoints to 'my/path/' at every epoch
>>> checkpoint_callback = ModelCheckpoint(filepath='my/path/')
>>> trainer = Trainer(checkpoint_callback=checkpoint_callback)

# save epoch and val_loss in name
# saves a file like: my/path/sample-mnist-epoch=02-val_loss=0.32.ckpt
>>> checkpoint_callback = ModelCheckpoint(monitor='val_loss',
```

(continues on next page)

(continued from previous page)

```

...     filepath='my/path/sample-mnist-{epoch:02d}-{val_loss:.2f}'
... )

# retrieve the best checkpoint after training
checkpoint_callback = ModelCheckpoint(filepath='my/path/')
trainer = Trainer(checkpoint_callback=checkpoint_callback)
model = ...
trainer.fit(model)
checkpoint_callback.best_model_path

```

format_checkpoint_name (*epoch, metrics, ver=None*)

Generate a filename according to the defined template.

Example:

```

>>> tmpdir = os.path.dirname(__file__)
>>> ckpt = ModelCheckpoint(os.path.join(tmpdir, '{epoch}'))
>>> os.path.basename(ckpt.format_checkpoint_name(0, {}))
'epoch=0.ckpt'
>>> ckpt = ModelCheckpoint(os.path.join(tmpdir, '{epoch:03d}'))
>>> os.path.basename(ckpt.format_checkpoint_name(5, {}))
'epoch=005.ckpt'
>>> ckpt = ModelCheckpoint(os.path.join(tmpdir, '{epoch}-{val_loss:.2f}'))
>>> os.path.basename(ckpt.format_checkpoint_name(2, dict(val_loss=0.123456)))
'epoch=2-val_loss=0.12.ckpt'
>>> ckpt = ModelCheckpoint(os.path.join(tmpdir, '{missing:d}'))
>>> os.path.basename(ckpt.format_checkpoint_name(0, {}))
'missing=0.ckpt'

```

Return type `str`

on_load_checkpoint (*checkpointed_state*)

Called when loading a model checkpoint, use to reload state.

on_pretrain_routine_start (*trainer, pl_module*)

When pretrain routine starts we build the ckpt dir on the fly

on_save_checkpoint (*trainer, pl_module*)

Called when saving a model checkpoint, use to persist state.

Return type `Dict[str, Any]`

on_validation_end (*trainer, pl_module*)

checkpoints can be saved at the end of the val loop

save_checkpoint (*trainer, pl_module*)

Performs the main logic around saving a checkpoint. This method runs on all ranks, it is the responsibility of *self.save_function* to handle correct behaviour in distributed training, i.e., saving only on rank 0.

to_yaml (*filepath=None*)

Saves the *best_k_models* dict containing the checkpoint paths with the corresponding scores to a YAML file.

8.2.7 ProgressBar

`class pytorch_lightning.callbacks.ProgressBar` (*refresh_rate=1, process_position=0*)
 Bases: `pytorch_lightning.callbacks.progress.ProgressBarBase`

This is the default progress bar used by Lightning. It prints to *stdout* using the `tqdm` package and shows up to four different bars:

- **sanity check progress:** the progress during the sanity check run
- **main progress:** shows training + validation progress combined. It also accounts for multiple validation runs during training when `val_check_interval` is used.
- **validation progress:** only visible during validation; shows total progress over all validation datasets.
- **test progress:** only active when testing; shows total progress over all test datasets.

For infinite datasets, the progress bar never ends.

If you want to customize the default `tqdm` progress bars used by Lightning, you can override specific methods of the callback class and pass your custom implementation to the `Trainer`:

Example:

```
class LitProgressBar(ProgressBar):

    def init_validation_tqdm(self):
        bar = super().init_validation_tqdm()
        bar.set_description('running validation ...')
        return bar

bar = LitProgressBar()
trainer = Trainer(callbacks=[bar])
```

Parameters

- **refresh_rate** (*int*) – Determines at which rate (in number of batches) the progress bars get updated. Set it to 0 to disable the display. By default, the `Trainer` uses this implementation of the progress bar and sets the refresh rate to the value provided to the `progress_bar_refresh_rate` argument in the `Trainer`.
- **process_position** (*int*) – Set this to a value greater than 0 to offset the progress bars by this many lines. This is useful when you have progress bars defined elsewhere and want to show all of them together. This corresponds to `process_position` in the `Trainer`.

`disable()`

You should provide a way to disable the progress bar. The `Trainer` will call this to disable the output on processes that have a rank different from 0, e.g., in multi-node training.

Return type `None`

`enable()`

You should provide a way to enable the progress bar. The `Trainer` will call this in e.g. pre-training routines like the *learning rate finder* to temporarily enable and disable the main progress bar.

Return type `None`

`init_sanity_tqdm()`

Override this to customize the `tqdm` bar for the validation sanity run.

Return type `tqdm`

`init_test_tqdm()`

Override this to customize the tqdm bar for testing.

Return type `tqdm`

`init_train_tqdm()`

Override this to customize the tqdm bar for training.

Return type `tqdm`

`init_validation_tqdm()`

Override this to customize the tqdm bar for validation.

Return type `tqdm`

`on_epoch_start(trainer, pl_module)`

Called when the epoch begins.

`on_sanity_check_end(trainer, pl_module)`

Called when the validation sanity check ends.

`on_sanity_check_start(trainer, pl_module)`

Called when the validation sanity check starts.

`on_test_batch_end(trainer, pl_module, outputs, batch, batch_idx, dataloader_idx)`

Called when the test batch ends.

`on_test_end(trainer, pl_module)`

Called when the test ends.

`on_test_start(trainer, pl_module)`

Called when the test begins.

`on_train_batch_end(trainer, pl_module, outputs, batch, batch_idx, dataloader_idx)`

Called when the train batch ends.

`on_train_end(trainer, pl_module)`

Called when the train ends.

`on_train_start(trainer, pl_module)`

Called when the train begins.

`on_validation_batch_end(trainer, pl_module, outputs, batch, batch_idx, dataloader_idx)`

Called when the validation batch ends.

`on_validation_end(trainer, pl_module)`

Called when the validation loop ends.

`on_validation_start(trainer, pl_module)`

Called when the validation loop begins.

8.2.8 ProgressBarBase

`class pytorch_lightning.callbacks.ProgressBarBase`

Bases: `pytorch_lightning.callbacks.base.Callback`

The base class for progress bars in Lightning. It is a [Callback](#) that keeps track of the batch progress in the Trainer. You should implement your highly custom progress bars with this as the base class.

Example:

```

class LitProgressBar(ProgressBarBase):

    def __init__(self):
        super().__init__() # don't forget this :)
        self.enable = True

    def disable(self):
        self.enable = False

    def on_train_batch_end(self, trainer, pl_module, outputs):
        super().on_train_batch_end(trainer, pl_module, outputs) # don't forget_
↪this :)
        percent = (self.train_batch_idx / self.total_train_batches) * 100
        sys.stdout.flush()
        sys.stdout.write(f'{percent:.01f} percent complete \r')

bar = LitProgressBar()
trainer = Trainer(callbacks=[bar])

```

disable()

You should provide a way to disable the progress bar. The `Trainer` will call this to disable the output on processes that have a rank different from 0, e.g., in multi-node training.

enable()

You should provide a way to enable the progress bar. The `Trainer` will call this in e.g. pre-training routines like the *learning rate finder* to temporarily enable and disable the main progress bar.

on_epoch_start(trainer, pl_module)

Called when the epoch begins.

on_init_end(trainer)

Called when the trainer initialization ends, model has not yet been set.

on_test_batch_end(trainer, pl_module, outputs, batch, batch_idx, dataloader_idx)

Called when the test batch ends.

on_test_start(trainer, pl_module)

Called when the test begins.

on_train_batch_end(trainer, pl_module, outputs, batch, batch_idx, dataloader_idx)

Called when the train batch ends.

on_train_start(trainer, pl_module)

Called when the train begins.

on_validation_batch_end(trainer, pl_module, outputs, batch, batch_idx, dataloader_idx)

Called when the validation batch ends.

on_validation_start(trainer, pl_module)

Called when the validation loop begins.

property test_batch_idx

The current batch index being processed during testing. Use this to update your progress bar.

Return type `int`

property total_test_batches

The total number of training batches during testing, which may change from epoch to epoch. Use this to set the total number of iterations in the progress bar. Can return `inf` if the test dataloader is of infinite size.

Return type `int`

property `total_train_batches`

The total number of training batches during training, which may change from epoch to epoch. Use this to set the total number of iterations in the progress bar. Can return `inf` if the training dataloader is of infinite size.

Return type `int`

property `total_val_batches`

The total number of training batches during validation, which may change from epoch to epoch. Use this to set the total number of iterations in the progress bar. Can return `inf` if the validation dataloader is of infinite size.

Return type `int`

property `train_batch_idx`

The current batch index being processed during training. Use this to update your progress bar.

Return type `int`

property `val_batch_idx`

The current batch index being processed during validation. Use this to update your progress bar.

Return type `int`

8.3 Persisting State

Some callbacks require internal state in order to function properly. You can optionally choose to persist your callback's state as part of model checkpoint files using the callback hooks `on_save_checkpoint()` and `on_load_checkpoint()`. However, you must follow two constraints:

1. Your returned state must be able to be pickled.
2. You can only use one instance of that class in the Trainer callbacks list. We don't support persisting state for multiple callbacks of the same class.

8.4 Best Practices

The following are best practices when using/designing callbacks.

1. Callbacks should be isolated in their functionality.
 2. Your callback should not rely on the behavior of other callbacks in order to work properly.
 3. Do not manually call methods from the callback.
 4. Directly calling methods (eg. `on_validation_end`) is strongly discouraged.
 5. Whenever possible, your callbacks should not depend on the order in which they are executed.
-

8.5 Available Callback hooks

8.5.1 on_epoch_start

`ModelHooks.on_epoch_start()`

Called in the training loop at the very beginning of the epoch.

Return type `None`

8.5.2 on_epoch_end

`ModelHooks.on_epoch_end()`

Called in the training loop at the very end of the epoch.

Return type `None`

8.5.3 on_fit_start

`ModelHooks.on_fit_start()`

Called at the very beginning of fit. If on DDP it is called on every process

8.5.4 on_fit_end

`ModelHooks.on_fit_end()`

Called at the very end of fit. If on DDP it is called on every process

8.5.5 on_save_checkpoint

`CheckpointHooks.on_save_checkpoint(checkpoint)`

Called by Lightning when saving a checkpoint to give you a chance to store anything else you might want to save.

Parameters `checkpoint` `(Dict[str, Any])` – Checkpoint to be saved

Example

```
def on_save_checkpoint(self, checkpoint):  
    # 99% of use cases you don't need to implement this method  
    checkpoint['something_cool_i_want_to_save'] = my_cool_pickable_object
```

Note: Lightning saves all aspects of training (epoch, global step, etc...) including amp scaling. There is no need for you to store anything about training.

Return type `None`

8.5.6 on_load_checkpoint

CheckpointHooks.on_load_checkpoint (checkpoint)

Called by Lightning to restore your model. If you saved something with on_save_checkpoint() this is your chance to restore this.

Parameters checkpoint (Dict[str, Any]) – Loaded checkpoint

Example

```
def on_load_checkpoint(self, checkpoint):
    # 99% of the time you don't need to implement this method
    self.something_cool_i_want_to_save = checkpoint['something_cool_i_want_to_save']
```

Note: Lightning auto-restores global step, epoch, and train state including amp scaling. There is no need for you to restore anything regarding training.

Return type None

8.5.7 on_pretrain_routine_start

ModelHooks.on_pretrain_routine_start()

Called at the beginning of the pretrain routine (between fit and train start).

- fit
- pretrain_routine start
- pretrain_routine end
- training_start

Return type None

8.5.8 on_pretrain_routine_end

ModelHooks.on_pretrain_routine_end()

Called at the end of the pretrain routine (between fit and train start).

- fit
- pretrain_routine start
- pretrain_routine end
- training_start

Return type None

8.5.9 on_test_batch_start

`ModelHooks.on_test_batch_start(batch, batch_idx, dataloader_idx)`

Called in the test loop before anything happens for that batch.

Parameters

- `batch` (Any) – The batched data as it is returned by the training DataLoader.
- `batch_idx` (int) – the index of the batch
- `dataloader_idx` (int) – the index of the dataloader

Return type None

8.5.10 on_test_batch_end

`ModelHooks.on_test_batch_end(outputs, batch, batch_idx, dataloader_idx)`

Called in the test loop after the batch.

Parameters

- `outputs` (Any) – The outputs of `test_step_end(test_step(x))`
- `batch` (Any) – The batched data as it is returned by the training DataLoader.
- `batch_idx` (int) – the index of the batch
- `dataloader_idx` (int) – the index of the dataloader

Return type None

8.5.11 on_test_epoch_start

`ModelHooks.on_test_epoch_start()`

Called in the test loop at the very beginning of the epoch.

Return type None

8.5.12 on_test_epoch_end

`ModelHooks.on_test_epoch_end()`

Called in the test loop at the very end of the epoch.

Return type None

8.5.13 on_test_model_train

`ModelHooks.on_test_model_train()`

Sets the model to train during the test loop

Return type None

8.5.14 on_test_model_eval

`ModelHooks.on_test_model_eval()`
Sets the model to eval during the test loop

Return type `None`

8.5.15 on_train_batch_start

`ModelHooks.on_train_batch_start(batch, batch_idx, dataloader_idx)`
Called in the training loop before anything happens for that batch.

If you return -1 here, you will skip training for the rest of the current epoch.

Parameters

- **batch** `(Any)` – The batched data as it is returned by the training DataLoader.
- **batch_idx** `(int)` – the index of the batch
- **dataloader_idx** `(int)` – the index of the dataloader

Return type `None`

8.5.16 on_train_batch_end

`ModelHooks.on_train_batch_end(outputs, batch, batch_idx, dataloader_idx)`
Called in the training loop after the batch.

Parameters

- **outputs** `(Any)` – The outputs of `validation_step_end(validation_step(x))`
- **batch** `(Any)` – The batched data as it is returned by the training DataLoader.
- **batch_idx** `(int)` – the index of the batch
- **dataloader_idx** `(int)` – the index of the dataloader

Return type `None`

8.5.17 on_train_start

`ModelHooks.on_train_start()`
Called at the beginning of training before sanity check.

Return type `None`

8.5.18 on_train_end

`ModelHooks.on_train_end()`

Called at the end of training before logger experiment is closed.

Return type `None`

8.5.19 on_train_epoch_start

`ModelHooks.on_train_epoch_start()`

Called in the training loop at the very beginning of the epoch.

Return type `None`

8.5.20 on_train_epoch_end

`ModelHooks.on_train_epoch_end(outputs)`

Called in the training loop at the very end of the epoch.

Return type `None`

8.5.21 on_validation_batch_start

`ModelHooks.on_validation_batch_start(batch, batch_idx, dataloader_idx)`

Called in the validation loop before anything happens for that batch.

Parameters

- **batch** `(Any)` – The batched data as it is returned by the training `Dataloader`.
- **batch_idx** `(int)` – the index of the batch
- **dataloader_idx** `(int)` – the index of the dataloader

Return type `None`

8.5.22 on_validation_batch_end

`ModelHooks.on_validation_batch_end(outputs, batch, batch_idx, dataloader_idx)`

Called in the validation loop after the batch.

Parameters

- **outputs** `(Any)` – The outputs of `validation_step_end(validation_step(x))`
- **batch** `(Any)` – The batched data as it is returned by the training `Dataloader`.
- **batch_idx** `(int)` – the index of the batch
- **dataloader_idx** `(int)` – the index of the dataloader

Return type `None`

8.5.23 on_validation_epoch_start

`ModelHooks.on_validation_epoch_start()`

Called in the validation loop at the very beginning of the epoch.

Return type `None`

8.5.24 on_validation_epoch_end

`ModelHooks.on_validation_epoch_end()`

Called in the validation loop at the very end of the epoch.

Return type `None`

8.5.25 on_validation_model_eval

`ModelHooks.on_validation_model_eval()`

Sets the model to eval during the val loop

Return type `None`

8.5.26 on_validation_model_train

`ModelHooks.on_validation_model_train()`

Sets the model to train during the val loop

Return type `None`

8.5.27 setup

`ModelHooks.setup(stage)`

Called at the beginning of fit and test. This is a good hook when you need to build models dynamically or adjust something about them. This hook is called on every process when using DDP.

Parameters `stage` (`str`) – either ‘fit’ or ‘test’

Example:

```
class LitModel(...):
    def __init__(self):
        self.ll = None

    def prepare_data(self):
        download_data()
        tokenize()

        # don't do this
        self.something = else

    def setup(stage):
        data = Load_data(...)
        self.ll = nn.Linear(28, data.num_classes)
```

8.5.28 teardown

`ModelHooks.teardown(stage)`

Called at the end of fit and test.

Parameters `stage` (`str`) – either ‘fit’ or ‘test’

LIGHTNINGDATAMODULE

A datamodule is a shareable, reusable class that encapsulates all the steps needed to process data:

A datamodule encapsulates the five steps involved in data processing in PyTorch:

1. Download / tokenize / process.
2. Clean and (maybe) save to disk.
3. Load inside `Dataset`.
4. Apply transforms (rotate, tokenize, etc...).
5. Wrap inside a `DataLoader`.

This class can then be shared and used anywhere:

```
from pl_bolts.datamodules import CIFAR10DataModule, ImagenetDataModule

model = LitClassifier()
trainer = Trainer()

imagenet = ImagenetDataModule()
trainer.fit(model, imagenet)

cifar10 = CIFAR10DataModule()
trainer.fit(model, cifar10)
```

9.1 Why do I need a DataModule?

In normal PyTorch code, the data cleaning/preparation is usually scattered across many files. This makes sharing and reusing the exact splits and transforms across projects impossible.

Datamodules are for you if you ever asked the questions:

- what splits did you use?
 - what transforms did you use?
 - what normalization did you use?
 - how did you prepare/tokenize the data?
-

9.2 What is a DataModule

A DataModule is simply a collection of a `train_dataloader`, `val_dataloader(s)`, `test_dataloader(s)` along with the matching transforms and data processing/downloads steps required.

Here's a simple PyTorch example:

```
# regular PyTorch
test_data = MNIST(PATH, train=False, download=True)
train_data = MNIST(PATH, train=True, download=True)
train_data, val_data = random_split(train_data, [55000, 5000])

train_loader = DataLoader(train_data, batch_size=32)
val_loader = DataLoader(val_data, batch_size=32)
test_loader = DataLoader(test_data, batch_size=32)
```

The equivalent DataModule just organizes the same exact code, but makes it reusable across projects.

```
class MNISTDataModule(pl.LightningDataModule):

    def __init__(self, data_dir: str = PATH, batch_size):
        super().__init__()
        self.batch_size = batch_size

    def setup(self, stage=None):
        self.mnist_test = MNIST(self.data_dir, train=False)
        mnist_full = MNIST(self.data_dir, train=True)
        self.mnist_train, self.mnist_val = random_split(mnist_full, [55000, 5000])

    def train_dataloader(self):
        return DataLoader(self.mnist_train, batch_size=self.batch_size)

    def val_dataloader(self):
        return DataLoader(self.mnist_val, batch_size=self.batch_size)

    def test_dataloader(self):
        return DataLoader(self.mnist_test, batch_size=self.batch_size)
```

But now, as the complexity of your processing grows (transforms, multiple-GPU training), you can let Lightning handle those details for you while making this dataset reusable so you can share with colleagues or use in different projects.

```
mnist = MNISTDataModule(PATH)
model = LitClassifier()

trainer = Trainer()
trainer.fit(model, mnist)
```

Here's a more realistic, complex DataModule that shows how much more reusable the datamodule is.

```
import pytorch_lightning as pl
from torch.utils.data import random_split, DataLoader

# Note - you must have torchvision installed for this example
from torchvision.datasets import MNIST
from torchvision import transforms

class MNISTDataModule(pl.LightningDataModule):

    def __init__(self, data_dir: str = '.'):
        super().__init__()
        self.data_dir = data_dir
        self.transform = transforms.Compose([
            transforms.ToTensor(),
            transforms.Normalize((0.1307,), (0.3081,))
        ])

        # self.dims is returned when you call dm.size()
        # Setting default dims here because we know them.
        # Could optionally be assigned dynamically in dm.setup()
        self.dims = (1, 28, 28)

    def prepare_data(self):
        # download
        MNIST(self.data_dir, train=True, download=True)
        MNIST(self.data_dir, train=False, download=True)

    def setup(self, stage=None):

        # Assign train/val datasets for use in dataloaders
        if stage == 'fit' or stage is None:
            mnist_full = MNIST(self.data_dir, train=True, transform=self.transform)
            self.mnist_train, self.mnist_val = random_split(mnist_full, [55000, 5000])

            # Optionally...
            # self.dims = tuple(self.mnist_train[0][0].shape)

        # Assign test dataset for use in dataloader(s)
        if stage == 'test' or stage is None:
            self.mnist_test = MNIST(self.data_dir, train=False, transform=self.
↪transform)

            # Optionally...
            # self.dims = tuple(self.mnist_test[0][0].shape)

    def train_dataloader(self):
        return DataLoader(self.mnist_train, batch_size=32)
```

(continues on next page)

(continued from previous page)

```
def val_dataloader(self):
    return DataLoader(self.mnist_val, batch_size=32)

def test_dataloader(self):
    return DataLoader(self.mnist_test, batch_size=32)
```

Note: `setup` expects a string arg `stage`. It is used to separate setup logic for `trainer.fit` and `trainer.test`.

9.3 LightningDataModule API

To define a `DataModule` define 5 methods:

- `prepare_data` (how to download(), tokenize, etc...)
- `setup` (how to split, etc...)
- `train_dataloader`
- `val_dataloader(s)`
- `test_dataloader(s)`

9.3.1 prepare_data

Use this method to do things that might write to disk or that need to be done only from a single GPU in distributed settings.

- download
- tokenize
- etc...

```
class MNISTDataModule(pl.LightningDataModule):
    def prepare_data(self):
        # download
        MNIST(os.getcwd(), train=True, download=True, transform=transforms.ToTensor())
        MNIST(os.getcwd(), train=False, download=True, transform=transforms.
↪ ToTensor())
```

Warning: `prepare_data` is called from a single GPU. Do not use it to assign state (`self.x = y`).

9.3.2 setup

There are also data operations you might want to perform on every GPU. Use `setup` to do things like:

- count number of classes
- build vocabulary
- perform train/val/test splits
- apply transforms (defined explicitly in your datamodule or assigned in `init`)
- etc...

```
import pytorch_lightning as pl

class MNISTDataModule(pl.LightningDataModule):

    def setup(self, stage: Optional[str] = None):

        # Assign Train/val split(s) for use in Dataloaders
        if stage == 'fit' or stage is None:
            mnist_full = MNIST(
                self.data_dir,
                train=True,
                download=True,
                transform=self.transform
            )
            self.mnist_train, self.mnist_val = random_split(mnist_full, [55000, 5000])
            self.dims = self.mnist_train[0][0].shape

        # Assign Test split(s) for use in Dataloaders
        if stage == 'test' or stage is None:
            self.mnist_test = MNIST(
                self.data_dir,
                train=False,
                download=True,
                transform=self.transform
            )
            self.dims = getattr(self, 'dims', self.mnist_test[0][0].shape)
```

Warning: `setup` is called from every GPU. Setting state here is okay.

9.3.3 train_dataloader

Use this method to generate the train dataloader. Usually you just wrap the dataset you defined in `setup`.

```
import pytorch_lightning as pl

class MNISTDataModule(pl.LightningDataModule):
    def train_dataloader(self):
        return DataLoader(self.mnist_train, batch_size=64)
```

9.3.4 val_dataloader

Use this method to generate the val dataloader. Usually you just wrap the dataset you defined in `setup`.

```
import pytorch_lightning as pl

class MNISTDataModule(pl.LightningDataModule):
    def val_dataloader(self):
        return DataLoader(self.mnist_val, batch_size=64)
```

9.3.5 test_dataloader

Use this method to generate the test dataloader. Usually you just wrap the dataset you defined in `setup`.

```
import pytorch_lightning as pl

class MNISTDataModule(pl.LightningDataModule):
    def test_dataloader(self):
        return DataLoader(self.mnist_test, batch_size=64)
```

9.3.6 transfer_batch_to_device

Override to define how you want to move an arbitrary batch to a device

```
import pytorch_lightning as pl

class MNISTDataModule(pl.LightningDataModule):
    def transfer_batch_to_device(self, batch, device):
        x = batch['x']
        x = CustomDataWrapper(x)
        batch['x'].to(device)
        return batch
```

Note: To decouple your data from transforms you can parametrize them via `__init__`.

```
class MNISTDataModule(pl.LightningDataModule):
    def __init__(self, train_transforms, val_transforms, test_transforms):
        super().__init__()
        self.train_transforms = train_transforms
        self.val_transforms = val_transforms
        self.test_transforms = test_transforms
```

9.4 Using a DataModule

The recommended way to use a DataModule is simply:

```
dm = MNISTDataModule()
model = Model()
trainer.fit(model, dm)

trainer.test(datamodule=dm)
```

If you need information from the dataset to build your model, then run *prepare_data* and *setup* manually (Lightning still ensures the method runs on the correct devices)

```
dm = MNISTDataModule()
dm.prepare_data()
dm.setup('fit')

model = Model(num_classes=dm.num_classes, width=dm.width, vocab=dm.vocab)
trainer.fit(model, dm)

dm.setup('test')
trainer.test(datamodule=dm)
```

9.5 Datamodules without Lightning

You can of course use DataModules in plain PyTorch code as well.

```
# download, etc...
dm = MNISTDataModule()
dm.prepare_data()

# splits/transforms
dm.setup('fit')

# use data
for batch in dm.train_dataloader():
    ...
for batch in dm.val_dataloader():
    ...

# lazy load test data
dm.setup('test')
for batch in dm.test_dataloader():
    ...
```

But overall, DataModules encourage reproducibility by allowing all details of a dataset to be specified in a unified structure.

LOGGING

Lightning supports the most popular logging frameworks (TensorBoard, Comet, etc...). To use a logger, simply pass it into the `Trainer`. Lightning uses TensorBoard by default.

```
from pytorch_lightning import loggers as pl_loggers

tb_logger = pl_loggers.TensorBoardLogger('logs/')
trainer = Trainer(logger=tb_logger)
```

Choose from any of the others such as MLflow, Comet, Neptune, WandB, ...

```
comet_logger = pl_loggers.CometLogger(save_dir='logs/')
trainer = Trainer(logger=comet_logger)
```

To use multiple loggers, simply pass in a list or tuple of loggers ...

```
tb_logger = pl_loggers.TensorBoardLogger('logs/')
comet_logger = pl_loggers.CometLogger(save_dir='logs/')
trainer = Trainer(logger=[tb_logger, comet_logger])
```

Note: By default, lightning logs every 50 steps. Use Trainer flags to *Control logging frequency*.

Note: All loggers log by default to `os.getcwd()`. To change the path without creating a logger set `Trainer(default_root_dir='/your/path/to/save/checkpoints')`

10.1 Logging from a LightningModule

Lightning offers automatic log functionalities for logging scalars, or manual logging for anything else.

10.1.1 Automatic logging

Use the `log()` method to log from anywhere in a *LightningModule*.

```
def training_step(self, batch, batch_idx):
    self.log('my_metric', x)
```

Depending on where log is called from, Lightning auto-determines the correct logging mode for you. But of course you can override the default behavior by manually setting the `log()` parameters.

```
def training_step(self, batch, batch_idx):
    self.log('my_loss', loss, on_step=True, on_epoch=True, prog_bar=True, logger=True)
```

The `log()` method has a few options:

- *on_step*: Logs the metric at the current step. Defaults to *True* in `training_step()`, and `training_step_end()`.
- *on_epoch*: Automatically accumulates and logs at the end of the epoch. Defaults to *True* anywhere in validation or test loops, and in `training_epoch_end()`.
- *prog_bar*: Logs to the progress bar.
- *logger*: Logs to the logger like Tensorboard, or any other custom logger passed to the Trainer.

Note: Setting *on_epoch=True* will accumulate your logged values over the full training epoch.

10.1.2 Manual logging

If you want to log anything that is not a scalar, like histograms, text, images, etc... you may need to use the logger object directly.

```
def training_step(...):
    ...
    # the logger you used (in this case tensorboard)
    tensorboard = self.logger.experiment
    tensorboard.add_image()
    tensorboard.add_histogram(...)
    tensorboard.add_figure(...)
```

10.1.3 Access your logs

Once your training starts, you can view the logs by using your favorite logger or booting up the Tensorboard logs:

```
tensorboard --logdir ./lightning_logs
```

10.2 Make a custom logger

You can implement your own logger by writing a class that inherits from `LightningLoggerBase`. Use the `rank_zero_only()` decorator to make sure that only the first process in DDP training logs data.

```
from pytorch_lightning.utilities import rank_zero_only
from pytorch_lightning.loggers import LightningLoggerBase

class MyLogger(LightningLoggerBase):

    def name(self):
        return 'MyLogger'

    def experiment(self):
        # Return the experiment object associated with this logger.
        pass

    def version(self):
        # Return the experiment version, int or str.
        return '0.1'

    @rank_zero_only
    def log_hyperparams(self, params):
        # params is an argparse.Namespace
        # your code to record hyperparameters goes here
        pass

    @rank_zero_only
    def log_metrics(self, metrics, step):
        # metrics is a dictionary of metric names and values
        # your code to record metrics goes here
        pass

    def save(self):
        # Optional. Any code necessary to save logger data goes here
        # If you implement this, remember to call `super().save()`
        # at the start of the method (important for aggregation of metrics)
        super().save()

    @rank_zero_only
    def finalize(self, status):
        # Optional. Any code that needs to be run after training
        # finishes goes here
        pass
```

If you write a logger that may be useful to others, please send a pull request to add it to Lightning!

10.3 Control logging frequency

10.3.1 Logging frequency

It may slow training down to log every single batch. By default, Lightning logs every 50 rows, or 50 training steps. To change this behaviour, set the `log_every_n_steps` Trainer flag.

```
k = 10
trainer = Trainer(log_every_n_steps=k)
```

10.3.2 Log writing frequency

Writing to a logger can be expensive, so by default Lightning write logs to disc or to the given logger every 100 training steps. To change this behaviour, set the interval at which you wish to flush logs to the filesystem using `log_every_n_steps` Trainer flag.

```
k = 100
trainer = Trainer(flush_logs_every_n_steps=k)
```

Unlike the `log_every_n_steps`, this argument does not apply to all loggers. The example shown here works with `TensorBoardLogger`, which is the default logger in Lightning.

10.4 Progress Bar

You can add any metric to the progress bar using `log()` method, setting `prog_bar=True`.

```
def training_step(self, batch, batch_idx):
    self.log('my_loss', loss, prog_bar=True)
```

10.4.1 Modifying the progress bar

The progress bar by default already includes the training loss and version number of the experiment if you are using a logger. These defaults can be customized by overriding the `get_progress_bar_dict()` hook in your module.

```
def get_progress_bar_dict(self):
    # don't show the version number
    items = super().get_progress_bar_dict()
    items.pop("v_num", None)
    return items
```

10.5 Configure console logging

Lightning logs useful information about the training process and user warnings to the console. You can retrieve the Lightning logger and change it to your liking. For example, increase the logging level to see fewer messages like so:

```
import logging
logging.getLogger("lightning").setLevel(logging.ERROR)
```

Read more about custom Python logging [here](#).

10.6 Logging hyperparameters

When training a model, it's useful to know what hyperparams went into that model. When Lightning creates a checkpoint, it stores a key "hparams" with the hyperparams.

```
lightning_checkpoint = torch.load(filepath, map_location=lambda storage, loc: storage)
hyperparams = lightning_checkpoint['hparams']
```

Some loggers also allow logging the hyperparams used in the experiment. For instance, when using the TestTubeLogger or the TensorBoardLogger, all hyperparams will show in the [hparams tab](#).

10.7 Snapshot code

Loggers also allow you to snapshot a copy of the code used in this experiment. For example, TestTubeLogger does this with a flag:

```
from pytorch_lightning.loggers import TestTubeLogger
logger = TestTubeLogger('.', create_git_tag=True)
```

10.8 Supported Loggers

The following are loggers we support

10.8.1 Comet

```
class pytorch_lightning.loggers.comet.CometLogger(
    api_key=None, save_dir=None,
    project_name=None,
    rest_api_key=None, experiment_name=None,
    experiment_key=None, offline=False,
    **kwargs)
Bases: pytorch_lightning.loggers.base.LightningLoggerBase
```

Log using [Comet.ml](#). Install it with pip:

```
pip install comet-ml
```

Comet requires either an API Key (online mode) or a local directory path (offline mode).

ONLINE MODE

```
import os
from pytorch_lightning import Trainer
from pytorch_lightning.loggers import CometLogger
# arguments made to CometLogger are passed on to the comet_ml.Experiment class
comet_logger = CometLogger(
    api_key=os.environ.get('COMET_API_KEY'),
    workspace=os.environ.get('COMET_WORKSPACE'), # Optional
    save_dir='.', # Optional
    project_name='default_project', # Optional
    rest_api_key=os.environ.get('COMET_REST_API_KEY'), # Optional
    experiment_name='default' # Optional
)
trainer = Trainer(logger=comet_logger)
```

OFFLINE MODE

```
from pytorch_lightning.loggers import CometLogger
# arguments made to CometLogger are passed on to the comet_ml.Experiment class
comet_logger = CometLogger(
    save_dir='.',
    workspace=os.environ.get('COMET_WORKSPACE'), # Optional
    project_name='default_project', # Optional
    rest_api_key=os.environ.get('COMET_REST_API_KEY'), # Optional
    experiment_name='default' # Optional
)
trainer = Trainer(logger=comet_logger)
```

Parameters

- **api_key** (Optional[str]) – Required in online mode. API key, found on Comet.ml. If not given, this will be loaded from the environment variable `COMET_API_KEY` or `~/comet.config` if either exists.
- **save_dir** (Optional[str]) – Required in offline mode. The path for the directory to save local comet logs. If given, this also sets the directory for saving checkpoints.
- **project_name** (Optional[str]) – Optional. Send your experiment to a specific project. Otherwise will be sent to Uncategorized Experiments. If the project name does not already exist, Comet.ml will create a new project.
- **rest_api_key** (Optional[str]) – Optional. Rest API key found in Comet.ml settings. This is used to determine version number
- **experiment_name** (Optional[str]) – Optional. String representing the name for this particular experiment on Comet.ml.
- **experiment_key** (Optional[str]) – Optional. If set, restores from existing experiment.
- **offline** (bool) – If `api_key` and `save_dir` are both given, this determines whether the experiment will be in online or offline mode. This is useful if you use `save_dir` to control the checkpoints directory and have a `~/comet.config` file but still want to run offline experiments.

- ****kwargs** – Additional arguments like *workspace*, *log_code*, etc. used by CometExperiment can be passed as keyword arguments in this logger.

finalize (*status*)

When calling `self.experiment.end()`, that experiment won't log any more data to Comet. That's why, if you need to log any more data, you need to create an ExistingCometExperiment. For example, to log data when testing your model after training, because when training is finalized CometLogger.`finalize()` is called.

This happens automatically in the `experiment()` property, when `self._experiment` is set to `None`, i.e. `self.reset_experiment()`.

Return type `None`

log_hyperparams (*params*)

Record hyperparameters.

Parameters **params** (`Union[Dict[str, Any], Namespace]`) – `Namespace` containing the hyperparameters

Return type `None`

log_metrics (*metrics*, *step=None*)

Records metrics. This method logs metrics as as soon as it received them. If you want to aggregate metrics for one specific *step*, use the `agg_and_log_metrics()` method.

Parameters

- **metrics** (`Dict[str, Union[Tensor, float]]`) – Dictionary with metric names as keys and measured quantities as values
- **step** (`Optional[int]`) – Step number at which the metrics should be recorded

Return type `None`

property experiment

Actual Comet object. To use Comet features in your LightningModule do the following.

Example:

```
self.logger.experiment.some_comet_function()
```

property name

Return the experiment name.

Return type `str`

property save_dir

Return the root directory where experiment logs get saved, or `None` if the logger does not save data locally.

Return type `Optional[str]`

property version

Return the experiment version.

Return type `str`

10.8.2 CSVLogger

class `pytorch_lightning.loggers.csv_logs.CSVLogger` (*save_dir*, *name*='default', *version*=None)
Bases: `pytorch_lightning.loggers.base.LightningLoggerBase`
Log to local file system in yaml and CSV format. Logs are saved to `os.path.join(save_dir, name, version)`.

Example

```
>>> from pytorch_lightning import Trainer
>>> from pytorch_lightning.loggers import CSVLogger
>>> logger = CSVLogger("logs", name="my_exp_name")
>>> trainer = Trainer(logger=logger)
```

Parameters

- **save_dir** *(str)* – Save directory
- **name** *(Optional[str])* – Experiment name. Defaults to 'default'.
- **version** *(Union[int, str, None])* – Experiment version. If version is not specified the logger inspects the save directory for existing versions, then automatically assigns the next available version.

finalize (*status*)

Do any processing that is necessary to finalize an experiment.

Parameters **status** *(str)* – Status that the experiment finished with (e.g. success, failed, aborted)

Return type `None`

log_hyperparams (*params*)

Record hyperparameters.

Parameters **params** *(Union[Dict[str, Any], Namespace])* – `Namespace` containing the hyperparameters

Return type `None`

log_metrics (*metrics*, *step*=None)

Records metrics. This method logs metrics as as soon as it received them. If you want to aggregate metrics for one specific *step*, use the `agg_and_log_metrics()` method.

Parameters

- **metrics** *(Dict[str, float])* – Dictionary with metric names as keys and measured quantities as values
- **step** *(Optional[int])* – Step number at which the metrics should be recorded

Return type `None`

save ()

Save log data.

Return type `None`

property experiment

Actual ExperimentWriter object. To use ExperimentWriter features in your LightningModule do the following.

Example:

```
self.logger.experiment.some_experiment_writer_function()
```

Return type ExperimentWriter

property log_dir

The log directory for this run. By default, it is named 'version_\${self.version}' but it can be overridden by passing a string value for the constructor's version parameter instead of None or an int.

Return type str

property name

Return the experiment name.

Return type str

property root_dir

Parent directory for all checkpoint subdirectories. If the experiment name parameter is None or the empty string, no experiment subdirectory is used and the checkpoint will be saved in "save_dir/version_dir"

Return type str

property save_dir

Return the root directory where experiment logs get saved, or None if the logger does not save data locally.

Return type Optional[str]

property version

Return the experiment version.

Return type int

10.8.3 MLFlow

```
class pytorch_lightning.loggers.mlflow.MLFlowLogger(experiment_name='default',
                                                    tracking_uri=None, tags=None,
                                                    save_dir='./mlruns')
```

Bases: pytorch_lightning.loggers.base.LightningLoggerBase

Log using MLflow. Install it with pip:

```
pip install mlflow
```

```
from pytorch_lightning import Trainer
from pytorch_lightning.loggers import MLFlowLogger
mlf_logger = MLFlowLogger(
    experiment_name="default",
    tracking_uri="file:./ml-runs"
)
trainer = Trainer(logger=mlf_logger)
```

Use the logger anywhere in your LightningModule as follows:

```
from pytorch_lightning import LightningModule
class LitModel(LightningModule):
    def training_step(self, batch, batch_idx):
        # example
        self.logger.experiment.whatever_ml_flow_supports(...)

    def any_lightning_module_function_or_hook(self):
        self.logger.experiment.whatever_ml_flow_supports(...)
```

Parameters

- **experiment_name** (str) – The name of the experiment
- **tracking_uri** (Optional[str]) – Address of local or remote tracking server. If not provided, defaults to `file:<save_dir>`.
- **tags** (Optional[Dict[str, Any]]) – A dictionary tags for the experiment.
- **save_dir** (Optional[str]) – A path to a local directory where the MLflow runs get saved. Defaults to `.mlflow` if `tracking_uri` is not provided. Has no effect if `tracking_uri` is provided.

finalize (status='FINISHED')

Do any processing that is necessary to finalize an experiment.

Parameters **status** (str) – Status that the experiment finished with (e.g. success, failed, aborted)

Return type None

log_hyperparams (params)

Record hyperparameters.

Parameters **params** (Union[Dict[str, Any], Namespace]) – Namespace containing the hyperparameters

Return type None

log_metrics (metrics, step=None)

Records metrics. This method logs metrics as as soon as it received them. If you want to aggregate metrics for one specific `step`, use the `agg_and_log_metrics()` method.

Parameters

- **metrics** (Dict[str, float]) – Dictionary with metric names as keys and measured quantities as values
- **step** (Optional[int]) – Step number at which the metrics should be recorded

Return type None

property **experiment**

Actual MLflow object. To use MLflow features in your `LightningModule` do the following.

Example:

```
self.logger.experiment.some_mlflow_function()
```

Return type MlflowClient

property **name**

Return the experiment name.

Return type `str`

property `save_dir`

The root file directory in which MLflow experiments are saved.

Return type `Optional[str]`

Returns Local path to the root experiment directory if the tracking uri is local. Otherwise returns *None*.

property `version`

Return the experiment version.

Return type `str`

10.8.4 Neptune

```
class pytorch_lightning.loggers.neptune.NeptuneLogger (api_key=None,
                                                       project_name=None,
                                                       close_after_fit=True,      of-
                                                       offline_mode=False,      ex-
                                                       experiment_name=None,
                                                       **kwargs)
```

Bases: `pytorch_lightning.loggers.base.LightningLoggerBase`

Log using [Neptune](#). Install it with `pip`:

```
pip install neptune-client
```

The Neptune logger can be used in the online mode or offline (silent) mode. To log experiment data in online mode, `NeptuneLogger` requires an API key. In offline mode, the logger does not connect to Neptune.

ONLINE MODE

```
from pytorch_lightning import Trainer
from pytorch_lightning.loggers import NeptuneLogger

# arguments made to NeptuneLogger are passed on to the neptune.experiments.
↪Experiment class
# We are using an api_key for the anonymous user "neptuner" but you can use your_
↪own.
neptune_logger = NeptuneLogger(
    api_key='ANONYMOUS',
    project_name='shared/pytorch-lightning-integration',
    experiment_name='default', # Optional,
    params={'max_epochs': 10}, # Optional,
    tags=['pytorch-lightning', 'mlp'] # Optional,
)
trainer = Trainer(max_epochs=10, logger=neptune_logger)
```

OFFLINE MODE

```
from pytorch_lightning.loggers import NeptuneLogger

# arguments made to NeptuneLogger are passed on to the neptune.experiments.
↪Experiment class
neptune_logger = NeptuneLogger(
    offline_mode=True,
```

(continues on next page)

(continued from previous page)

```

project_name='USER_NAME/PROJECT_NAME',
experiment_name='default', # Optional,
params={'max_epochs': 10}, # Optional,
tags=['pytorch-lightning', 'mlp'] # Optional,
)
trainer = Trainer(max_epochs=10, logger=neptune_logger)

```

Use the logger anywhere in you LightningModule as follows:

```

class LitModel(LightningModule):
    def training_step(self, batch, batch_idx):
        # log metrics
        self.logger.experiment.log_metric('acc_train', ...)
        # log images
        self.logger.experiment.log_image('worse_predictions', ...)
        # log model checkpoint
        self.logger.experiment.log_artifact('model_checkpoint.pt', ...)
        self.logger.experiment.whatever_neptune_supports(...)

    def any_lightning_module_function_or_hook(self):
        self.logger.experiment.log_metric('acc_train', ...)
        self.logger.experiment.log_image('worse_predictions', ...)
        self.logger.experiment.log_artifact('model_checkpoint.pt', ...)
        self.logger.experiment.whatever_neptune_supports(...)

```

If you want to log objects after the training is finished use `close_after_fit=False`:

```

neptune_logger = NeptuneLogger(
    ...
    close_after_fit=False,
    ...
)
trainer = Trainer(logger=neptune_logger)
trainer.fit()

# Log test metrics
trainer.test(model)

# Log additional metrics
from sklearn.metrics import accuracy_score

accuracy = accuracy_score(y_true, y_pred)
neptune_logger.experiment.log_metric('test_accuracy', accuracy)

# Log charts
from scikitplot.metrics import plot_confusion_matrix
import matplotlib.pyplot as plt

fig, ax = plt.subplots(figsize=(16, 12))
plot_confusion_matrix(y_true, y_pred, ax=ax)
neptune_logger.experiment.log_image('confusion_matrix', fig)

# Save checkpoints folder
neptune_logger.experiment.log_artifact('my/checkpoints')

# When you are done, stop the experiment
neptune_logger.experiment.stop()

```

See also:

- An [Example experiment](#) showing the UI of Neptune.
- [Tutorial](#) on how to use Pytorch Lightning with Neptune.

Parameters

- **api_key** (Optional[str]) – Required in online mode. Neptune API token, found on <https://neptune.ai>. Read how to get your [API key](#). It is recommended to keep it in the `NEPTUNE_API_TOKEN` environment variable and then you can leave `api_key=None`.
- **project_name** (Optional[str]) – Required in online mode. Qualified name of a project in a form of “namespace/project_name” for example “tom/minst-classification”. If `None`, the value of `NEPTUNE_PROJECT` environment variable will be taken. You need to create the project in <https://neptune.ai> first.
- **offline_mode** (bool) – Optional default `False`. If `True` no logs will be sent to Neptune. Usually used for debug purposes.
- **close_after_fit** (Optional[bool]) – Optional default `True`. If `False` the experiment will not be closed after training and additional metrics, images or artifacts can be logged. Also, remember to close the experiment explicitly by running `neptune_logger.experiment.stop()`.
- **experiment_name** (Optional[str]) – Optional. Editable name of the experiment. Name is displayed in the experiment’s Details (Metadata section) and in experiments view as a column.
- ****kwargs** – Additional arguments like *params*, *tags*, *properties*, etc. used by `neptune.Session.create_experiment()` can be passed as keyword arguments in this logger.

append_tags (tags)

Appends tags to the neptune experiment.

Parameters **tags** (Union[str, Iterable[str]]) – Tags to add to the current experiment. If str is passed, a single tag is added. If multiple - comma separated - str are passed, all of them are added as tags. If list of str is passed, all elements of the list are added as tags.

Return type `None`

finalize (status)

Do any processing that is necessary to finalize an experiment.

Parameters **status** (str) – Status that the experiment finished with (e.g. success, failed, aborted)

Return type `None`

log_artifact (artifact, destination=None)

Save an artifact (file) in Neptune experiment storage.

Parameters

- **artifact** (str) – A path to the file in local filesystem.
- **destination** (Optional[str]) – Optional. Default is `None`. A destination path. If `None` is passed, an artifact file name will be used.

Return type `None`

log_hyperparams (*params*)

Record hyperparameters.

Parameters **params** `(Union[Dict[str, Any], Namespace])` – `Namespace` containing the hyperparameters

Return type `None`

log_image (*log_name*, *image*, *step=None*)

Log image data in Neptune experiment

Parameters

- **log_name** `(str)` – The name of log, i.e. bboxes, visualisations, sample_images.
- **image** `(Union[str, Any])` – The value of the log (data-point). Can be one of the following types: PIL image, `matplotlib.figure.Figure`, path to image file (`str`)
- **step** `(Optional[int])` – Step number at which the metrics should be recorded, must be strictly increasing

Return type `None`

log_metric (*metric_name*, *metric_value*, *step=None*)

Log metrics (numeric values) in Neptune experiments.

Parameters

- **metric_name** `(str)` – The name of log, i.e. mse, loss, accuracy.
- **metric_value** `(Union[Tensor, float, str])` – The value of the log (data-point).
- **step** `(Optional[int])` – Step number at which the metrics should be recorded, must be strictly increasing

Return type `None`

log_metrics (*metrics*, *step=None*)

Log metrics (numeric values) in Neptune experiments.

Parameters

- **metrics** `(Dict[str, Union[Tensor, float]])` – Dictionary with metric names as keys and measured quantities as values
- **step** `(Optional[int])` – Step number at which the metrics should be recorded, must be strictly increasing

Return type `None`

log_text (*log_name*, *text*, *step=None*)

Log text data in Neptune experiments.

Parameters

- **log_name** `(str)` – The name of log, i.e. mse, my_text_data, timing_info.
- **text** `(str)` – The value of the log (data-point).
- **step** `(Optional[int])` – Step number at which the metrics should be recorded, must be strictly increasing

Return type `None`

set_property (*key*, *value*)

Set key-value pair as Neptune experiment property.

Parameters

- **key** `(str)` – Property key.
- **value** `(Any)` – New value of a property.

Return type `None`**property experiment**Actual Neptune object. To use neptune features in your `LightningModule` do the following.

Example:

```
self.logger.experiment.some_neptune_function()
```

Return type `Experiment`**property name**

Return the experiment name.

Return type `str`**property save_dir**Return the root directory where experiment logs get saved, or `None` if the logger does not save data locally.**Return type** `Optional[str]`**property version**

Return the experiment version.

Return type `str`

10.8.5 Tensorboard

```
class pytorch_lightning.loggers.tensorboard.TensorBoardLogger(save_dir,
                                                               name='default',
                                                               version=None,
                                                               log_graph=False,
                                                               default_hp_metric=True,
                                                               **kwargs)
```

Bases: `pytorch_lightning.loggers.base.LightningLoggerBase`

Log to local file system in `TensorBoard` format. Implemented using `SummaryWriter`. Logs are saved to `os.path.join(save_dir, name, version)`. This is the default logger in Lightning, it comes pre-installed.

Example

```
>>> from pytorch_lightning import Trainer
>>> from pytorch_lightning.loggers import TensorBoardLogger
>>> logger = TensorBoardLogger("tb_logs", name="my_model")
>>> trainer = Trainer(logger=logger)
```

Parameters

- **save_dir** `(str)` – Save directory

- **name** (Optional[str]) – Experiment name. Defaults to 'default'. If it is the empty string then no per-experiment subdirectory is used.
- **version** (Union[int, str, None]) – Experiment version. If version is not specified the logger inspects the save directory for existing versions, then automatically assigns the next available version. If it is a string then it is used as the run-specific subdirectory name, otherwise 'version_\${version}' is used.
- **log_graph** (bool) – Adds the computational graph to tensorboard. This requires that the user has defined the *self.example_input_array* attribute in their model.
- **default_hp_metric** (bool) – Enables a placeholder metric with key *hp_metric* when *log_hyperparams* is called without a metric (otherwise calls to *log_hyperparams* without a metric are ignored).
- ****kwargs** – Additional arguments like *comment*, *filename_suffix*, etc. used by *SummaryWriter* can be passed as keyword arguments in this logger.

finalize (status)

Do any processing that is necessary to finalize an experiment.

Parameters **status** (str) – Status that the experiment finished with (e.g. success, failed, aborted)

Return type None

log_graph (model, input_array=None)

Record model graph

Parameters

- **model** (LightningModule) – lightning model
- **input_array** – input passes to *model.forward*

log_hyperparams (params, metrics=None)

Record hyperparameters.

Parameters **params** (Union[Dict[str, Any], Namespace]) – Namespace containing the hyperparameters

Return type None

log_metrics (metrics, step=None)

Records metrics. This method logs metrics as soon as it received them. If you want to aggregate metrics for one specific *step*, use the *agg_and_log_metrics()* method.

Parameters

- **metrics** (Dict[str, float]) – Dictionary with metric names as keys and measured quantities as values
- **step** (Optional[int]) – Step number at which the metrics should be recorded

Return type None

save ()

Save log data.

Return type None

property experiment

Actual tensorboard object. To use TensorBoard features in your *LightningModule* do the following.

Example:


```
self.logger.experiment.some_tensorboard_function()
```

Return type `SummaryWriter`

property log_dir

The directory for this run's tensorboard checkpoint. By default, it is named 'version_\${self.version}' but it can be overridden by passing a string value for the constructor's version parameter instead of `None` or an `int`.

Return type `str`

property name

Return the experiment name.

Return type `str`

property root_dir

Parent directory for all tensorboard checkpoint subdirectories. If the experiment name parameter is `None` or the empty string, no experiment subdirectory is used and the checkpoint will be saved in "save_dir/version_dir"

Return type `str`

property save_dir

Return the root directory where experiment logs get saved, or `None` if the logger does not save data locally.

Return type `Optional[str]`

property version

Return the experiment version.

Return type `int`

10.8.6 Test-tube

```
class pytorch_lightning.loggers.test_tube.TestTubeLogger (save_dir, name='default',
                                                         description=None,
                                                         debug=False,      ver-
                                                         sion=None,          cre-
                                                         ate_git_tag=False,
                                                         log_graph=False)
```

Bases: `pytorch_lightning.loggers.base.LightningLoggerBase`

Log to local file system in [TensorBoard](#) format but using a nicer folder structure (see [full docs](#)). Install it with `pip`:

```
pip install test_tube
```

```
from pytorch_lightning import Trainer
from pytorch_lightning.loggers import TestTubeLogger
logger = TestTubeLogger("tt_logs", name="my_exp_name")
trainer = Trainer(logger=logger)
```

Use the logger anywhere in your `LightningModule` as follows:

```
from pytorch_lightning import LightningModule
class LitModel(LightningModule):
    def training_step(self, batch, batch_idx):
```

(continues on next page)

(continued from previous page)

```
# example
self.logger.experiment.whatever_method_summary_writer_supports(...)

def any_lightning_module_function_or_hook(self):
    self.logger.experiment.add_histogram(...)
```

Parameters

- **save_dir** (str) – Save directory
- **name** (str) – Experiment name. Defaults to 'default'.
- **description** (Optional[str]) – A short snippet about this experiment
- **debug** (bool) – If True, it doesn't log anything.
- **version** (Optional[int]) – Experiment version. If version is not specified the logger inspects the save directory for existing versions, then automatically assigns the next available version.
- **create_git_tag** (bool) – If True creates a git tag to save the code used in this experiment.
- **log_graph** (bool) – Adds the computational graph to tensorboard. This requires that the user has defined the *self.example_input_array* attribute in their model.

close()

Do any cleanup that is necessary to close an experiment.

Return type None**finalize(status)**

Do any processing that is necessary to finalize an experiment.

Parameters **status** (str) – Status that the experiment finished with (e.g. success, failed, aborted)**Return type** None**log_graph(model, input_array=None)**

Record model graph

Parameters

- **model** (LightningModule) – lightning model
- **input_array** – input passes to *model.forward*

log_hyperparams(params)

Record hyperparameters.

Parameters **params** (Union[Dict[str, Any], Namespace]) – Namespace containing the hyperparameters**Return type** None**log_metrics(metrics, step=None)**Records metrics. This method logs metrics as as soon as it received them. If you want to aggregate metrics for one specific *step*, use the *agg_and_log_metrics()* method.**Parameters**

- **metrics** (Dict[str, float]) – Dictionary with metric names as keys and measured quantities as values

- **step** (Optional[int]) – Step number at which the metrics should be recorded

Return type None

save()

Save log data.

Return type None

property experiment

Actual TestTube object. To use TestTube features in your LightningModule do the following.

Example:

```
self.logger.experiment.some_test_tube_function()
```

Return type Experiment

property name

Return the experiment name.

Return type str

property save_dir

Return the root directory where experiment logs get saved, or None if the logger does not save data locally.

Return type Optional[str]

property version

Return the experiment version.

Return type int

10.8.7 Weights and Biases

```
class pytorch_lightning.loggers.wandb.WandbLogger (name=None, save_dir=None,
                                                    offline=False, id=None, anonymous=False, version=None,
                                                    project=None, log_model=False,
                                                    experiment=None, **kwargs)
```

Bases: pytorch_lightning.loggers.base.LightningLoggerBase

Log using [Weights and Biases](#). Install it with pip:

```
pip install wandb
```

Parameters

- **name** (Optional[str]) – Display name for the run.
- **save_dir** (Optional[str]) – Path where data is saved.
- **offline** (bool) – Run offline (data can be streamed later to wandb servers).
- **id** (Optional[str]) – Sets the version, mainly used to resume a previous run.
- **anonymous** (bool) – Enables or explicitly disables anonymous logging.
- **version** (Optional[str]) – Sets the version, mainly used to resume a previous run.
- **project** (Optional[str]) – The name of the project to which this run will belong.
- **log_model** (bool) – Save checkpoints in wandb dir to upload on W&B servers.

- **experiment** – WandB experiment object.
- ****kwargs** – Additional arguments like *entity*, *group*, *tags*, etc. used by `wandb.init()` can be passed as keyword arguments in this logger.

Example

```
from pytorch_lightning.loggers import WandbLogger from pytorch_lightning import Trainer wandb_logger = WandbLogger() trainer = Trainer(logger=wandb_logger)
```

See also:

- [Tutorial](#) on how to use W&B with Pytorch Lightning.

log_hyperparams (*params*)
Record hyperparameters.

Parameters *params* (Union[Dict[str, Any], Namespace]) – Namespace containing the hyperparameters

Return type None

log_metrics (*metrics*, *step=None*)

Records metrics. This method logs metrics as soon as it received them. If you want to aggregate metrics for one specific *step*, use the `agg_and_log_metrics()` method.

Parameters

- **metrics** (Dict[str, float]) – Dictionary with metric names as keys and measured quantities as values
- **step** (Optional[int]) – Step number at which the metrics should be recorded

Return type None

property experiment

Actual wandb object. To use wandb features in your LightningModule do the following.

Example:

```
self.logger.experiment.some_wandb_function()
```

Return type Run

property name

Return the experiment name.

Return type Optional[str]

property save_dir

Return the root directory where experiment logs get saved, or *None* if the logger does not save data locally.

Return type Optional[str]

property version

Return the experiment version.

Return type Optional[str]

METRICS

`pytorch_lightning.metrics` is a Metrics API created for easy metric development and usage in PyTorch and PyTorch Lightning. It is rigorously tested for all edge cases and includes a growing list of common metric implementations.

The metrics API provides `update()`, `compute()`, `reset()` functions to the user. The metric base class inherits `nn.Module` which allows us to call `metric(...)` directly. The `forward()` method of the base `Metric` class serves the dual purpose of calling `update()` on its input and simultaneously returning the value of the metric over the provided input.

These metrics work with DDP in PyTorch and PyTorch Lightning by default. When `.compute()` is called in distributed mode, the internal state of each metric is synced and reduced across each process, so that the logic present in `.compute()` is applied to state information from all processes.

The example below shows how to use a metric in your `LightningModule`:

```
def __init__(self):
    ...
    self.accuracy = pl.metrics.Accuracy()

def training_step(self, batch, batch_idx):
    logits = self(x)
    ...
    # log step metric
    self.log('train_acc_step', self.accuracy(logits, y))
    ...

def training_epoch_end(self, outs):
    # log epoch metric
    self.log('train_acc_epoch', self.accuracy.compute())
```

Metric objects can also be directly logged, in which case Lightning will log the metric based on `on_step` and `on_epoch` flags present in `self.log(...)`. If `on_epoch` is `True`, the logger automatically logs the end of epoch metric value by calling `.compute()`.

Note: `sync_dist`, `sync_dist_op`, `sync_dist_group`, `reduce_fx` and `tbptt_reduce_fx` flags from `self.log(...)` don't affect the metric logging in any manner. The metric class contains its own distributed synchronization logic.

This however is only true for metrics that inherit the base class `Metric`, and thus the functional metric API provides no support for in-built distributed synchronization or reduction functions.

```
def __init__(self):
    ...
    self.train_acc = pl.metrics.Accuracy()
    self.valid_acc = pl.metrics.Accuracy()

def training_step(self, batch, batch_idx):
    logits = self(x)
    ...
    self.train_acc(logits, y)
    self.log('train_acc', self.train_acc, on_step=True, on_epoch=False)

def validation_step(self, batch, batch_idx):
    logits = self(x)
    ...
    self.valid_acc(logits, y)
    self.log('valid_acc', self.valid_acc, on_step=True, on_epoch=True)
```

This metrics API is independent of PyTorch Lightning. Metrics can directly be used in PyTorch as shown in the example:

```
from pytorch_lightning import metrics

train_accuracy = metrics.Accuracy()
valid_accuracy = metrics.Accuracy(compute_on_step=False)

for epoch in range(epochs):
    for x, y in train_data:
        y_hat = model(x)

        # training step accuracy
        batch_acc = train_accuracy(y_hat, y)

    for x, y in valid_data:
        y_hat = model(x)
        valid_accuracy(y_hat, y)

# total accuracy over all training batches
total_train_accuracy = train_accuracy.compute()

# total accuracy over all validation batches
total_valid_accuracy = valid_accuracy.compute()
```

11.1 Implementing a Metric

To implement your custom metric, subclass the base `Metric` class and implement the following methods:

- `__init__()`: Each state variable should be called using `self.add_state(...)`.
- `update()`: Any code needed to update the state given any inputs to the metric.
- `compute()`: Computes a final value from the state of the metric.

All you need to do is call `add_state` correctly to implement a custom metric with DDP. `reset()` is called on metric state variables added using `add_state()`.

To see how metric states are synchronized across distributed processes, refer to `add_state()` docs from the base `Metric` class.

Example implementation:

```
from pytorch_lightning.metrics import Metric

class MyAccuracy(Metric):
    def __init__(self, dist_sync_on_step=False):
        super().__init__(dist_sync_on_step=dist_sync_on_step)

        self.add_state("correct", default=torch.tensor(0), dist_reduce_fx="sum")
        self.add_state("total", default=torch.tensor(0), dist_reduce_fx="sum")

    def update(self, preds: torch.Tensor, target: torch.Tensor):
        preds, target = self._input_format(preds, target)
        assert preds.shape == target.shape

        self.correct += torch.sum(preds == target)
        self.total += target.numel()

    def compute(self):
        return self.correct.float() / self.total
```

11.2 Metric API

class `pytorch_lightning.metrics.Metric` (*compute_on_step=True*, *dist_sync_on_step=False*, *process_group=None*)

Bases: `torch.nn.Module`, `abc.ABC`

Base class for all metrics present in the Metrics API.

Implements `add_state()`, `forward()`, `reset()` and a few other things to handle distributed synchronization and per-step metric computation.

Override `update()` and `compute()` functions to implement your own metric. Use `add_state()` to register metric state variables which keep track of state on each call of `update()` and are synchronized across processes when `compute()` is called.

Note: Metric state variables can either be `torch.Tensors` or an empty list which can be used to store `torch.Tensors`.

Note: Different metrics only override `update()` and not `forward()`. A call to `update()` is valid, but it won't return the metric value at the current step. A call to `forward()` automatically calls `update()` and also returns the metric value at the current step.

Parameters

- **`compute_on_step`** (`bool`) – Forward only calls `update()` and returns `None` if this is set to `False`. default: `True`
- **`dist_sync_on_step`** (`bool`) – Synchronize metric state across processes at each `forward()` before returning the value at the step. default: `False`
- **`process_group`** (`Optional[Any]`) – Specify the process group on which synchronization is called. default: `None` (which selects the entire world)

add_state (*name*, *default*, *dist_reduce_fx=None*)

Adds metric state variable. Only used by subclasses.

Parameters

- **name** `⚡` (*str*) – The name of the state variable. The variable will then be accessible at `self.name`.
- **default** `⚡` – Default value of the state; can either be a `torch.Tensor` or an empty list. The state will be reset to this value when `self.reset()` is called.
- **dist_reduce_fx** `⚡` (*Optional*) – Function to reduce state across multiple processes in distributed mode. If value is "sum", "mean", or "cat", we will use `torch.sum`, `torch.mean`, and `torch.cat` respectively, each with argument `dim=0`. The user can also pass a custom function in this parameter.

Note: Setting `dist_reduce_fx` to `None` will return the metric state synchronized across different processes. However, there won't be any reduction function applied to the synchronized metric state.

The metric states would be synced as follows

- If the metric state is `torch.Tensor`, the synced value will be a stacked `torch.Tensor` across the process dimension if the metric state was a `torch.Tensor`. The original `torch.Tensor` metric state retains dimension and hence the synchronized output will be of shape `(num_process, ...)`.
- If the metric state is a `list`, the synced value will be a `list` containing the combined elements from all processes.

Note: When passing a custom function to `dist_reduce_fx`, expect the synchronized metric state to follow the format discussed in the above note.

abstract compute ()

Override this method to compute the final metric value from state variables synchronized across the distributed backend.

forward (**args*, ***kwargs*)

Automatically calls `update()`. Returns the metric value over inputs if `compute_on_step` is `True`.

reset ()

This method automatically resets the metric state variables to their default value.

abstract update ()

Override this method to update the state variables of your metric class.

Return type `None`

11.3 Class metrics

11.3.1 Classification Metrics

Accuracy

```
class pytorch_lightning.metrics.classification.Accuracy (threshold=0.5, compute_on_step=True,  
dist_sync_on_step=False,  
process_group=None)
```

Bases: `pytorch_lightning.metrics.metric.Metric`

Computes accuracy. Works with binary, multiclass, and multilabel data. Accepts logits from a model output or integer class values in prediction. Works with multi-dimensional preds and target.

Forward accepts

- `preds` (float or long tensor): (N, \dots) or (N, C, \dots) where C is the number of classes
- `target` (long tensor): (N, \dots)

If `preds` and `target` are the same shape and `preds` is a float tensor, we use the `self.threshold` argument. This is the case for binary and multi-label logits.

If `preds` has an extra dimension as in the case of multi-class scores we perform an `argmax` on `dim=1`.

Parameters

- **`threshold`** (float) – Threshold value for binary or multi-label logits. default: 0.5
- **`compute_on_step`** (bool) – Forward only calls `update()` and return `None` if this is set to `False`. default: `True`
- **`dist_sync_on_step`** (bool) – Synchronize metric state across processes at each `forward()` before returning the value at the step. default: `False`
- **`process_group`** (Optional[Any]) – Specify the process group on which synchronization is called. default: `None` (which selects the entire world)

Example

```
>>> from pytorch_lightning.metrics import Accuracy
>>> target = torch.tensor([0, 1, 2, 3])
>>> preds = torch.tensor([0, 2, 1, 3])
>>> accuracy = Accuracy()
>>> accuracy(preds, target)
tensor(0.5000)
```

`compute()`

Computes accuracy over state.

`update(preds, target)`

Update state with predictions and targets.

Parameters

- **`preds`** (Tensor) – Predictions from model
- **`target`** (Tensor) – Ground truth values

Precision

```
class pytorch_lightning.metrics.classification.Precision (num_classes=1,
                                                         threshold=0.5,      av-
                                                         erage='micro',      mul-
                                                         tilabel=False,      com-
                                                         pute_on_step=True,
                                                         dist_sync_on_step=False,
                                                         process_group=None)
```

Bases: `pytorch_lightning.metrics.metric.Metric`

Computes the precision metric.

Works with binary, multiclass, and multilabel data. Accepts logits from a model output or integer class values in prediction. Works with multi-dimensional preds and target.

Forward accepts

- `preds` (float or long tensor): (N, \dots) or (N, C, \dots) where C is the number of classes
- `target` (long tensor): (N, \dots)

If `preds` and `target` are the same shape and `preds` is a float tensor, we use the `self.threshold` argument. This is the case for binary and multi-label logits.

If `preds` has an extra dimension as in the case of multi-class scores we perform an `argmax` on `dim=1`.

Parameters

- **`num_classes`** (int) – Number of classes in the dataset.
- **`beta`** – Beta coefficient in the F measure.
- **`threshold`** (float) – Threshold value for binary or multi-label logits. default: 0.5
- **`average`** (str) –
 - ‘micro’ computes metric globally
 - ‘macro’ computes metric for each class and then takes the mean
- **`multilabel`** (bool) – If predictions are from multilabel classification.
- **`compute_on_step`** (bool) – Forward only calls `update()` and return None if this is set to False. default: True
- **`dist_sync_on_step`** (bool) – Synchronize metric state across processes at each `forward()` before returning the value at the step. default: False
- **`process_group`** (Optional[Any]) – Specify the process group on which synchronization is called. default: None (which selects the entire world)

Example

```
>>> from pytorch_lightning.metrics import Precision
>>> target = torch.tensor([0, 1, 2, 0, 1, 2])
>>> preds = torch.tensor([0, 2, 1, 0, 0, 1])
>>> precision = Precision(num_classes=3)
>>> precision(preds, target)
tensor(0.3333)
```

compute()

Override this method to compute the final metric value from state variables synchronized across the distributed backend.

update(preds, target)

Override this method to update the state variables of your metric class.

Recall

```
class pytorch_lightning.metrics.classification.Recall(num_classes=1,      thresh-
                                                    old=0.5,      average='micro',
                                                    multilabel=False,      com-
                                                    pute_on_step=True,
                                                    dist_sync_on_step=False,
                                                    process_group=None)
```

Bases: `pytorch_lightning.metrics.metric.Metric`

Computes the recall metric.

Works with binary, multiclass, and multilabel data. Accepts logits from a model output or integer class values in prediction. Works with multi-dimensional preds and target.

Forward accepts

- `preds` (float or long tensor): (N, ...) or (N, C, ...) where C is the number of classes
- `target` (long tensor): (N, ...)

If `preds` and `target` are the same shape and `preds` is a float tensor, we use the `self.threshold` argument. This is the case for binary and multi-label logits.

If `preds` has an extra dimension as in the case of multi-class scores we perform an `argmax` on `dim=1`.

Parameters

- **num_classes** (int) – Number of classes in the dataset.
- **beta** – Beta coefficient in the F measure.
- **threshold** (float) – Threshold value for binary or multi-label logits. default: 0.5
- **average** (str) –
 - ‘micro’ computes metric globally
 - ‘macro’ computes metric for each class and then takes the mean
- **multilabel** (bool) – If predictions are from multilabel classification.
- **compute_on_step** (bool) – Forward only calls `update()` and return None if this is set to False. default: True
- **dist_sync_on_step** (bool) – Synchronize metric state across processes at each `forward()` before returning the value at the step. default: False
- **process_group** (Optional[Any]) – Specify the process group on which synchronization is called. default: None (which selects the entire world)

Example

```
>>> from pytorch_lightning.metrics import Recall
>>> target = torch.tensor([0, 1, 2, 0, 1, 2])
>>> preds = torch.tensor([0, 2, 1, 0, 0, 1])
>>> recall = Recall(num_classes=3)
>>> recall(preds, target)
tensor(0.3333)
```

compute()

Computes accuracy over state.

update(preds, target)

Update state with predictions and targets.

Parameters

- **preds** `(Tensor)` – Predictions from model
- **target** `(Tensor)` – Ground truth values

Fbeta

```
class pytorch_lightning.metrics.classification.Fbeta(num_classes=1,      beta=1.0,
                                                    threshold=0.5,      aver-
                                                    age='micro', multilabel=False,
                                                    compute_on_step=True,
                                                    dist_sync_on_step=False,
                                                    process_group=None)
```

Bases: `pytorch_lightning.metrics.metric.Metric`

Computes f_beta metric.

Works with binary, multiclass, and multilabel data. Accepts logits from a model output or integer class values in prediction. Works with multi-dimensional preds and target.

Forward accepts

- **preds** (float or long tensor): (N, ...) or (N, C, ...) where C is the number of classes
- **target** (long tensor): (N, ...)

If preds and target are the same shape and preds is a float tensor, we use the `self.threshold` argument. This is the case for binary and multi-label logits.

If preds has an extra dimension as in the case of multi-class scores we perform an `argmax` on `dim=1`.

Parameters

- **num_classes** `(int)` – Number of classes in the dataset.
- **beta** `(float)` – Beta coefficient in the F measure.
- **threshold** `(float)` – Threshold value for binary or multi-label logits. default: 0.5
- **average** `(str)` –
 - ‘micro’ computes metric globally
 - ‘macro’ computes metric for each class and then takes the mean
- **multilabel** `(bool)` – If predictions are from multilabel classification.

- **compute_on_step** (bool) – Forward only calls `update()` and return `None` if this is set to `False`. default: `True`
- **dist_sync_on_step** (bool) – Synchronize metric state across processes at each `forward()` before returning the value at the step. default: `False`
- **process_group** (Optional[Any]) – Specify the process group on which synchronization is called. default: `None` (which selects the entire world)

Example

```
>>> from pytorch_lightning.metrics import Fbeta
>>> target = torch.tensor([0, 1, 2, 0, 1, 2])
>>> preds = torch.tensor([0, 2, 1, 0, 0, 1])
>>> f_beta = Fbeta(num_classes=3, beta=0.5)
>>> f_beta(preds, target)
tensor(0.3333)
```

compute()

Computes accuracy over state.

update(preds, target)

Update state with predictions and targets.

Parameters

- **preds** (Tensor) – Predictions from model
- **target** (Tensor) – Ground truth values

11.3.2 Regression Metrics

MeanSquaredError

```
class pytorch_lightning.metrics.regression.MeanSquaredError (compute_on_step=True,
                                                             dist_sync_on_step=False,
                                                             process_group=None)
```

Bases: `pytorch_lightning.metrics.metric.Metric`

Computes mean squared error.

Parameters

- **compute_on_step** (bool) – Forward only calls `update()` and return `None` if this is set to `False`. default: `True`
- **dist_sync_on_step** (bool) – Synchronize metric state across processes at each `forward()` before returning the value at the step. default: `False`
- **process_group** (Optional[Any]) – Specify the process group on which synchronization is called. default: `None` (which selects the entire world)

Example

```
>>> from pytorch_lightning.metrics import MeanSquaredError
>>> target = torch.tensor([2.5, 5.0, 4.0, 8.0])
>>> preds = torch.tensor([3.0, 5.0, 2.5, 7.0])
>>> mean_squared_error = MeanSquaredError()
>>> mean_squared_error(preds, target)
tensor(0.8750)
```

compute()

Computes mean squared error over state.

update(preds, target)

Update state with predictions and targets.

Parameters

- **preds** *(Tensor)* – Predictions from model
- **target** *(Tensor)* – Ground truth values

MeanAbsoluteError

```
class pytorch_lightning.metrics.regression.MeanAbsoluteError(compute_on_step=True,  
                                                             dist_sync_on_step=False,  
                                                             pro-  
                                                             cess_group=None)
```

Bases: `pytorch_lightning.metrics.metric.Metric`

Computes mean absolute error.

Parameters

- **compute_on_step** *(bool)* – Forward only calls `update()` and return `None` if this is set to `False`. default: `True`
- **dist_sync_on_step** *(bool)* – Synchronize metric state across processes at each `forward()` before returning the value at the step. default: `False`
- **process_group** *(Optional[Any])* – Specify the process group on which synchronization is called. default: `None` (which selects the entire world)

Example

```
>>> from pytorch_lightning.metrics import MeanAbsoluteError
>>> target = torch.tensor([3.0, -0.5, 2.0, 7.0])
>>> preds = torch.tensor([2.5, 0.0, 2.0, 8.0])
>>> mean_absolute_error = MeanAbsoluteError()
>>> mean_absolute_error(preds, target)
tensor(0.5000)
```

compute()

Computes mean absolute error over state.

update(preds, target)

Update state with predictions and targets.

Parameters

- **preds** (Tensor) – Predictions from model
- **target** (Tensor) – Ground truth values

MeanSquaredLogError

```
class pytorch_lightning.metrics.regression.MeanSquaredLogError (compute_on_step=True,
                                                                dist_sync_on_step=False,
                                                                process_group=None)
```

Bases: `pytorch_lightning.metrics.metric.Metric`

Computes mean squared logarithmic error.

Parameters

- **compute_on_step** (bool) – Forward only calls `update()` and return `None` if this is set to `False`. default: `True`
- **dist_sync_on_step** (bool) – Synchronize metric state across processes at each `forward()` before returning the value at the step. default: `False`
- **process_group** (Optional[Any]) – Specify the process group on which synchronization is called. default: `None` (which selects the entire world)

Example

```
>>> from pytorch_lightning.metrics import MeanSquaredLogError
>>> target = torch.tensor([2.5, 5, 4, 8])
>>> preds = torch.tensor([3, 5, 2.5, 7])
>>> mean_squared_log_error = MeanSquaredLogError()
>>> mean_squared_log_error(preds, target)
tensor(0.0397)
```

compute()

Compute mean squared logarithmic error over state.

update(preds, target)

Update state with predictions and targets.

Parameters

- **preds** (Tensor) – Predictions from model
- **target** (Tensor) – Ground truth values

ExplainedVariance

```
class pytorch_lightning.metrics.regression.ExplainedVariance (multioutput='uniform_average',
                                                                compute_on_step=True,
                                                                dist_sync_on_step=False,
                                                                process_group=None)
```

Bases: `pytorch_lightning.metrics.metric.Metric`

Computes explained variance.

Forward accepts

- `preds` (float tensor): (N,) or (N, ...) (multioutput)
- `target` (long tensor): (N,) or (N, ...) (multioutput)

In the case of multioutput, as default the variances will be uniformly averaged over the additional dimensions. Please see argument *multioutput* for changing this behavior.

Parameters

- **`multioutput`** *(str)* – Defines aggregation in the case of multiple output scores. Can be one of the following strings (default is *'uniform_average'*):
 - *'raw_values'* returns full set of scores
 - *'uniform_average'* scores are uniformly averaged
 - *'variance_weighted'* scores are weighted by their individual variances
- **`compute_on_step`** *(bool)* – Forward only calls `update()` and return None if this is set to False. default: True
- **`dist_sync_on_step`** *(bool)* – Synchronize metric state across processes at each `forward()` before returning the value at the step. default: False
- **`process_group`** *(Optional[Any])* – Specify the process group on which synchronization is called. default: None (which selects the entire world)

Example

```
>>> from pytorch_lightning.metrics import ExplainedVariance
>>> target = torch.tensor([3, -0.5, 2, 7])
>>> preds = torch.tensor([2.5, 0.0, 2, 8])
>>> explained_variance = ExplainedVariance()
>>> explained_variance(preds, target)
tensor(0.9572)
```

```
>>> target = torch.tensor([[0.5, 1], [-1, 1], [7, -6]])
>>> preds = torch.tensor([[0, 2], [-1, 2], [8, -5]])
>>> explained_variance = ExplainedVariance(multioutput='raw_values')
>>> explained_variance(preds, target)
tensor([0.9677, 1.0000])
```

`compute()`

Computes explained variance over state.

`update(preds, target)`

Update state with predictions and targets.

Parameters

- **`preds`** *(Tensor)* – Predictions from model
- **`target`** *(Tensor)* – Ground truth values

11.4 Functional Metrics

The functional metrics follow the simple paradigm input in, output out. This means, they don't provide any advanced mechanisms for syncing across DDP nodes or aggregation over batches. They simply compute the metric value based on the given inputs.

Also the integration within other parts of PyTorch Lightning will never be as tight as with the class-based interface. If you look for just computing the values, the functional metrics are the way to go. However, if you are looking for the best integration and user experience, please consider also to use the class interface.

11.4.1 Classification

accuracy [func]

```
pytorch_lightning.metrics.functional.classification.accuracy(pred, target,
                                                             num_classes=None,
                                                             class_reduction='micro',
                                                             re-
                                                             turn_state=False)
```

Computes the accuracy classification score

Parameters

- **pred** (Tensor) – predicted labels
- **target** (Tensor) – ground truth labels
- **num_classes** (Optional[int]) – number of classes
- **class_reduction** (str) – method to reduce metric score over labels
 - 'micro': calculate metrics globally (default)
 - 'macro': calculate metrics for each label, and find their unweighted mean.
 - 'weighted': calculate metrics for each label, and find their weighted mean.
 - 'none': returns calculated metric per class
- **return_state** (bool) – returns a internal state that can be ddp reduced before doing the final calculation

Return type Tensor

Returns A Tensor with the accuracy score.

Example

```
>>> x = torch.tensor([0, 1, 2, 3])
>>> y = torch.tensor([0, 1, 2, 2])
>>> accuracy(x, y)
tensor(0.7500)
```

auc [func]

`pytorch_lightning.metrics.functional.classification.auc(x, y, reorder=True)`
Computes Area Under the Curve (AUC) using the trapezoidal rule

Parameters

- **x** `(Tensor)` – x-coordinates
- **y** `(Tensor)` – y-coordinates
- **reorder** `(bool)` – reorder coordinates, so they are increasing

Return type `Tensor`

Returns Tensor containing AUC score (float)

Example

```
>>> x = torch.tensor([0, 1, 2, 3])
>>> y = torch.tensor([0, 1, 2, 2])
>>> auc(x, y)
tensor(4.)
```

auroc [func]

`pytorch_lightning.metrics.functional.classification.auroc(pred, target, sample_weight=None, pos_label=1.0)`
Compute Area Under the Receiver Operating Characteristic Curve (ROC AUC) from prediction scores

Parameters

- **pred** `(Tensor)` – estimated probabilities
- **target** `(Tensor)` – ground-truth labels
- **sample_weight** `(Optional[Sequence])` – sample weights
- **pos_label** `(int)` – the label for the positive class

Return type `Tensor`

Returns Tensor containing ROCAUC score

Example

```
>>> x = torch.tensor([0, 1, 2, 3])
>>> y = torch.tensor([0, 1, 1, 0])
>>> auroc(x, y)
tensor(0.5000)
```

average_precision [func]

```
pytorch_lightning.metrics.functional.classification.average_precision(pred,
                                                                    tar-
                                                                    get,
                                                                    sam-
                                                                    ple_weight=None,
                                                                    pos_label=1.0)
```

Compute average precision from prediction scores

Parameters

- **pred** (Tensor) – estimated probabilities
- **target** (Tensor) – ground-truth labels
- **sample_weight** (Optional[Sequence]) – sample weights
- **pos_label** (int) – the label for the positive class

Return type Tensor

Returns Tensor containing average precision score

Example

```
>>> x = torch.tensor([0, 1, 2, 3])
>>> y = torch.tensor([0, 1, 2, 2])
>>> average_precision(x, y)
tensor(0.3333)
```

confusion_matrix [func]

```
pytorch_lightning.metrics.functional.classification.confusion_matrix(pred,
                                                                    target,
                                                                    nor-
                                                                    mal-
                                                                    ize=False,
                                                                    num_classes=None)
```

Computes the confusion matrix C where each entry $C_{\{i,j\}}$ is the number of observations in group i that were predicted in group j.

Parameters

- **pred** (Tensor) – estimated targets
- **target** (Tensor) – ground truth labels
- **normalize** (bool) – normalizes confusion matrix
- **num_classes** (Optional[int]) – number of classes

Return type Tensor

Returns Tensor, confusion matrix C [num_classes, num_classes]

Example

```
>>> x = torch.tensor([1, 2, 3])
>>> y = torch.tensor([0, 2, 3])
>>> confusion_matrix(x, y)
tensor([[0., 1., 0., 0.],
        [0., 0., 0., 0.],
        [0., 0., 1., 0.],
        [0., 0., 0., 1.]])
```

dice_score [func]

```
pytorch_lightning.metrics.functional.classification.dice_score(pred, target,
                                                                bg=False,
                                                                nan_score=0.0,
                                                                no_fg_score=0.0,
                                                                reduction='elementwise_mean')
```

Compute dice score from prediction scores

Parameters

- **pred** `(Tensor)` – estimated probabilities
- **target** `(Tensor)` – ground-truth labels
- **bg** `(bool)` – whether to also compute dice for the background
- **nan_score** `(float)` – score to return, if a NaN occurs during computation
- **no_fg_score** `(float)` – score to return, if no foreground pixel was found in target
- **reduction** `(str)` – a method to reduce metric score over labels.
 - 'elementwise_mean': takes the mean (default)
 - 'sum': takes the sum
 - 'none': no reduction will be applied

Return type `Tensor`

Returns Tensor containing dice score

Example

```
>>> pred = torch.tensor([[0.85, 0.05, 0.05, 0.05],
...                      [0.05, 0.85, 0.05, 0.05],
...                      [0.05, 0.05, 0.85, 0.05],
...                      [0.05, 0.05, 0.05, 0.85]])
>>> target = torch.tensor([0, 1, 3, 2])
>>> dice_score(pred, target)
tensor(0.3333)
```

f1_score [func]

```
pytorch_lightning.metrics.functional.classification.f1_score(pred, target,
                                                            num_classes=None,
                                                            class_reduction='micro')
```

Computes the F1-score (a.k.a F-measure), which is the harmonic mean of the precision and recall. It ranges between 1 and 0, where 1 is perfect and the worst value is 0.

Parameters

- **pred** (Tensor) – estimated probabilities
- **target** (Tensor) – ground-truth labels
- **num_classes** (Optional[int]) – number of classes
- **class_reduction** (str) – method to reduce metric score over labels
 - 'micro': calculate metrics globally (default)
 - 'macro': calculate metrics for each label, and find their unweighted mean.
 - 'weighted': calculate metrics for each label, and find their weighted mean.
 - 'none': returns calculated metric per class

Return type Tensor

Returns Tensor containing F1-score

Example

```
>>> x = torch.tensor([0, 1, 2, 3])
>>> y = torch.tensor([0, 1, 2, 2])
>>> f1_score(x, y)
tensor(0.7500)
```

fbeta_score [func]

```
pytorch_lightning.metrics.functional.classification.fbeta_score(pred, target, beta,
                                                                num_classes=None,
                                                                class_reduction='micro')
```

Computes the F-beta score which is a weighted harmonic mean of precision and recall. It ranges between 1 and 0, where 1 is perfect and the worst value is 0.

Parameters

- **pred** (Tensor) – estimated probabilities
- **target** (Tensor) – ground-truth labels
- **beta** (float) – weights recall when combining the score. beta < 1: more weight to precision. beta > 1 more weight to recall beta = 0: only precision beta -> inf: only recall
- **num_classes** (Optional[int]) – number of classes
- **class_reduction** (str) – method to reduce metric score over labels
 - 'micro': calculate metrics globally (default)
 - 'macro': calculate metrics for each label, and find their unweighted mean.

- 'weighted': calculate metrics for each label, and find their weighted mean.
- 'none': returns calculated metric per class

Return type `Tensor`

Returns Tensor with the value of F-score. It is a value between 0-1.

Example

```
>>> x = torch.tensor([0, 1, 2, 3])
>>> y = torch.tensor([0, 1, 2, 2])
>>> fbeta_score(x, y, 0.2)
tensor(0.7500)
```

iou [func]

`pytorch_lightning.metrics.functional.classification.iou(pred, target, ignore_index=None, absent_score=0.0, num_classes=None, reduction='elementwise_mean')`

Intersection over union, or Jaccard index calculation.

Parameters

- **pred** `(Tensor)` – Tensor containing predictions
- **target** `(Tensor)` – Tensor containing targets
- **ignore_index** `(Optional[int])` – optional int specifying a target class to ignore. If given, this class index does not contribute to the returned score, regardless of reduction method. Has no effect if given an int that is not in the range `[0, num_classes-1]`, where `num_classes` is either given or derived from `pred` and `target`. By default, no index is ignored, and all classes are used.
- **absent_score** `(float)` – score to use for an individual class, if no instances of the class index were present in `pred` AND no instances of the class index were present in `target`. For example, if we have 3 classes, `[0, 0]` for `pred`, and `[0, 2]` for `target`, then class 1 would be assigned the `absent_score`. Default is 0.0.
- **num_classes** `(Optional[int])` – Optionally specify the number of classes
- **reduction** `(str)` – a method to reduce metric score over labels.
 - 'elementwise_mean': takes the mean (default)
 - 'sum': takes the sum
 - 'none': no reduction will be applied

Returns Tensor containing single value if reduction is 'elementwise_mean', or number of classes if reduction is 'none'

Return type IoU score

Example

```
>>> target = torch.randint(0, 1, (10, 25, 25))
>>> pred = torch.tensor(target)
>>> pred[2:5, 7:13, 9:15] = 1 - pred[2:5, 7:13, 9:15]
>>> iou(pred, target)
tensor(0.4914)
```

multiclass_roc [func]

```
pytorch_lightning.metrics.functional.classification.multiclass_roc(pred, target, sample_weight=None, num_classes=None)
```

Computes the Receiver Operating Characteristic (ROC) for multiclass predictors.

Parameters

- **pred** (Tensor) – estimated probabilities
- **target** (Tensor) – ground-truth labels
- **sample_weight** (Optional[Sequence]) – sample weights
- **num_classes** (Optional[int]) – number of classes (default: None, computes automatically from data)

Return type Tuple[Tuple[Tensor, Tensor, Tensor]]

Returns returns roc for each class. Number of classes, false-positive rate (fpr), true-positive rate (tpr), thresholds

Example

```
>>> pred = torch.tensor([[0.85, 0.05, 0.05, 0.05],
...                      [0.05, 0.85, 0.05, 0.05],
...                      [0.05, 0.05, 0.85, 0.05],
...                      [0.05, 0.05, 0.05, 0.85]])
>>> target = torch.tensor([0, 1, 3, 2])
>>> multiclass_roc(pred, target)
((tensor([0., 0., 1.]), tensor([0., 1., 1.]), tensor([1.8500, 0.8500, 0.0500])),
 (tensor([0., 0., 1.]), tensor([0., 1., 1.]), tensor([1.8500, 0.8500, 0.0500])),
 (tensor([0.0000, 0.3333, 1.0000]), tensor([0., 0., 1.]), tensor([1.8500, 0.8500, ↵
↵0.0500])),
 (tensor([0.0000, 0.3333, 1.0000]), tensor([0., 0., 1.]), tensor([1.8500, 0.8500, ↵
↵0.0500])))
```

precision [func]

```
pytorch_lightning.metrics.functional.classification.precision(pred, target,
                                                             num_classes=None,
                                                             class_reduction='micro')
```

Computes precision score.

Parameters

- **pred** (Tensor) – estimated probabilities
- **target** (Tensor) – ground-truth labels
- **num_classes** (Optional[int]) – number of classes
- **class_reduction** (str) – method to reduce metric score over labels
 - 'micro': calculate metrics globally (default)
 - 'macro': calculate metrics for each label, and find their unweighted mean.
 - 'weighted': calculate metrics for each label, and find their weighted mean.
 - 'none': returns calculated metric per class

Return type Tensor

Returns Tensor with precision.

Example

```
>>> x = torch.tensor([0, 1, 2, 3])
>>> y = torch.tensor([0, 1, 2, 2])
>>> precision(x, y)
tensor(0.7500)
```

precision_recall [func]

```
pytorch_lightning.metrics.functional.classification.precision_recall(pred,
                                                                      target,
                                                                      num_classes=None,
                                                                      class_reduction='micro',
                                                                      re-
                                                                      turn_support=False,
                                                                      re-
                                                                      turn_state=False)
```

Computes precision and recall for different thresholds

Parameters

- **pred** (Tensor) – estimated probabilities
- **target** (Tensor) – ground-truth labels
- **num_classes** (Optional[int]) – number of classes
- **class_reduction** (str) – method to reduce metric score over labels
 - 'micro': calculate metrics globally (default)
 - 'macro': calculate metrics for each label, and find their unweighted mean.

- 'weighted': calculate metrics for each label, and find their weighted mean.
- 'none': returns calculated metric per class
- **return_support** *(bool)* – returns the support for each class, need for fbeta/f1 calculations
- **return_state** *(bool)* – returns a internal state that can be ddp reduced before doing the final calculation

Return type `Tuple[Tensor, Tensor]`

Returns Tensor with precision and recall

Example

```
>>> x = torch.tensor([0, 1, 2, 3])
>>> y = torch.tensor([0, 2, 2, 2])
>>> precision_recall(x, y, class_reduction='macro')
(tensor(0.5000), tensor(0.3333))
```

precision_recall_curve [func]

`pytorch_lightning.metrics.functional.classification.precision_recall_curve` (*pred, target, sample_weight=None, pos_label=1.0*)

Computes precision-recall pairs for different thresholds.

Parameters

- **pred** *(Tensor)* – estimated probabilities
- **target** *(Tensor)* – ground-truth labels
- **sample_weight** *(Optional[Sequence])* – sample weights
- **pos_label** *(int)* – the label for the positive class

Return type `Tuple[Tensor, Tensor, Tensor]`

Returns precision, recall, thresholds

Example

```
>>> pred = torch.tensor([0, 1, 2, 3])
>>> target = torch.tensor([0, 1, 1, 0])
>>> precision, recall, thresholds = precision_recall_curve(pred, target)
>>> precision
tensor([0.6667, 0.5000, 0.0000, 1.0000])
>>> recall
tensor([1.0000, 0.5000, 0.0000, 0.0000])
>>> thresholds
tensor([1, 2, 3])
```

recall [func]

```
pytorch_lightning.metrics.functional.classification.recall(pred, target,
                                                            num_classes=None,
                                                            class_reduction='micro')
```

Computes recall score.

Parameters

- **pred** (Tensor) – estimated probabilities
- **target** (Tensor) – ground-truth labels
- **num_classes** (Optional[int]) – number of classes
- **class_reduction** (str) – method to reduce metric score over labels
 - 'micro': calculate metrics globally (default)
 - 'macro': calculate metrics for each label, and find their unweighted mean.
 - 'weighted': calculate metrics for each label, and find their weighted mean.
 - 'none': returns calculated metric per class

Return type Tensor

Returns Tensor with recall.

Example

```
>>> x = torch.tensor([0, 1, 2, 3])
>>> y = torch.tensor([0, 1, 2, 2])
>>> recall(x, y)
tensor(0.7500)
```

roc [func]

```
pytorch_lightning.metrics.functional.classification.roc(pred, target, sample_weight=None,
                                                         pos_label=1.0)
```

Computes the Receiver Operating Characteristic (ROC). It assumes classifier is binary.

Parameters

- **pred** (Tensor) – estimated probabilities
- **target** (Tensor) – ground-truth labels
- **sample_weight** (Optional[Sequence]) – sample weights
- **pos_label** (int) – the label for the positive class

Return type Tuple[Tensor, Tensor, Tensor]

Returns false-positive rate (fpr), true-positive rate (tpr), thresholds

Example

```
>>> x = torch.tensor([0, 1, 2, 3])
>>> y = torch.tensor([0, 1, 1, 1])
>>> fpr, tpr, thresholds = roc(x, y)
>>> fpr
tensor([0., 0., 0., 0., 1.])
>>> tpr
tensor([0.0000, 0.3333, 0.6667, 1.0000, 1.0000])
>>> thresholds
tensor([4, 3, 2, 1, 0])
```

stat_scores [func]

`pytorch_lightning.metrics.functional.classification.stat_scores` (*pred*, *target*, *class_index*, *argmax_dim=1*)

Calculates the number of true positive, false positive, true negative and false negative for a specific class

Parameters

- **pred** (Tensor) – prediction tensor
- **target** (Tensor) – target tensor
- **class_index** (int) – class to calculate over
- **argmax_dim** (int) – if pred is a tensor of probabilities, this indicates the axis the argmax transformation will be applied over

Return type Tuple[`Tensor`, `Tensor`, `Tensor`, `Tensor`, `Tensor`]

Returns True Positive, False Positive, True Negative, False Negative, Support

Example

```
>>> x = torch.tensor([1, 2, 3])
>>> y = torch.tensor([0, 2, 3])
>>> tp, fp, tn, fn, sup = stat_scores(x, y, class_index=1)
>>> tp, fp, tn, fn, sup
(tensor(0), tensor(1), tensor(2), tensor(0), tensor(0))
```

stat_scores_multiple_classes [func]

`pytorch_lightning.metrics.functional.classification.stat_scores_multiple_classes` (*pred*, *target*, *num_classes*, *argmax_dim*, *reduction='none'*)

Calculates the number of true positive, false positive, true negative and false negative for each class

Parameters

- **pred** (Tensor) – prediction tensor
- **target** (Tensor) – target tensor
- **num_classes** (Optional[int]) – number of classes if known
- **argmax_dim** (int) – if pred is a tensor of probabilities, this indicates the axis the argmax transformation will be applied over
- **reduction** (str) – a method to reduce metric score over labels (default: none) Available reduction methods:
 - `elementwise_mean`: takes the mean
 - `none`: pass array
 - `sum`: add elements

Return type Tuple[Tensor, Tensor, Tensor, Tensor, Tensor]

Returns True Positive, False Positive, True Negative, False Negative, Support

Example

```
>>> x = torch.tensor([1, 2, 3])
>>> y = torch.tensor([0, 2, 3])
>>> tps, fps, tns, fns, sups = stat_scores_multiple_classes(x, y)
>>> tps
tensor([0., 0., 1., 1.])
>>> fps
tensor([0., 1., 0., 0.])
>>> tns
tensor([2., 2., 2., 2.])
>>> fns
tensor([1., 0., 0., 0.])
>>> sups
tensor([1., 0., 1., 1.])
```

to_categorical [func]

pytorch_lightning.metrics.functional.classification.to_categorical(tensor, *argmax_dim=1*)

Converts a tensor of probabilities to a dense label tensor

Parameters

- **tensor** (Tensor) – probabilities to get the categorical label [N, d1, d2, ...]
- **argmax_dim** (int) – dimension to apply

Return type Tensor

Returns A tensor with categorical labels [N, d2, ...]

Example

```
>>> x = torch.tensor([[0.2, 0.5], [0.9, 0.1]])
>>> to_categorical(x)
tensor([1, 0])
```

to_onehot [func]

`pytorch_lightning.metrics.functional.classification.to_onehot` (*tensor*,
num_classes=None)

Converts a dense label tensor to one-hot format

Parameters

- **tensor** (Tensor) – dense label tensor, with shape [N, d1, d2, ...]
- **num_classes** (Optional[int]) – number of classes C

Output: A sparse label tensor with shape [N, C, d1, d2, ...]

Example

```
>>> x = torch.tensor([1, 2, 3])
>>> to_onehot(x)
tensor([[0, 1, 0, 0],
        [0, 0, 1, 0],
        [0, 0, 0, 1]])
```

Return type Tensor

11.4.2 Regression

mae [func]

`pytorch_lightning.metrics.functional.regression.mae` (*pred*, *target*, *reduction='elementwise_mean'*,
return_state=False)

Computes mean absolute error

Parameters

- **pred** (Tensor) – estimated labels
- **target** (Tensor) – ground truth labels
- **reduction** (str) – a method to reduce metric score over labels.
 - 'elementwise_mean': takes the mean (default)
 - 'sum': takes the sum
 - 'none': no reduction will be applied
- **return_state** (bool) – returns a internal state that can be ddp reduced before doing the final calculation

Return type Tensor

Returns Tensor with MAE

Example

```
>>> x = torch.tensor([0., 1, 2, 3])
>>> y = torch.tensor([0., 1, 2, 2])
>>> mae(x, y)
tensor(0.2500)
```

mse [func]

`pytorch_lightning.metrics.functional.regression.mse(pred, target, reduction='elementwise_mean', return_state=False)`

Computes mean squared error

Parameters

- **pred** (Tensor) – estimated labels
- **target** (Tensor) – ground truth labels
- **reduction** (str) – a method to reduce metric score over labels.
 - 'elementwise_mean': takes the mean (default)
 - 'sum': takes the sum
 - 'none': no reduction will be applied
- **return_state** (bool) – returns a internal state that can be ddp reduced before doing the final calculation

Return type Tensor

Returns Tensor with MSE

Example

```
>>> x = torch.tensor([0., 1, 2, 3])
>>> y = torch.tensor([0., 1, 2, 2])
>>> mse(x, y)
tensor(0.2500)
```

psnr [func]

`pytorch_lightning.metrics.functional.regression.psnr(pred, target, data_range=None, base=10.0, reduction='elementwise_mean', return_state=False)`

Computes the peak signal-to-noise ratio

Parameters

- **pred** (Tensor) – estimated signal
- **target** (Tensor) – ground truth signal

- **data_range** (Optional[float]) – the range of the data. If None, it is determined from the data (max - min)
- **base** (float) – a base of a logarithm to use (default: 10)
- **reduction** (str) – a method to reduce metric score over labels.
 - 'elementwise_mean': takes the mean (default)
 - 'sum': takes the sum
 - 'none': no reduction will be applied
- **return_state** (bool) – returns a internal state that can be ddp reduced before doing the final calculation

Return type Tensor

Returns Tensor with PSNR score

Example

```
>>> pred = torch.tensor([[0.0, 1.0], [2.0, 3.0]])
>>> target = torch.tensor([[3.0, 2.0], [1.0, 0.0]])
>>> psnr(pred, target)
tensor(2.5527)
```

rmse [func]

pytorch_lightning.metrics.functional.regression.**rmse** (*pred*, *target*, *reduction*='elementwise_mean', *return_state*=False)

Computes root mean squared error

Parameters

- **pred** (Tensor) – estimated labels
- **target** (Tensor) – ground truth labels
- **reduction** (str) – a method to reduce metric score over labels.
 - 'elementwise_mean': takes the mean (default)
 - 'sum': takes the sum
 - 'none': no reduction will be applied
- **return_state** (bool) – returns a internal state that can be ddp reduced before doing the final calculation

Return type Tensor

Returns

Tensor with RMSE

```
>>> x = torch.tensor([0., 1, 2, 3])
>>> y = torch.tensor([0., 1, 2, 2])
>>> rmse(x, y)
tensor(0.5000)
```

rmsle [func]

`pytorch_lightning.metrics.functional.regression.rmsle(pred, target, reduction='elementwise_mean')`

Computes root mean squared log error

Parameters

- **pred** `(Tensor)` – estimated labels
- **target** `(Tensor)` – ground truth labels
- **reduction** `(str)` – a method to reduce metric score over labels.
 - 'elementwise_mean': takes the mean (default)
 - 'sum': takes the sum
 - 'none': no reduction will be applied

Return type `Tensor`

Returns Tensor with RMSLE

Example

```
>>> x = torch.tensor([0., 1, 2, 3])
>>> y = torch.tensor([0., 1, 2, 2])
>>> rmsle(x, y)
tensor(0.1438)
```

ssim [func]

`pytorch_lightning.metrics.functional.regression.mae(pred, target, reduction='elementwise_mean', return_state=False)`

Computes mean absolute error

Parameters

- **pred** `(Tensor)` – estimated labels
- **target** `(Tensor)` – ground truth labels
- **reduction** `(str)` – a method to reduce metric score over labels.
 - 'elementwise_mean': takes the mean (default)
 - 'sum': takes the sum
 - 'none': no reduction will be applied
- **return_state** `(bool)` – returns a internal state that can be ddp reduced before doing the final calculation

Return type `Tensor`

Returns Tensor with MAE

Example

```
>>> x = torch.tensor([0., 1, 2, 3])
>>> y = torch.tensor([0., 1, 2, 2])
>>> mae(x, y)
tensor(0.2500)
```

11.4.3 NLP

bleu_score [func]

`pytorch_lightning.metrics.functional.nlp.bleu_score` (*translate_corpus*, *reference_corpus*, *n_gram=4*, *smooth=False*)

Calculate BLEU score of machine translated text with one or more references

Parameters

- **translate_corpus** *(Sequence[str])* – An iterable of machine translated corpus
- **reference_corpus** *(Sequence[str])* – An iterable of iterables of reference corpus
- **n_gram** *(int)* – Gram value ranged from 1 to 4 (Default 4)
- **smooth** *(bool)* – Whether or not to apply smoothing – Lin et al. 2004

Return type `Tensor`

Returns Tensor with BLEU Score

Example

```
>>> translate_corpus = ['the cat is on the mat'.split()]
>>> reference_corpus = [['there is a cat on the mat'.split(), 'a cat is on the mat
↪'.split()]]
>>> bleu_score(translate_corpus, reference_corpus)
tensor(0.7598)
```

11.4.4 Pairwise

embedding_similarity [func]

`pytorch_lightning.metrics.functional.self_supervised.embedding_similarity` (*batch*, *similarity='cosine'*, *reduction='none'*, *zero_diagonal=True*)

Computes representation similarity

Example

```
>>> embeddings = torch.tensor([[1., 2., 3., 4.], [1., 2., 3., 4.], [4., 5., 6., 7.
↵]])
>>> embedding_similarity(embeddings)
tensor([[0.0000, 1.0000, 0.9759],
        [1.0000, 0.0000, 0.9759],
        [0.9759, 0.9759, 0.0000]])
```

Parameters

- **batch** *(Tensor)* – (batch, dim)
- **similarity** *(str)* – ‘dot’ or ‘cosine’
- **reduction** *(str)* – ‘none’, ‘sum’, ‘mean’ (all along dim -1)
- **zero_diagonal** *(bool)* – if True, the diagonals are set to zero

Return type *Tensor*

Returns A square matrix (batch, batch) with the similarity scores between all elements. If sum or mean are used, then returns (b, 1) with the reduced value for each row.

STEP-BY-STEP WALK-THROUGH

This guide will walk you through the core pieces of PyTorch Lightning.

We'll accomplish the following:

- Implement an MNIST classifier.
- Use inheritance to implement an AutoEncoder

Note: Any DL/ML PyTorch project fits into the Lightning structure. Here we just focus on 3 types of research to illustrate.

12.1 From MNIST to AutoEncoders

12.1.1 Installing Lightning

Lightning is trivial to install. We recommend using conda environments

```
conda activate my_env
pip install pytorch-lightning
```

Or without conda environments, use pip.

```
pip install pytorch-lightning
```

Or conda.

```
conda install pytorch-lightning -c conda-forge
```

12.1.2 The research

The Model

The *LightningModule* holds all the core research ingredients:

- The model
- The optimizers
- The train/ val/ test steps

Let's first start with the model. In this case, we'll design a 3-layer neural network.

```
import torch
from torch.nn import functional as F
from torch import nn
from pytorch_lightning.core.lightning import LightningModule

class LitMNIST(LightningModule):

    def __init__(self):
        super().__init__()

        # mnist images are (1, 28, 28) (channels, width, height)
        self.layer_1 = torch.nn.Linear(28 * 28, 128)
        self.layer_2 = torch.nn.Linear(128, 256)
        self.layer_3 = torch.nn.Linear(256, 10)

    def forward(self, x):
        batch_size, channels, width, height = x.size()

        # (b, 1, 28, 28) -> (b, 1*28*28)
        x = x.view(batch_size, -1)
        x = self.layer_1(x)
        x = F.relu(x)
        x = self.layer_2(x)
        x = F.relu(x)
        x = self.layer_3(x)

        x = F.log_softmax(x, dim=1)
        return x
```

Notice this is a *LightningModule* instead of a `torch.nn.Module`. A `LightningModule` is equivalent to a pure PyTorch Module except it has added functionality. However, you can use it **EXACTLY** the same as you would a PyTorch Module.

```
net = LitMNIST()
x = torch.randn(1, 1, 28, 28)
out = net(x)
```

Out:

```
torch.Size([1, 10])
```

Now we add the `training_step` which has all our training loop logic

```
class LitMNIST(LightningModule):
```

(continues on next page)

(continued from previous page)

```
def training_step(self, batch, batch_idx):
    x, y = batch
    logits = self(x)
    loss = F.nll_loss(logits, y)
    return loss
```

Data

Lightning operates on pure dataloaders. Here's the PyTorch code for loading MNIST.

```
from torch.utils.data import DataLoader, random_split
from torchvision.datasets import MNIST
import os
from torchvision import datasets, transforms

# transforms
# prepare transforms standard to MNIST
transform=transforms.Compose([transforms.ToTensor(),
                             transforms.Normalize((0.1307,), (0.3081,))])

# data
mnist_train = MNIST(os.getcwd(), train=True, download=True, transform=transform)
mnist_train = DataLoader(mnist_train, batch_size=64)
```

You can use DataLoaders in 3 ways:

1. Pass DataLoaders to .fit()

Pass in the dataloaders to the .fit() function.

```
model = LitMNIST()
trainer = Trainer()
trainer.fit(model, mnist_train)
```

2. LightningModule DataLoaders

For fast research prototyping, it might be easier to link the model with the dataloaders.

```
class LitMNIST(pl.LightningModule):

    def train_dataloader(self):
        # transforms
        # prepare transforms standard to MNIST
        transform=transforms.Compose([transforms.ToTensor(),
                                     transforms.Normalize((0.1307,), (0.3081,))])

        # data
        mnist_train = MNIST(os.getcwd(), train=True, download=True,
                             transform=transform)
        return DataLoader(mnist_train, batch_size=64)

    def val_dataloader(self):
        transforms = ...
```

(continues on next page)

(continued from previous page)

```

mnist_val = ...
return DataLoader(mnist_val, batch_size=64)

def test_dataloader(self):
    transforms = ...
    mnist_test = ...
    return DataLoader(mnist_test, batch_size=64)

```

DataLoaders are already in the model, no need to specify on `.fit()`.

```

model = LitMNIST()
trainer = Trainer()
trainer.fit(model)

```

3. DataModules (recommended)

Defining free-floating dataloaders, splits, download instructions, and such can get messy. In this case, it's better to group the full definition of a dataset into a *DataModule* which includes:

- Download instructions
- Processing instructions
- Split instructions
- Train dataloader
- Val dataloader(s)
- Test dataloader(s)

```

class MyDataModule(LightningDataModule):

    def __init__(self):
        super().__init__()
        self.train_dims = None
        self.vocab_size = 0

    def prepare_data(self):
        # called only on 1 GPU
        download_dataset()
        tokenize()
        build_vocab()

    def setup(self):
        # called on every GPU
        vocab = load_vocab()
        self.vocab_size = len(vocab)

        self.train, self.val, self.test = load_datasets()
        self.train_dims = self.train.next_batch.size()

    def train_dataloader(self):
        transforms = ...
        return DataLoader(self.train, batch_size=64)

    def val_dataloader(self):

```

(continues on next page)

(continued from previous page)

```

        transforms = ...
        return DataLoader(self.val, batch_size=64)

    def test_dataloader(self):
        transforms = ...
        return DataLoader(self.test, batch_size=64)

```

Using DataModules allows easier sharing of full dataset definitions.

```

# use an MNIST dataset
mnist_dm = MNISTDataModule()
model = LitModel(num_classes=mnist_dm.num_classes)
trainer.fit(model, mnist_dm)

# or other datasets with the same model
imagenet_dm = ImagenetDataModule()
model = LitModel(num_classes=imagenet_dm.num_classes)
trainer.fit(model, imagenet_dm)

```

Note: `prepare_data()` is called on only one GPU in distributed training (automatically)

Note: `setup()` is called on every GPU (automatically)

Models defined by data

When your models need to know about the data, it's best to process the data before passing it to the model.

```

# init dm AND call the processing manually
dm = ImagenetDataModule()
dm.prepare_data()
dm.setup()

model = LitModel(out_features=dm.num_classes, img_width=dm.img_width, img_height=dm.
    ↪img_height)
trainer.fit(model, dm)

```

1. use `prepare_data()` to download and process the dataset.
2. use `setup()` to do splits, and build your model internals

An alternative to using a DataModule is to defer initialization of the models modules to the `setup` method of your LightningModule as follows:

```

class LitMNIST(LightningModule):

    def __init__(self):
        self.ll = None

```

(continues on next page)

(continued from previous page)

```
def prepare_data(self):
    download_data()
    tokenize()

def setup(self, step):
    # step is either 'fit' or 'test' 90% of the time not relevant
    data = load_data()
    num_classes = data.classes
    self.ll = nn.Linear(..., num_classes)
```

Optimizer

Next we choose what optimizer to use for training our system. In PyTorch we do it as follows:

```
from torch.optim import Adam
optimizer = Adam(LitMNIST().parameters(), lr=1e-3)
```

In Lightning we do the same but organize it under the `configure_optimizers()` method.

```
class LitMNIST(LightningModule):

    def configure_optimizers(self):
        return Adam(self.parameters(), lr=1e-3)
```

Note: The LightningModule itself has the parameters, so pass in `self.parameters()`

However, if you have multiple optimizers use the matching parameters

```
class LitMNIST(LightningModule):

    def configure_optimizers(self):
        return Adam(self.generator(), lr=1e-3), Adam(self.discriminator(), lr=1e-3)
```

Training step

The training step is what happens inside the training loop.

```
for epoch in epochs:
    for batch in data:
        # TRAINING STEP
        # ....
        # TRAINING STEP
        loss.backward()
        optimizer.step()
        optimizer.zero_grad()
```

In the case of MNIST, we do the following

```
for epoch in epochs:
    for batch in data:
        # ----- TRAINING STEP START -----
```

(continues on next page)

(continued from previous page)

```

x, y = batch
logits = model(x)
loss = F.nll_loss(logits, y)
# ----- TRAINING STEP END -----

loss.backward()
optimizer.step()
optimizer.zero_grad()

```

In Lightning, everything that is in the training step gets organized under the `training_step()` function in the `LightningModule`.

```

class LitMNIST(LightningModule):

    def training_step(self, batch, batch_idx):
        x, y = batch
        logits = self(x)
        loss = F.nll_loss(logits, y)
        return loss

```

Again, this is the same PyTorch code except that it has been organized by the `LightningModule`. This code is not restricted which means it can be as complicated as a full seq-2-seq, RL loop, GAN, etc...

12.1.3 The engineering

Training

So far we defined 4 key ingredients in pure PyTorch but organized the code with the `LightningModule`.

1. Model.
2. Training data.
3. Optimizer.
4. What happens in the training loop.

For clarity, we'll recall that the full `LightningModule` now looks like this.

```

class LitMNIST(LightningModule):
    def __init__(self):
        super().__init__()
        self.layer_1 = torch.nn.Linear(28 * 28, 128)
        self.layer_2 = torch.nn.Linear(128, 256)
        self.layer_3 = torch.nn.Linear(256, 10)

    def forward(self, x):
        batch_size, channels, width, height = x.size()
        x = x.view(batch_size, -1)
        x = self.layer_1(x)

```

(continues on next page)

(continued from previous page)

```
x = F.relu(x)
x = self.layer_2(x)
x = F.relu(x)
x = self.layer_3(x)
x = F.log_softmax(x, dim=1)
return x

def training_step(self, batch, batch_idx):
    x, y = batch
    logits = self(x)
    loss = F.nll_loss(logits, y)
    return loss
```

Again, this is the same PyTorch code, except that it's organized by the LightningModule.

Logging

To log to Tensorboard, your favorite logger, and/or the progress bar, use the `log()` method which can be called from any method in the LightningModule.

```
def training_step(self, batch, batch_idx):
    self.log('my_metric', x)
```

The `log()` method has a few options:

- `on_step` (logs the metric at that step in training)
- `on_epoch` (automatically accumulates and logs at the end of the epoch)
- `prog_bar` (logs to the progress bar)
- `logger` (logs to the logger like Tensorboard)

Depending on where the log is called from, Lightning auto-determines the correct mode for you. But of course you can override the default behavior by manually setting the flags.

Note: Setting `on_epoch=True` will accumulate your logged values over the full training epoch.

```
def training_step(self, batch, batch_idx):
    self.log('my_loss', loss, on_step=True, on_epoch=True, prog_bar=True, logger=True)
```

You can also use any method of your logger directly:

```
def training_step(self, batch, batch_idx):
    tensorboard = self.logger.experiment
    tensorboard.any_summary_writer_method_you_want()
```

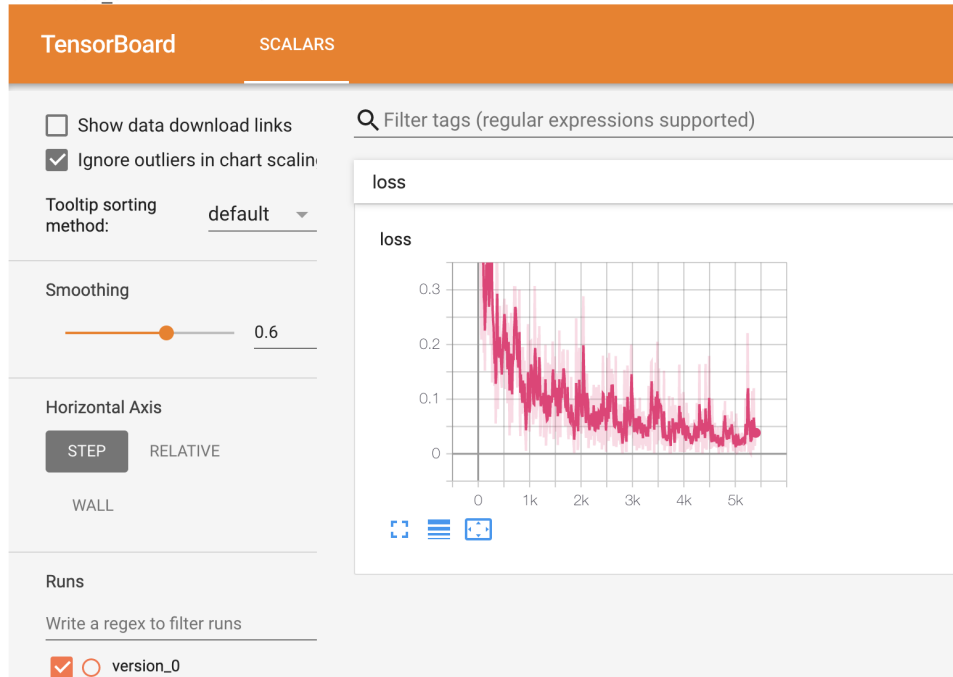
Once your training starts, you can view the logs by using your favorite logger or booting up the Tensorboard logs:

```
tensorboard --logdir ./lightning_logs
```

Which will generate automatic tensorboard logs (or with the logger of your choice).

```
[31] # Start tensorboard.
      %load_ext tensorboard
      %tensorboard --logdir lightning_logs/
```

↗ The tensorboard extension is already loaded. To reload it, use:
 %reload_ext tensorboard



But you can also use any of the *number of other loggers* we support.

Train on CPU

```
from pytorch_lightning import Trainer

model = LitMNIST()
trainer = Trainer()
trainer.fit(model, train_loader)
```

You should see the following weights summary and progress bar

	Name	Type	Params
0	layer_1	Linear	100 K
1	layer_2	Linear	33 K
2	layer_3	Linear	2 K

Epoch 1: 27% 250/938 [00:08<00:22, 30.10it/s, loss=0.353, v_num=2]


Train on GPU

But the beauty is all the magic you can do with the trainer flags. For instance, to run this model on a GPU:

```
model = LitMNIST()
trainer = Trainer(gpus=1)
trainer.fit(model, train_loader)
```

```
INFO:root:GPU available: True, used: True
INFO:root:VISIBLE_GPUS: 0
INFO:root:
```

	Name	Type	Params
0	layer_1	Linear	100 K
1	layer_2	Linear	33 K
2	layer_3	Linear	2 K

```
Epoch 1: 53%  500/938 [00:07<00:07, 61.53it/s, loss=0.206, v_num=3]
```

Train on Multi-GPU

Or you can also train on multiple GPUs.

```
model = LitMNIST()
trainer = Trainer(gpus=8)
trainer.fit(model, train_loader)
```

Or multiple nodes

```
# (32 GPUs)
model = LitMNIST()
trainer = Trainer(gpus=8, num_nodes=4, distributed_backend='ddp')
trainer.fit(model, train_loader)
```

Refer to the [distributed computing guide for more details](#).

Train on TPUs

Did you know you can use PyTorch on TPUs? It's very hard to do, but we've worked with the xla team to use their awesome library to get this to work out of the box!

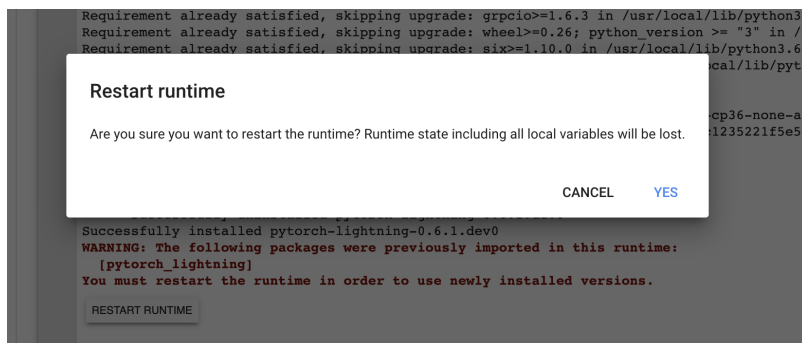
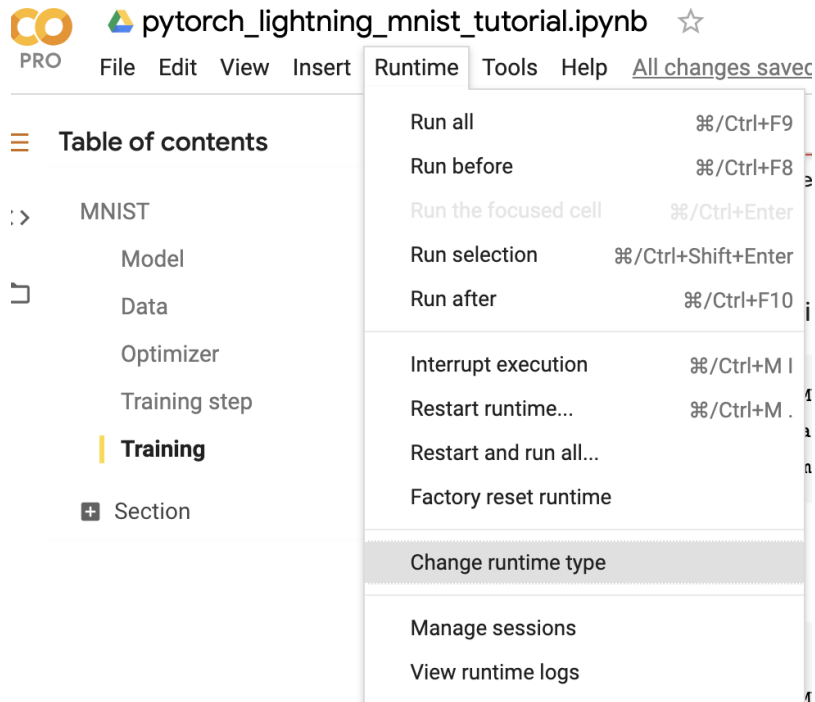
Let's train on Colab ([full demo available here](#))

First, change the runtime to TPU (and reinstall lightning).

Next, install the required xla library (adds support for PyTorch on TPUs)

```
!curl https://raw.githubusercontent.com/pytorch/xla/master/contrib/scripts/env-setup.
↪py -o pytorch-xla-env-setup.py

!python pytorch-xla-env-setup.py --version nightly --apt-packages libomp5 libopenblas-
↪dev
```



In distributed training (multiple GPUs and multiple TPU cores) each GPU or TPU core will run a copy of this program. This means that without taking any care you will download the dataset N times which will cause all sorts of issues.

To solve this problem, make sure your download code is in the `prepare_data` method in the `DataModule`. In this method we do all the preparation we need to do once (instead of on every GPU).

`prepare_data` can be called in two ways, once per node or only on the root node (`Trainer(prepare_data_per_node=False)`).

```
class MNISTDataModule(LightningDataModule):
    def __init__(self, batch_size=64):
        super().__init__()
        self.batch_size = batch_size

    def prepare_data(self):
        # download only
        MNIST(os.getcwd(), train=True, download=True, transform=transforms.ToTensor())
        MNIST(os.getcwd(), train=False, download=True, transform=transforms.
↪ToTensor())

    def setup(self, stage):
        # transform
        transform=transforms.Compose([transforms.ToTensor()])
        MNIST(os.getcwd(), train=True, download=False, transform=transform)
        MNIST(os.getcwd(), train=False, download=False, transform=transform)

        # train/val split
        mnist_train, mnist_val = random_split(mnist_train, [55000, 5000])

        # assign to use in dataloaders
        self.train_dataset = mnist_train
        self.val_dataset = mnist_val
        self.test_dataset = mnist_test

    def train_dataloader(self):
        return DataLoader(self.train_dataset, batch_size=self.batch_size)

    def val_dataloader(self):
        return DataLoader(self.val_dataset, batch_size=self.batch_size)

    def test_dataloader(self):
        return DataLoader(self.test_dataset, batch_size=self.batch_size)
```

The `prepare_data` method is also a good place to do any data processing that needs to be done only once (ie: download or tokenize, etc...).

Note: Lightning inserts the correct `DistributedSampler` for distributed training. No need to add yourself!

Now we can train the `LightningModule` on a TPU without doing anything else!

```
dm = MNISTDataModule()
model = LitMNIST()
trainer = Trainer(tpu_cores=8)
trainer.fit(model, dm)
```

You'll now see the TPU cores booting up.

Notice the epoch is MUCH faster!

```
INFO:root:training on 8 TPU cores
INFO:root:INIT TPU local core: 0, global rank: 0
INFO:root:INIT TPU local core: 3, global rank: 3
INFO:root:INIT TPU local core: 1, global rank: 1
```

```
INFO:root:
| Name      | Type   | Params
-----
0 | layer_1   | Linear | 100 K
1 | layer_2   | Linear | 33 K
2 | layer_3   | Linear | 2 K
Using downloaded and verified file: /content/MNIST/raw/train-images-idx3-ubyte.gz
Extracting /content/MNIST/raw/train-images-idx3-ubyte.gz to /content/MNIST/raw
Using downloaded and verified file: /content/MNIST/raw/train-labels-idx1-ubyte.gz
Extracting /content/MNIST/raw/train-labels-idx1-ubyte.gz to /content/MNIST/raw
Using downloaded and verified file: /content/MNIST/raw/t10k-images-idx3-ubyte.gz
Extracting /content/MNIST/raw/t10k-images-idx3-ubyte.gz to /content/MNIST/raw
Using downloaded and verified file: /content/MNIST/raw/t10k-labels-idx1-ubyte.gz
Extracting /content/MNIST/raw/t10k-labels-idx1-ubyte.gz to /content/MNIST/raw
Processing...
Done!
Epoch 6: 42%  50/118 [00:00<00:01, 62.22it/s, loss=0.067, v_num=10]
```

Hyperparameters

Lightning has utilities to interact seamlessly with the command line `ArgumentParser` and plays well with the hyperparameter optimization framework of your choice.

ArgumentParser

Lightning is designed to augment a lot of the functionality of the built-in Python `ArgumentParser`

```
from argparse import ArgumentParser
parser = ArgumentParser()
parser.add_argument('--layer_1_dim', type=int, default=128)
args = parser.parse_args()
```

This allows you to call your program like so:

```
python trainer.py --layer_1_dim 64
```

Argparser Best Practices

It is best practice to layer your arguments in three sections.

1. Trainer args (gpus, num_nodes, etc...)
2. Model specific arguments (layer_dim, num_layers, learning_rate, etc...)
3. Program arguments (data_path, cluster_email, etc...)

We can do this as follows. First, in your `LightningModule`, define the arguments specific to that module. Remember that data splits or data paths may also be specific to a module (i.e.: if your project has a model that trains on Imagenet and another on CIFAR-10).

```
class LitModel(LightningModule):  
  
    @staticmethod  
    def add_model_specific_args(parent_parser):  
        parser = ArgumentParser(parents=[parent_parser], add_help=False)  
        parser.add_argument('--encoder_layers', type=int, default=12)  
        parser.add_argument('--data_path', type=str, default='/some/path')  
        return parser
```

Now in your main trainer file, add the `Trainer` args, the program args, and add the model args

```
# -----  
# trainer_main.py  
# -----  
from argparse import ArgumentParser  
parser = ArgumentParser()  
  
# add PROGRAM level args  
parser.add_argument('--conda_env', type=str, default='some_name')  
parser.add_argument('--notification_email', type=str, default='will@email.com')  
  
# add model specific args  
parser = LitModel.add_model_specific_args(parser)  
  
# add all the available trainer options to argparse  
# ie: now --gpus --num_nodes ... --fast_dev_run all work in the cli  
parser = Trainer.add_argparse_args(parser)  
  
args = parser.parse_args()
```

Now you can call run your program like so:

```
python trainer_main.py --gpus 2 --num_nodes 2 --conda_env 'my_env' --encoder_layers 12
```

Finally, make sure to start the training like so:

```
# init the trainer like this  
trainer = Trainer.from_argparse_args(args, early_stopping_callback=...)  
  
# NOT like this  
trainer = Trainer(gpus=hparams.gpus, ...)  
  
# init the model with Namespace directly  
model = LitModel(args)  
  
# or init the model with all the key-value pairs  
dict_args = vars(args)  
model = LitModel(**dict_args)
```


LightningModule hyperparameters

Often times we train many versions of a model. You might share that model or come back to it a few months later at which point it is very useful to know how that model was trained (i.e.: what learning rate, neural network, etc...).

Lightning has a few ways of saving that information for you in checkpoints and yaml files. The goal here is to improve readability and reproducibility

1. The first way is to ask lightning to save the values of anything in the `__init__` for you to the checkpoint. This also makes those values available via `self.hparams`.

```
class LitMNIST(LightningModule):

    def __init__(self, layer_1_dim=128, learning_rate=1e-2, **kwargs):
        super().__init__()
        # call this to save (layer_1_dim=128, learning_rate=1e-4) to the checkpoint
        self.save_hyperparameters()

        # equivalent
        self.save_hyperparameters('layer_1_dim', 'learning_rate')

        # this now works
        self.hparams.layer_1_dim
```

2. Sometimes your init might have objects or other parameters you might not want to save. In that case, choose only a few

```
class LitMNIST(LightningModule):

    def __init__(self, loss_fx, generator_network, layer_1_dim=128 **kwargs):
        super().__init__()
        self.layer_1_dim = layer_1_dim
        self.loss_fx = loss_fx

        # call this to save (layer_1_dim=128) to the checkpoint
        self.save_hyperparameters('layer_1_dim')

# to load specify the other args
model = LitMNIST.load_from_checkpoint(PATH, loss_fx=torch.nn.SomeOtherLoss, generator_
↪network=MyGenerator())
```

3. Assign to `self.hparams`. Anything assigned to `self.hparams` will also be saved automatically

```
# using a argparse.Namespace
class LitMNIST(LightningModule):

    def __init__(self, hparams, *args, **kwargs):
        super().__init__()
        self.hparams = hparams

        self.layer_1 = torch.nn.Linear(28 * 28, self.hparams.layer_1_dim)
        self.layer_2 = torch.nn.Linear(self.hparams.layer_1_dim, self.hparams.layer_2_
↪dim)
        self.layer_3 = torch.nn.Linear(self.hparams.layer_2_dim, 10)

    def train_dataloader(self):
        return DataLoader(mnist_train, batch_size=self.hparams.batch_size)
```

4. You can also save full objects such as `dict` or `Namespace` to the checkpoint.

```
# using a argparse.Namespace
class LitMNIST(LightningModule):

    def __init__(self, conf, *args, **kwargs):
        super().__init__()
        self.hparams = conf

        # equivalent
        self.save_hyperparameters(conf)

        self.layer_1 = torch.nn.Linear(28 * 28, self.hparams.layer_1_dim)
        self.layer_2 = torch.nn.Linear(self.hparams.layer_1_dim, self.hparams.layer_2_
↪ dim)
        self.layer_3 = torch.nn.Linear(self.hparams.layer_2_dim, 10)

conf = OmegaConf.create(...)
model = LitMNIST(conf)

# this works
model.hparams.anything
```

Trainer args

To recap, add ALL possible trainer flags to the argparse and init the Trainer this way

```
parser = ArgumentParser()
parser = Trainer.add_argparse_args(parser)
hparams = parser.parse_args()

trainer = Trainer.from_argparse_args(hparams)

# or if you need to pass in callbacks
trainer = Trainer.from_argparse_args(hparams, checkpoint_callback=..., callbacks=[...
↪ ])
```

Multiple Lightning Modules

We often have multiple Lightning Modules where each one has different arguments. Instead of polluting the main.py file, the LightningModule lets you define arguments for each one.

```
class LitMNIST(LightningModule):

    def __init__(self, layer_1_dim, **kwargs):
        super().__init__()
        self.layer_1 = torch.nn.Linear(28 * 28, layer_1_dim)

    @staticmethod
    def add_model_specific_args(parent_parser):
        parser = ArgumentParser(parents=[parent_parser], add_help=False)
        parser.add_argument('--layer_1_dim', type=int, default=128)
        return parser
```

```

class GoodGAN(LightningModule):

    def __init__(self, encoder_layers, **kwargs):
        super().__init__()
        self.encoder = Encoder(layers=encoder_layers)

    @staticmethod
    def add_model_specific_args(parent_parser):
        parser = ArgumentParser(parents=[parent_parser], add_help=False)
        parser.add_argument('--encoder_layers', type=int, default=12)
        return parser

```

Now we can allow each model to inject the arguments it needs in the main.py

```

def main(args):
    dict_args = vars(args)

    # pick model
    if args.model_name == 'gan':
        model = GoodGAN(**dict_args)
    elif args.model_name == 'mnist':
        model = LitMNIST(**dict_args)

    trainer = Trainer.from_argparse_args(args)
    trainer.fit(model)

if __name__ == '__main__':
    parser = ArgumentParser()
    parser = Trainer.add_argparse_args(parser)

    # figure out which model to use
    parser.add_argument('--model_name', type=str, default='gan', help='gan or mnist')

    # THIS LINE IS KEY TO PULL THE MODEL NAME
    temp_args, _ = parser.parse_known_args()

    # let the model add what it wants
    if temp_args.model_name == 'gan':
        parser = GoodGAN.add_model_specific_args(parser)
    elif temp_args.model_name == 'mnist':
        parser = LitMNIST.add_model_specific_args(parser)

    args = parser.parse_args()

    # train
    main(args)

```

and now we can train MNIST or the GAN using the command line interface!

```

$ python main.py --model_name gan --encoder_layers 24
$ python main.py --model_name mnist --layer_1_dim 128

```

Validating

For most cases, we stop training the model when the performance on a validation split of the data reaches a minimum.

Just like the `training_step`, we can define a `validation_step` to check whatever metrics we care about, generate samples, or add more to our logs.

```
def validation_step(self, batch, batch_idx):
    loss = MSE_loss(...)
    self.log('val_loss', loss)
```

Now we can train with a validation loop as well.

```
from pytorch_lightning import Trainer

model = LitMNIST()
trainer = Trainer(tpu_cores=8)
trainer.fit(model, train_loader, val_loader)
```

You may have noticed the words **Validation sanity check** logged. This is because Lightning runs 2 batches of validation before starting to train. This is a kind of unit test to make sure that if you have a bug in the validation loop, you won't need to potentially wait for a full epoch to find out.

Note: Lightning disables gradients, puts model in eval mode, and does everything needed for validation.

Val loop under the hood

Under the hood, Lightning does the following:

```
model = Model()
model.train()
torch.set_grad_enabled(True)

for epoch in epochs:
    for batch in data:
        # ...
        # train

    # validate
    model.eval()
    torch.set_grad_enabled(False)

    outputs = []
    for batch in val_data:
        x, y = batch
        y_hat = model(x)
        loss = loss(y_hat, x)
        outputs.append({'val_loss': loss})

    total_loss = outputs.mean()
```

Optional methods

If you still need even more fine-grain control, define the other optional methods for the loop.

```
def validation_step(self, batch, batch_idx):
    preds = ...
    return preds

def validation_epoch_end(self, val_step_outputs):
    for pred in val_step_outputs:
        # do something with all the predictions from each validation_step
```

Testing

Once our research is done and we're about to publish or deploy a model, we normally want to figure out how it will generalize in the "real world." For this, we use a held-out split of the data for testing.

Just like the validation loop, we define a test loop

```
class LitMNIST(LightningModule):
    def test_step(self, batch, batch_idx):
        x, y = batch
        logits = self(x)
        loss = F.nll_loss(logits, y)
        self.log('test_loss', loss)
```

However, to make sure the test set isn't used inadvertently, Lightning has a separate API to run tests. Once you train your model simply call `.test()`.

```
from pytorch_lightning import Trainer

model = LitMNIST()
trainer = Trainer(tpu_cores=8)
trainer.fit(model)

# run test set
result = trainer.test()
print(result)
```

Out:

```
-----
TEST RESULTS
{'test_loss': tensor(1.1703, device='cuda:0')}
-----
```

You can also run the test from a saved lightning model

```
model = LitMNIST.load_from_checkpoint(PATH)
trainer = Trainer(tpu_cores=8)
trainer.test(model)
```

Note: Lightning disables gradients, puts model in eval mode, and does everything needed for testing.

Warning: `.test()` is not stable yet on TPUs. We're working on getting around the multiprocessing challenges.

Predicting

Again, a `LightningModule` is exactly the same as a PyTorch module. This means you can load it and use it for prediction.

```
model = LitMNIST.load_from_checkpoint(PATH)
x = torch.randn(1, 1, 28, 28)
out = model(x)
```

On the surface, it looks like `forward` and `training_step` are similar. Generally, we want to make sure that what we want the model to do is what happens in the `forward`, whereas the `training_step` likely calls `forward` from within it.

```
class MNISTClassifier(LightningModule):

    def forward(self, x):
        batch_size, channels, width, height = x.size()
        x = x.view(batch_size, -1)
        x = self.layer_1(x)
        x = F.relu(x)
        x = self.layer_2(x)
        x = F.relu(x)
        x = self.layer_3(x)
        x = F.log_softmax(x, dim=1)
        return x

    def training_step(self, batch, batch_idx):
        x, y = batch
        logits = self(x)
        loss = F.nll_loss(logits, y)
        return loss
```

```
model = MNISTClassifier()
x = mnist_image()
logits = model(x)
```

In this case, we've set this `LightningModel` to predict logits. But we could also have it predict feature maps:

```
class MNISTRepresentator(LightningModule):

    def forward(self, x):
        batch_size, channels, width, height = x.size()
        x = x.view(batch_size, -1)
        x = self.layer_1(x)
        x1 = F.relu(x)
        x = self.layer_2(x1)
        x2 = F.relu(x)
        x3 = self.layer_3(x2)
        return [x, x1, x2, x3]
```

(continues on next page)

(continued from previous page)

```
def training_step(self, batch, batch_idx):
    x, y = batch
    out, l1_feats, l2_feats, l3_feats = self(x)
    logits = F.log_softmax(out, dim=1)
    ce_loss = F.nll_loss(logits, y)
    loss = perceptual_loss(l1_feats, l2_feats, l3_feats) + ce_loss
    return loss
```

```
model = MNISTRepresentator.load_from_checkpoint(PATH)
x = mnist_image()
feature_maps = model(x)
```

Or maybe we have a model that we use to do generation

```
class LitMNISTDreamer(LightningModule):

    def forward(self, z):
        imgs = self.decoder(z)
        return imgs

    def training_step(self, batch, batch_idx):
        x, y = batch
        representation = self.encoder(x)
        imgs = self(representation)

        loss = perceptual_loss(imgs, x)
        return loss
```

```
model = LitMNISTDreamer.load_from_checkpoint(PATH)
z = sample_noise()
generated_imgs = model(z)
```

How you split up what goes in `forward` vs `training_step` depends on how you want to use this model for prediction.

12.1.4 The nonessentials

Extensibility

Although lightning makes everything super simple, it doesn't sacrifice any flexibility or control. Lightning offers multiple ways of managing the training state.

Training overrides

Any part of the training, validation, and testing loop can be modified. For instance, if you wanted to do your own backward pass, you would override the default implementation

```
def backward(self, use_amp, loss, optimizer):
    loss.backward()
```

With your own

```
class LitMNIST(LightningModule):

    def backward(self, use_amp, loss, optimizer, optimizer_idx):
        # do a custom way of backward
        loss.backward(retain_graph=True)
```

Every single part of training is configurable this way. For a full list look at [LightningModule](#).

Callbacks

Another way to add arbitrary functionality is to add a custom callback for hooks that you might care about

```
from pytorch_lightning.callbacks import Callback

class MyPrintingCallback(Callback):

    def on_init_start(self, trainer):
        print('Starting to init trainer!')

    def on_init_end(self, trainer):
        print('Trainer is init now')

    def on_train_end(self, trainer, pl_module):
        print('do something when training ends')
```

And pass the callbacks into the trainer

```
trainer = Trainer(callbacks=[MyPrintingCallback()])
```

Tip: See full list of 12+ hooks in the [Callback](#).

Child Modules

Research projects tend to test different approaches to the same dataset. This is very easy to do in Lightning with inheritance.

For example, imagine we now want to train an Autoencoder to use as a feature extractor for MNIST images. We are extending our Autoencoder from the *LitMNIST*-module which already defines all the dataloading. The only things that change in the *Autoencoder* model are the init, forward, training, validation and test step.

```
class Encoder(torch.nn.Module):
    pass

class Decoder(torch.nn.Module):
    pass

class AutoEncoder(LitMNIST):

    def __init__(self):
        super().__init__()
        self.encoder = Encoder()
        self.decoder = Decoder()
        self.metric = MSE()

    def forward(self, x):
        return self.encoder(x)

    def training_step(self, batch, batch_idx):
        x, _ = batch

        representation = self.encoder(x)
        x_hat = self.decoder(representation)

        loss = self.metric(x, x_hat)
        return loss

    def validation_step(self, batch, batch_idx):
        self._shared_eval(batch, batch_idx, 'val')

    def test_step(self, batch, batch_idx):
        self._shared_eval(batch, batch_idx, 'test')

    def _shared_eval(self, batch, batch_idx, prefix):
        x, _ = batch
        representation = self.encoder(x)
        x_hat = self.decoder(representation)

        loss = self.metric(x, x_hat)
        self.log(f'{prefix}_loss', loss)
```

and we can train this using the same trainer

```
autoencoder = AutoEncoder()
trainer = Trainer()
trainer.fit(autoencoder)
```

And remember that the forward method should define the practical use of a LightningModule. In this case, we want to use the *AutoEncoder* to extract image representations

```
some_images = torch.Tensor(32, 1, 28, 28)
representations = autoencoder(some_images)
```

Transfer Learning

Using Pretrained Models

Sometimes we want to use a `LightningModule` as a pretrained model. This is fine because a `LightningModule` is just a `torch.nn.Module`!

Note: Remember that a `LightningModule` is EXACTLY a `torch.nn.Module` but with more capabilities.

Let's use the *AutoEncoder* as a feature extractor in a separate model.

```
class Encoder(torch.nn.Module):
    ...

class AutoEncoder(LightningModule):
    def __init__(self):
        self.encoder = Encoder()
        self.decoder = Decoder()

class CIFAR10Classifier(LightningModule):
    def __init__(self):
        # init the pretrained LightningModule
        self.feature_extractor = AutoEncoder.load_from_checkpoint(PATH)
        self.feature_extractor.freeze()

        # the autoencoder outputs a 100-dim representation and CIFAR-10 has 10 classes
        self.classifier = nn.Linear(100, 10)

    def forward(self, x):
        representations = self.feature_extractor(x)
        x = self.classifier(representations)
        ...
```

We used our pretrained Autoencoder (a `LightningModule`) for transfer learning!

Example: Imagenet (computer Vision)

```
import torchvision.models as models

class ImagenetTransferLearning(LightningModule):
    def __init__(self):
        # init a pretrained resnet
        num_target_classes = 10
        self.feature_extractor = models.resnet50(pretrained=True)
        self.feature_extractor.eval()

        # use the pretrained model to classify cifar-10 (10 image classes)
```

(continues on next page)

(continued from previous page)

```

self.classifier = nn.Linear(2048, num_target_classes)

def forward(self, x):
    representations = self.feature_extractor(x)
    x = self.classifier(representations)
    ...

```

Finetune

```

model = ImagenetTransferLearning()
trainer = Trainer()
trainer.fit(model)

```

And use it to predict your data of interest

```

model = ImagenetTransferLearning.load_from_checkpoint(PATH)
model.freeze()

x = some_images_from_cifar10()
predictions = model(x)

```

We used a pretrained model on imagenet, finetuned on CIFAR-10 to predict on CIFAR-10. In the non-academic world we would finetune on a tiny dataset you have and predict on your dataset.

Example: BERT (NLP)

Lightning is completely agnostic to what's used for transfer learning so long as it is a *torch.nn.Module* subclass.

Here's a model that uses [Huggingface transformers](#).

```

class BertMNLIFinetuner(LightningModule):

    def __init__(self):
        super().__init__()

        self.bert = BertModel.from_pretrained('bert-base-cased', output_
↪attentions=True)
        self.W = nn.Linear(bert.config.hidden_size, 3)
        self.num_classes = 3

    def forward(self, input_ids, attention_mask, token_type_ids):

        h, _, attn = self.bert(input_ids=input_ids,
                                attention_mask=attention_mask,
                                token_type_ids=token_type_ids)

        h_cls = h[:, 0]
        logits = self.W(h_cls)
        return logits, attn

```

12.2 Why PyTorch Lightning

12.2.1 a. Less boilerplate

Research and production code starts with simple code, but quickly grows in complexity once you add GPU training, 16-bit, checkpointing, logging, etc...

PyTorch Lightning implements these features for you and tests them rigorously to make sure you can instead focus on the research idea.

Writing less engineering/boilerplate code means:

- fewer bugs
- faster iteration
- faster prototyping

12.2.2 b. More functionality

In PyTorch Lightning you leverage code written by hundreds of AI researchers, research engs and PhDs from the world's top AI labs, implementing all the latest best practices and SOTA features such as

- GPU, Multi GPU, TPU training
- Multi-node training
- Auto logging
- ...
- Gradient accumulation

12.2.3 c. Less error-prone

Why re-invent the wheel?

Use PyTorch Lightning to enjoy a deep learning structure that is rigorously tested (500+ tests) across CPUs/multi-GPUs/multi-TPUs on every pull-request.

We promise our collective team of 20+ from the top labs has thought about training more than you :)

12.2.4 d. Not a new library

PyTorch Lightning is organized PyTorch - no need to learn a new framework.

Switching your model to Lightning is straight forward - here's a 2-minute video on how to do it.

Your projects WILL grow in complexity and you WILL end up engineering more than trying out new ideas... Defer the hardest parts to Lightning!

12.3 Lightning Philosophy

Lightning structures your deep learning code in 4 parts:

- Research code
- Engineering code
- Non-essential code
- Data code

12.3.1 Research code

In the MNIST generation example, the research code would be the particular system and how it's trained (ie: A GAN or VAE or GPT).

```
l1 = nn.Linear(...)
l2 = nn.Linear(...)
decoder = Decoder()

x1 = l1(x)
x2 = l2(x2)
out = decoder(features, x)

loss = perceptual_loss(x1, x2, x) + CE(out, x)
```

In Lightning, this code is organized into a *LightningModule*.

12.3.2 Engineering code

The Engineering code is all the code related to training this system. Things such as early stopping, distribution over GPUs, 16-bit precision, etc. This is normally code that is THE SAME across most projects.

```
model.cuda(0)
x = x.cuda(0)

distributed = DistributedParallel(model)

with gpu_zero:
    download_data()

dist.barrier()
```

In Lightning, this code is abstracted out by the *Trainer*.

12.3.3 Non-essential code

This is code that helps the research but isn't relevant to the research code. Some examples might be:

1. Inspect gradients
2. Log to tensorboard.

```
# log samples
z = Q.rsample()
generated = decoder(z)
self.experiment.log('images', generated)
```

In Lightning this code is organized into *Callback*.

12.3.4 Data code

Lightning uses standard PyTorch DataLoaders or anything that gives a batch of data. This code tends to end up getting messy with transforms, normalization constants, and data splitting spread all over files.

```
# data
train = MNIST(...)
train, val = split(train, val)
test = MNIST(...)

# transforms
train_transforms = ...
val_transforms = ...
test_transforms = ...

# dataloader ...
# download with dist.barrier() for multi-gpu, etc...
```

This code gets especially complicated once you start doing multi-GPU training or needing info about the data to build your models.

In Lightning this code is organized inside a *LightningDataModule*.

Tip: DataModules are optional but encouraged, otherwise you can use standard DataLoaders

BOLTS

PyTorch Lightning Bolts, is our official collection of prebuilt models across many research domains.

```
pip install pytorch-lightning-bolts
```

In bolts we have:

- A collection of pretrained state-of-the-art models.
- A collection of models designed to bootstrap your research.
- A collection of callbacks, transforms, full datasets.
- All models work on CPUs, TPUs, GPUs and 16-bit precision.

13.1 Quality control

The Lightning community builds bolts and contributes them to Bolts. The lightning team guarantees that contributions are:

- Rigorously Tested (CPUs, GPUs, TPUs).
- Rigorously Documented.
- Standardized via PyTorch Lightning.
- Optimized for speed.
- Checked for correctness.

13.2 Example 1: Pretrained, prebuilt models

```
from pl_bolts.models import VAE, GPT2, ImageGPT, PixelCNN
from pl_bolts.models.self_supervised import AMDIM, CPCV2, SimCLR, MocoV2
from pl_bolts.models import LinearRegression, LogisticRegression
from pl_bolts.models.gans import GAN
from pl_bolts.callbacks import PrintTableMetricsCallback
from pl_bolts.datamodules import FashionMNISTDataModule, CIFAR10DataModule,   
↳ ImagenetDataModule
```

13.3 Example 2: Extend for faster research

Bolts are contributed with benchmarks and continuous-integration tests. This means you can trust the implementations and use them to bootstrap your research much faster.

```
from pl_bolts.models import ImageGPT
from pl_bolts.self_supervised import SimCLR

class VideoGPT(ImageGPT):

    def training_step(self, batch, batch_idx):
        x, y = batch
        x = _shape_input(x)

        logits = self.gpt(x)
        simclr_features = self.simclr(x)

        # -----
        # do something new with GPT logits + simclr_features
        # -----

        loss = self.criterion(logits.view(-1, logits.size(-1)), x.view(-1).long())

        logs = {"loss": loss}
        return {"loss": loss, "log": logs}
```

13.4 Example 3: Callbacks

We also have a collection of callbacks.

```
from pl_bolts.callbacks import PrintTableMetricsCallback
import pytorch_lightning as pl

trainer = pl.Trainer(callbacks=[PrintTableMetricsCallback()])

# loss|train_loss|val_loss|epoch
# -----
# 2.2541470527648926|2.2541470527648926|2.2158432006835938|0
```


COMMUNITY EXAMPLES

- Contextual Emotion Detection (DoubleDistilBert).
- Cotatron: Transcription-Guided Speech Encoder.
- FasterRCNN object detection + Hydra.
- Image Inpainting using Partial Convolutions.
- MNIST on TPU.
- NER (transformers, TPU).
- NeuralTexture (CVPR).
- Recurrent Attentive Neural Process.
- Siamese Nets for One-shot Image Recognition.
- Speech Transformers.
- Transformers transfer learning (Huggingface).
- Transformers text classification.
- VAE Library of over 18+ VAE flavors.
- Transformers Question Answering (SQuAD).
- Atlas: End-to-End 3D Scene Reconstruction from Posed Images.
- Self-Supervised Representation Learning (MoCo and BYOL).

AWS/GCP TRAINING

Lightning has a native solution for training on AWS/GCP at scale (Lightning-Grid). Grid is in private early-access now but you can request access at grid.ai.

We've designed Grid to work for Lightning users without needing to make ANY changes to their code.

To use grid, take your regular command:

```
python my_model.py --learning_rate 1e-6 --layers 2 --gpus 4
```

And change it to use the grid train command:

```
grid train --grid_gpus 4 my_model.py --learning_rate 'uniform(1e-6, 1e-1, 20)' --  
→layers '[2, 4, 8, 16]'
```

The above command will launch (20 * 4) experiments each running on 4 GPUs (320 GPUs!) - by making ZERO changes to your code.

The *uniform* command is part of our new expressive syntax which lets you construct hyperparameter combinations using over 20+ distributions, lists, etc. Of course, you can also configure all of this using yamls which can be dynamically assembled at runtime.

Hint: Grid supports the search strategy of your choice! (and much more than just sweeps)

16-BIT TRAINING

Lightning offers 16-bit training for CPUs, GPUs, and TPUs.

16.1 GPU 16-bit

16-bit precision can cut your memory footprint by half. If using volta architecture GPUs it can give a dramatic training speed-up as well.

Note: PyTorch 1.6+ is recommended for 16-bit

16.1.1 Native torch

When using PyTorch 1.6+ Lightning uses the native amp implementation to support 16-bit.

```
# turn on 16-bit
trainer = Trainer(precision=16)
```

16.1.2 Apex 16-bit

If you are using an earlier version of PyTorch Lightning uses Apex to support 16-bit.

Follow these instructions to install Apex. To use 16-bit precision, do two things:

1. Install Apex
2. Set the “precision” trainer flag.

```
$ git clone https://github.com/NVIDIA/apex
$ cd apex

# -----
# OPTIONAL: on your cluster you might need to load CUDA 10 or 9
```

(continues on next page)

(continued from previous page)

```
# depending on how you installed PyTorch

# see available modules
module avail

# load correct CUDA before install
module load cuda-10.0
# -----

# make sure you've loaded a cuda version > 4.0 and < 7.0
module load gcc-6.1.0

$ pip install -v --no-cache-dir --global-option="--cpp_ext" --global-option="--cuda_
↪ext" ./
```

Warning: NVIDIA Apex and DDP have instability problems. We recommend native 16-bit in PyTorch 1.6+

16.1.3 Enable 16-bit

```
# turn on 16-bit
trainer = Trainer(amp_level='O2', precision=16)
```

If you need to configure the apex init for your particular use case or want to use a different way of doing 16-bit training, override `pytorch_lightning.core.LightningModule.configure_apex()`.

16.2 TPU 16-bit

16-bit on TPUs is much simpler. To use 16-bit with TPUs set precision to 16 when using the TPU flag

```
# DEFAULT
trainer = Trainer(tpu_cores=8, precision=32)

# turn on 16-bit
trainer = Trainer(tpu_cores=8, precision=16)
```

COMPUTING CLUSTER (SLURM)

Lightning automates the details behind training on a SLURM-powered cluster.

17.1 Multi-node training

To train a model using multiple nodes, do the following:

1. Design your *LightningModule*.
2. Enable DDP in the trainer

```
# train on 32 GPUs across 4 nodes
trainer = Trainer(gpus=8, num_nodes=4, distributed_backend='ddp')
```

3. It's a good idea to structure your training script like this:

```
# train.py
def main(hparams):
    model = LightningTemplateModel(hparams)

    trainer = pl.Trainer(
        gpus=8,
        num_nodes=4,
        distributed_backend='ddp'
    )

    trainer.fit(model)

if __name__ == '__main__':
    root_dir = os.path.dirname(os.path.realpath(__file__))
    parent_parser = ArgumentParser(add_help=False)
    hyperparams = parser.parse_args()

    # TRAIN
    main(hyperparams)
```

4. Create the appropriate SLURM job:

```
# (submit.sh)
#!/bin/bash -l
```

(continues on next page)

(continued from previous page)

```
# SLURM SUBMIT SCRIPT
#SBATCH --nodes=4
#SBATCH --gres=gpu:8
#SBATCH --ntasks-per-node=8
#SBATCH --mem=0
#SBATCH --time=0-02:00:00

# activate conda env
source activate $1

# debugging flags (optional)
export NCCL_DEBUG=INFO
export PYTHONFAULTHANDLER=1

# on your cluster you might need these:
# set the network interface
# export NCCL_SOCKET_IFNAME=^docker0,lo

# might need the latest CUDA
# module load NCCL/2.4.7-1-cuda.10.0

# run script from above
srun python3 train.py
```

5. If you want auto-resubmit (read below), add this line to the submit.sh script

```
#SBATCH --signal=SIGUSR1@90
```

6. Submit the SLURM job

```
sbatch submit.sh
```

Note: When running in DDP mode, any errors in your code will show up as an NCCL issue. Set the `NCCL_DEBUG=INFO` flag to see the ACTUAL error.

Normally now you would need to add a `DistributedSampler` to your dataset, however Lightning automates this for you. But if you still need to set a sampler set the Trainer flag `replace_sampler_ddp` to `False`.

Here's an example of how to add your own sampler (again, not needed with Lightning).

```
# in your LightningModule
def train_dataloader(self):
    dataset = MyDataset()
    dist_sampler = torch.utils.data.distributed.DistributedSampler(dataset)
    dataloader = DataLoader(dataset, sampler=dist_sampler)
    return dataloader

# in your training script
trainer = Trainer(replace_sampler_ddp=False)
```


17.2 Wall time auto-resubmit

When you use Lightning in a SLURM cluster, it automatically detects when it is about to run into the wall time and does the following:

1. Saves a temporary checkpoint.
2. Requeues the job.
3. When the job starts, it loads the temporary checkpoint.

To get this behavior make sure to add the correct signal to your SLURM script

```
# 90 seconds before training ends
SBATCH --signal=SIGUSR1@90
```

17.3 Building SLURM scripts

Instead of manually building SLURM scripts, you can use the `SlurmCluster` object to do this for you. The `SlurmCluster` can also run a grid search if you pass in a `HyperOptArgumentParser`.

Here is an example where you run a grid search of 9 combinations of hyperparameters. See also the multi-node examples [here](#).

```
# grid search 3 values of learning rate and 3 values of number of layers for your net
# this generates 9 experiments (lr=1e-3, layers=16), (lr=1e-3, layers=32),
# (lr=1e-3, layers=64), ... (lr=1e-1, layers=64)
parser = HyperOptArgumentParser(strategy='grid_search', add_help=False)
parser.opt_list('--learning_rate', default=0.001, type=float,
                options=[1e-3, 1e-2, 1e-1], tunable=True)
parser.opt_list('--layers', default=1, type=float, options=[16, 32, 64], tunable=True)
hyperparams = parser.parse_args()

# Slurm cluster submits 9 jobs, each with a set of hyperparams
cluster = SlurmCluster(
    hyperparam_optimizer=hyperparams,
    log_path='/some/path/to/save',
)

# OPTIONAL FLAGS WHICH MAY BE CLUSTER DEPENDENT
# which interface your nodes use for communication
cluster.add_command('export NCCL_SOCKET_IFNAME=docker0,lo')

# see the output of the NCCL connection process
# NCCL is how the nodes talk to each other
cluster.add_command('export NCCL_DEBUG=INFO')

# setting a master port here is a good idea.
cluster.add_command('export MASTER_PORT=%r' % PORT)

# ***** DON'T FORGET THIS *****
# MUST load the latest NCCL version
cluster.load_modules(['NCCL/2.4.7-1-cuda.10.0'])
```

(continues on next page)

(continued from previous page)

```
# configure cluster
cluster.per_experiment_nb_nodes = 12
cluster.per_experiment_nb_gpus = 8

cluster.add_slurm_cmd(cmd='ntasks-per-node', value=8, comment='1 task per gpu')

# submit a script with 9 combinations of hyper params
# (lr=1e-3, layers=16), (lr=1e-3, layers=32), (lr=1e-3, layers=64), ... (lr=1e-1,
↳ layers=64)
cluster.optimize_parallel_cluster_gpu(
    main,
    nb_trials=9, # how many permutations of the grid search to run
    job_name='name_for_queue'
)
```

The other option is that you generate scripts on your own via a bash command or use another library.

17.4 Self-balancing architecture (COMING SOON)

Here Lightning distributes parts of your module across available GPUs to optimize for speed and memory.

CHILD MODULES

Research projects tend to test different approaches to the same dataset. This is very easy to do in Lightning with inheritance.

For example, imagine we now want to train an Autoencoder to use as a feature extractor for MNIST images. We are extending our Autoencoder from the *LitMNIST*-module which already defines all the dataloading. The only things that change in the *Autoencoder* model are the init, forward, training, validation and test step.

```
class Encoder(torch.nn.Module):
    pass

class Decoder(torch.nn.Module):
    pass

class AutoEncoder(LitMNIST):

    def __init__(self):
        super().__init__()
        self.encoder = Encoder()
        self.decoder = Decoder()
        self.metric = MSE()

    def forward(self, x):
        return self.encoder(x)

    def training_step(self, batch, batch_idx):
        x, _ = batch

        representation = self.encoder(x)
        x_hat = self.decoder(representation)

        loss = self.metric(x, x_hat)
        return loss

    def validation_step(self, batch, batch_idx):
        self._shared_eval(batch, batch_idx, 'val')

    def test_step(self, batch, batch_idx):
        self._shared_eval(batch, batch_idx, 'test')

    def _shared_eval(self, batch, batch_idx, prefix):
        x, _ = batch
        representation = self.encoder(x)
        x_hat = self.decoder(representation)
```

(continues on next page)

(continued from previous page)

```
loss = self.metric(x, x_hat)
self.log(f'{prefix}_loss', loss)
```

and we can train this using the same trainer

```
autoencoder = AutoEncoder()
trainer = Trainer()
trainer.fit(autoencoder)
```

And remember that the forward method should define the practical use of a `LightningModule`. In this case, we want to use the *AutoEncoder* to extract image representations

```
some_images = torch.Tensor(32, 1, 28, 28)
representations = autoencoder(some_images)
```

DEBUGGING

The following are flags that make debugging much easier.

19.1 fast_dev_run

This flag runs a “unit test” by running 1 training batch and 1 validation batch. The point is to detect any bugs in the training/validation loop without having to wait for a full epoch to crash.

(See: `fast_dev_run` argument of `Trainer`)

```
trainer = Trainer(fast_dev_run=True)
```

19.2 Inspect gradient norms

Logs (to a logger), the norm of each weight matrix.

(See: `track_grad_norm` argument of `Trainer`)

```
# the 2-norm
trainer = Trainer(track_grad_norm=2)
```

19.3 Log GPU usage

Logs (to a logger) the GPU usage for each GPU on the master machine.

(See: `log_gpu_memory` argument of `Trainer`)

```
trainer = Trainer(log_gpu_memory=True)
```

19.4 Make model overfit on subset of data

A good debugging technique is to take a tiny portion of your data (say 2 samples per class), and try to get your model to overfit. If it can't, it's a sign it won't work with large datasets.

(See: `overfit_batches` argument of `Trainer`)

```
# use only 1% of training data (and use the same training dataloader (with shuffle_
↪off) in val and test)
trainer = Trainer(overfit_batches=0.01)

# similar, but with a fixed 10 batches no matter the size of the dataset
trainer = Trainer(overfit_batches=10)
```

With this flag, the train, val, and test sets will all be the same train set. We will also replace the sampler in the training set to turn off shuffle for you.

19.5 Print a summary of your LightningModule

Whenever the `.fit()` function gets called, the `Trainer` will print the weights summary for the `LightningModule`. By default it only prints the top-level modules. If you want to show all submodules in your network, use the `'full'` option:

```
trainer = Trainer(weights_summary='full')
```

You can also display the intermediate input- and output sizes of all your layers by setting the `example_input_array` attribute in your `LightningModule`. It will print a table like this

	Name	Type	Params	In sizes	Out sizes
0	net	Sequential	132 K	[10, 256]	[10, 512]
1	net.0	Linear	131 K	[10, 256]	[10, 512]
2	net.1	BatchNorm1d	1 K	[10, 512]	[10, 512]

when you call `.fit()` on the `Trainer`. This can help you find bugs in the composition of your layers.

See Also:

- `weights_summary` `Trainer` argument
 - `ModelSummary`
-

19.6 Shorten epochs

Sometimes it's helpful to only use a percentage of your training, val or test data (or a set number of batches). For example, you can use 20% of the training set and 1% of the validation set.

On larger datasets like Imagenet, this can help you debug or test a few things faster than waiting for a full epoch.

```
# use only 10% of training data and 1% of val data
trainer = Trainer(limit_train_batches=0.1, limit_val_batches=0.01)

# use 10 batches of train and 5 batches of val
trainer = Trainer(limit_train_batches=10, limit_val_batches=5)
```

19.7 Set the number of validation sanity steps

Lightning runs a few steps of validation in the beginning of training. This avoids crashing in the validation loop sometime deep into a lengthy training loop.

(See: `num_sanity_val_steps` argument of `Trainer`)

```
# DEFAULT
trainer = Trainer(num_sanity_val_steps=2)
```


LOGGERS

Lightning supports the most popular logging frameworks (TensorBoard, Comet, etc...). TensorBoard is used by default, but you can pass to the `Trainer` any combination of the following loggers.

Note: All loggers log by default to `os.getcwd()`. To change the path without creating a logger set `Trainer(default_root_dir='/your/path/to/save/checkpoints')`

Read more about [Logging](#) options.

20.1 Comet.ml

[Comet.ml](#) is a third-party logger. To use `CometLogger` as your logger do the following. First, install the package:

```
pip install comet-ml
```

Then configure the logger and pass it to the `Trainer`:

```
import os
from pytorch_lightning.loggers import CometLogger
comet_logger = CometLogger(
    api_key=os.environ.get('COMET_API_KEY'),
    workspace=os.environ.get('COMET_WORKSPACE'), # Optional
    save_dir='.', # Optional
    project_name='default_project', # Optional
    rest_api_key=os.environ.get('COMET_REST_API_KEY'), # Optional
    experiment_name='default' # Optional
)
trainer = Trainer(logger=comet_logger)
```

The `CometLogger` is available anywhere except `__init__` in your `LightningModule`.

```
class MyModule(LightningModule):
    def any_lightning_module_function_or_hook(self):
        some_img = fake_image()
        self.logger.experiment.add_image('generated_images', some_img, 0)
```

See also:

[CometLogger docs](#).

20.2 MLflow

MLflow is a third-party logger. To use MLFlowLogger as your logger do the following. First, install the package:

```
pip install mlflow
```

Then configure the logger and pass it to the Trainer:

```
from pytorch_lightning.loggers import MLFlowLogger
mlf_logger = MLFlowLogger(
    experiment_name="default",
    tracking_uri="file:./ml-runs"
)
trainer = Trainer(logger=mlf_logger)
```

See also:

MLFlowLogger docs.

20.3 Neptune.ai

Neptune.ai is a third-party logger. To use NeptuneLogger as your logger do the following. First, install the package:

```
pip install neptune-client
```

Then configure the logger and pass it to the Trainer:

```
from pytorch_lightning.loggers import NeptuneLogger

neptune_logger = NeptuneLogger(
    api_key='ANONYMOUS', # replace with your own
    project_name='shared/pytorch-lightning-integration',
    experiment_name='default', # Optional,
    params={'max_epochs': 10}, # Optional,
    tags=['pytorch-lightning', 'mlp'], # Optional,
)
trainer = Trainer(logger=neptune_logger)
```

The NeptuneLogger is available anywhere except `__init__` in your LightningModule.

```
class MyModule(LightningModule):
    def any_lightning_module_function_or_hook(self):
        some_img = fake_image()
        self.logger.experiment.add_image('generated_images', some_img, 0)
```

See also:

NeptuneLogger docs.

20.4 Tensorboard

To use `TensorBoard` as your logger do the following.

```
from pytorch_lightning.loggers import TensorBoardLogger
logger = TensorBoardLogger('tb_logs', name='my_model')
trainer = Trainer(logger=logger)
```

The `TensorBoardLogger` is available anywhere except `__init__` in your `LightningModule`.

```
class MyModule(LightningModule):
    def any_lightning_module_function_or_hook(self):
        some_img = fake_image()
        self.logger.experiment.add_image('generated_images', some_img, 0)
```

See also:

`TensorBoardLogger docs`.

20.5 Test Tube

`Test Tube` is a `TensorBoard` logger but with nicer file structure. To use `TestTubeLogger` as your logger do the following. First, install the package:

```
pip install test_tube
```

Then configure the logger and pass it to the `Trainer`:

```
from pytorch_lightning.loggers import TestTubeLogger
logger = TestTubeLogger('tb_logs', name='my_model')
trainer = Trainer(logger=logger)
```

The `TestTubeLogger` is available anywhere except `__init__` in your `LightningModule`.

```
class MyModule(LightningModule):
    def any_lightning_module_function_or_hook(self):
        some_img = fake_image()
        self.logger.experiment.add_image('generated_images', some_img, 0)
```

See also:

`TestTubeLogger docs`.

20.6 Weights and Biases

Weights and Biases is a third-party logger. To use WandbLogger as your logger do the following. First, install the package:

```
pip install wandb
```

Then configure the logger and pass it to the Trainer:

```
from pytorch_lightning.loggers import WandbLogger
wandb_logger = WandbLogger(offline=True)
trainer = Trainer(logger=wandb_logger)
```

The WandbLogger is available anywhere except `__init__` in your LightningModule.

```
class MyModule(LightningModule):
    def any_lightning_module_function_or_hook(self):
        some_img = fake_image()
        self.logger.experiment.log({
            "generated_images": [wandb.Image(some_img, caption="...")]
        })
```

See also:

WandbLogger docs.

20.7 Multiple Loggers

Lightning supports the use of multiple loggers, just pass a list to the Trainer.

```
from pytorch_lightning.loggers import TensorBoardLogger, TestTubeLogger
logger1 = TensorBoardLogger('tb_logs', name='my_model')
logger2 = TestTubeLogger('tb_logs', name='my_model')
trainer = Trainer(logger=[logger1, logger2])
```

The loggers are available as a list anywhere except `__init__` in your LightningModule.

```
class MyModule(LightningModule):
    def any_lightning_module_function_or_hook(self):
        some_img = fake_image()
        # Option 1
        self.logger.experiment[0].add_image('generated_images', some_img, 0)
        # Option 2
        self.logger[0].experiment.add_image('generated_images', some_img, 0)
```

EARLY STOPPING

21.1 Stopping an epoch early

You can stop an epoch early by overriding `on_train_batch_start()` to return `-1` when some condition is met. If you do this repeatedly, for every epoch you had originally requested, then this will stop your entire run.

21.2 Early stopping based on metric using the EarlyStopping Callback

The `EarlyStopping` callback can be used to monitor a validation metric and stop the training when no improvement is observed.

To enable it:

- Import `EarlyStopping` callback.
- Log the metric you want to monitor using `log()` method.
- Init the callback, and set *monitor* to the logged metric of your choice.
- Pass the `EarlyStopping` callback to the `Trainer` callbacks flag.

```
from pytorch_lightning.callbacks.early_stopping import EarlyStopping

def validation_step(...):
    self.log('val_loss', loss)

trainer = Trainer(callbacks=[EarlyStopping(monitor='val_loss')])
```

- You can customize the callbacks behaviour by changing its parameters.

```
early_stop_callback = EarlyStopping(
    monitor='val_accuracy',
    min_delta=0.00,
    patience=3,
    verbose=False,
    mode='max'
```

(continues on next page)

(continued from previous page)

```
)  
trainer = Trainer(callbacks=[early_stop_callback])
```

In case you need early stopping in a different part of training, subclass `EarlyStopping` and change where it is called:

```
class MyEarlyStopping(EarlyStopping):  
  
    def on_validation_end(self, trainer, pl_module):  
        # override this to disable early stopping at the end of val loop  
        pass  
  
    def on_train_end(self, trainer, pl_module):  
        # instead, do it at the end of training loop  
        self._run_early_stopping_check(trainer, pl_module)
```

Note: The `EarlyStopping` callback runs at the end of every validation epoch, which, under the default configuration, happen after every training epoch. However, the frequency of validation can be modified by setting various parameters in the `Trainer`, for example `check_val_every_n_epoch` and `val_check_interval`. It must be noted that the *patience* parameter counts the number of validation epochs with no improvement, and not the number of training epochs. Therefore, with parameters `check_val_every_n_epoch=10` and `patience=3`, the trainer will perform at least 40 training epochs before being stopped.

See also:

- `Trainer`
 - `EarlyStopping`
-

See also:

- `Trainer`
- `EarlyStopping`

FAST TRAINING

There are multiple options to speed up different parts of the training by choosing to train on a subset of data. This could be done for speed or debugging purposes.

22.1 Check validation every n epochs

If you have a small dataset you might want to check validation every n epochs

```
# DEFAULT
trainer = Trainer(check_val_every_n_epoch=1)
```

22.2 Force training for min or max epochs

It can be useful to force training for a minimum number of epochs or limit to a max number.

See also:

Trainer

```
# DEFAULT
trainer = Trainer(min_epochs=1, max_epochs=1000)
```

22.3 Set validation check frequency within 1 training epoch

For large datasets it's often desirable to check validation multiple times within a training loop. Pass in a float to check that often within 1 training epoch. Pass in an int *k* to check every *k* training batches. Must use an *int* if using an *IterableDataset*.

```
# DEFAULT
trainer = Trainer(val_check_interval=0.95)

# check every .25 of an epoch
trainer = Trainer(val_check_interval=0.25)
```

(continues on next page)

(continued from previous page)

```
# check every 100 train batches (ie: for `IterableDatasets` or fixed frequency)
trainer = Trainer(val_check_interval=100)
```

22.4 Use data subset for training, validation, and test

If you don't want to check 100% of the training/validation/test set (for debugging or if it's huge), set these flags.

```
# DEFAULT
trainer = Trainer(
    limit_train_batches=1.0,
    limit_val_batches=1.0,
    limit_test_batches=1.0
)

# check 10%, 20%, 30% only, respectively for training, validation and test set
trainer = Trainer(
    limit_train_batches=0.1,
    limit_val_batches=0.2,
    limit_test_batches=0.3
)
```

If you also pass `shuffle=True` to the dataloader, a different random subset of your dataset will be used for each epoch; otherwise the same subset will be used for all epochs.

Note: `limit_train_batches`, `limit_val_batches` and `limit_test_batches` will be overwritten by `overfit_batches` if `overfit_batches > 0`. `limit_val_batches` will be ignored if `fast_dev_run=True`.

Note: If you set `limit_val_batches=0`, validation will be disabled.

HYPERPARAMETERS

Lightning has utilities to interact seamlessly with the command line `ArgumentParser` and plays well with the hyperparameter optimization framework of your choice.

23.1 ArgumentParser

Lightning is designed to augment a lot of the functionality of the built-in Python `ArgumentParser`

```
from argparse import ArgumentParser
parser = ArgumentParser()
parser.add_argument('--layer_1_dim', type=int, default=128)
args = parser.parse_args()
```

This allows you to call your program like so:

```
python trainer.py --layer_1_dim 64
```

23.2 Argparser Best Practices

It is best practice to layer your arguments in three sections.

1. Trainer args (gpus, num_nodes, etc...)
2. Model specific arguments (layer_dim, num_layers, learning_rate, etc...)
3. Program arguments (data_path, cluster_email, etc...)

We can do this as follows. First, in your `LightningModule`, define the arguments specific to that module. Remember that data splits or data paths may also be specific to a module (i.e.: if your project has a model that trains on Imagenet and another on CIFAR-10).

```
class LitModel(LightningModule):

    @staticmethod
    def add_model_specific_args(parent_parser):
        parser = ArgumentParser(parents=[parent_parser], add_help=False)
        parser.add_argument('--encoder_layers', type=int, default=12)
        parser.add_argument('--data_path', type=str, default='/some/path')
        return parser
```

Now in your main trainer file, add the Trainer args, the program args, and add the model args

```
# -----
# trainer_main.py
# -----
from argparse import ArgumentParser
parser = ArgumentParser()

# add PROGRAM level args
parser.add_argument('--conda_env', type=str, default='some_name')
parser.add_argument('--notification_email', type=str, default='will@email.com')

# add model specific args
parser = LitModel.add_model_specific_args(parser)

# add all the available trainer options to argparse
# ie: now --gpus --num_nodes ... --fast_dev_run all work in the cli
parser = Trainer.add_argparse_args(parser)

args = parser.parse_args()
```

Now you can call run your program like so:

```
python trainer_main.py --gpus 2 --num_nodes 2 --conda_env 'my_env' --encoder_layers 12
```

Finally, make sure to start the training like so:

```
# init the trainer like this
trainer = Trainer.from_argparse_args(args, early_stopping_callback=...)

# NOT like this
trainer = Trainer(gpus=hparams.gpus, ...)

# init the model with Namespace directly
model = LitModel(args)

# or init the model with all the key-value pairs
dict_args = vars(args)
model = LitModel(**dict_args)
```

23.3 LightningModule hyperparameters

Often times we train many versions of a model. You might share that model or come back to it a few months later at which point it is very useful to know how that model was trained (i.e.: what learning rate, neural network, etc...).

Lightning has a few ways of saving that information for you in checkpoints and yaml files. The goal here is to improve readability and reproducibility

1. The first way is to ask lightning to save the values of anything in the `__init__` for you to the checkpoint. This also makes those values available via `self.hparams`.

```
class LitMNIST(LightningModule):

    def __init__(self, layer_1_dim=128, learning_rate=1e-2, **kwargs):
        super().__init__()
        # call this to save (layer_1_dim=128, learning_rate=1e-4) to the checkpoint
        self.save_hyperparameters()

        # equivalent
        self.save_hyperparameters('layer_1_dim', 'learning_rate')

        # this now works
        self.hparams.layer_1_dim
```

2. Sometimes your init might have objects or other parameters you might not want to save. In that case, choose only a few

```
class LitMNIST(LightningModule):

    def __init__(self, loss_fx, generator_network, layer_1_dim=128 **kwargs):
        super().__init__()
        self.layer_1_dim = layer_1_dim
        self.loss_fx = loss_fx

        # call this to save (layer_1_dim=128) to the checkpoint
        self.save_hyperparameters('layer_1_dim')

# to load specify the other args
model = LitMNIST.load_from_checkpoint(PATH, loss_fx=torch.nn.SomeOtherLoss, generator_
↪network=MyGenerator())
```

3. Assign to `self.hparams`. Anything assigned to `self.hparams` will also be saved automatically

```
# using a argparse.Namespace
class LitMNIST(LightningModule):

    def __init__(self, hparams, *args, **kwargs):
        super().__init__()
        self.hparams = hparams

        self.layer_1 = torch.nn.Linear(28 * 28, self.hparams.layer_1_dim)
        self.layer_2 = torch.nn.Linear(self.hparams.layer_1_dim, self.hparams.layer_2_
↪dim)
        self.layer_3 = torch.nn.Linear(self.hparams.layer_2_dim, 10)

    def train_dataloader(self):
        return DataLoader(mnist_train, batch_size=self.hparams.batch_size)
```

4. You can also save full objects such as *dict* or *Namespace* to the checkpoint.

```
# using a argparse.Namespace
class LitMNIST(LightningModule):

    def __init__(self, conf, *args, **kwargs):
        super().__init__()
        self.hparams = conf

        # equivalent
        self.save_hyperparameters(conf)

        self.layer_1 = torch.nn.Linear(28 * 28, self.hparams.layer_1_dim)
        self.layer_2 = torch.nn.Linear(self.hparams.layer_1_dim, self.hparams.layer_2_
→ dim)
        self.layer_3 = torch.nn.Linear(self.hparams.layer_2_dim, 10)

conf = OmegaConf.create(...)
model = LitMNIST(conf)

# this works
model.hparams.anything
```

23.4 Trainer args

To recap, add ALL possible trainer flags to the argparser and init the Trainer this way

```
parser = ArgumentParser()
parser = Trainer.add_argparse_args(parser)
hparams = parser.parse_args()

trainer = Trainer.from_argparse_args(hparams)

# or if you need to pass in callbacks
trainer = Trainer.from_argparse_args(hparams, checkpoint_callback=..., callbacks=[...
→ ])
```

23.5 Multiple Lightning Modules

We often have multiple Lightning Modules where each one has different arguments. Instead of polluting the `main.py` file, the `LightningModule` lets you define arguments for each one.

```
class LitMNIST(LightningModule):

    def __init__(self, layer_1_dim, **kwargs):
        super().__init__()
        self.layer_1 = torch.nn.Linear(28 * 28, layer_1_dim)

    @staticmethod
    def add_model_specific_args(parent_parser):
```

(continues on next page)

(continued from previous page)

```

parser = ArgumentParser(parents=[parent_parser], add_help=False)
parser.add_argument('--layer_1_dim', type=int, default=128)
return parser

```

```

class GoodGAN(LightningModule):

    def __init__(self, encoder_layers, **kwargs):
        super().__init__()
        self.encoder = Encoder(layers=encoder_layers)

    @staticmethod
    def add_model_specific_args(parent_parser):
        parser = ArgumentParser(parents=[parent_parser], add_help=False)
        parser.add_argument('--encoder_layers', type=int, default=12)
        return parser

```

Now we can allow each model to inject the arguments it needs in the `main.py`

```

def main(args):
    dict_args = vars(args)

    # pick model
    if args.model_name == 'gan':
        model = GoodGAN(**dict_args)
    elif args.model_name == 'mnist':
        model = LitMNIST(**dict_args)

    trainer = Trainer.from_argparse_args(args)
    trainer.fit(model)

if __name__ == '__main__':
    parser = ArgumentParser()
    parser = Trainer.add_argparse_args(parser)

    # figure out which model to use
    parser.add_argument('--model_name', type=str, default='gan', help='gan or mnist')

    # THIS LINE IS KEY TO PULL THE MODEL NAME
    temp_args, _ = parser.parse_known_args()

    # let the model add what it wants
    if temp_args.model_name == 'gan':
        parser = GoodGAN.add_model_specific_args(parser)
    elif temp_args.model_name == 'mnist':
        parser = LitMNIST.add_model_specific_args(parser)

    args = parser.parse_args()

    # train
    main(args)

```

and now we can train MNIST or the GAN using the command line interface!

```

$ python main.py --model_name gan --encoder_layers 24
$ python main.py --model_name mnist --layer_1_dim 128

```


LEARNING RATE FINDER

For training deep neural networks, selecting a good learning rate is essential for both better performance and faster convergence. Even optimizers such as *Adam* that are self-adjusting the learning rate can benefit from more optimal choices.

To reduce the amount of guesswork concerning choosing a good initial learning rate, a *learning rate finder* can be used. As described in this [paper](#) a learning rate finder does a small run where the learning rate is increased after each processed batch and the corresponding loss is logged. The result of this is a *lr* vs. *loss* plot that can be used as guidance for choosing a optimal initial lr.

Warning: For the moment, this feature only works with models having a single optimizer. LR Finder support for DDP is not implemented yet, it is coming soon.

24.1 Using Lightning's built-in LR finder

To enable the learning rate finder, your *LightningModule* needs to have a `learning_rate` or `lr` property. Then, set `Trainer(auto_lr_find=True)` during trainer construction, and then call `trainer.tune(model)` to run the LR finder. The suggested `learning_rate` will be written to the console and will be automatically set to your *LightningModule*, which can be accessed via `self.learning_rate` or `self.lr`.

```
class LitModel(LightningModule):

    def __init__(self, learning_rate):
        self.learning_rate = learning_rate

    def configure_optimizers(self):
        return Adam(self.parameters(), lr=(self.lr or self.learning_rate))

model = LitModel()

# finds learning rate automatically
# sets hparams.lr or hparams.learning_rate to that learning rate
trainer = Trainer(auto_lr_find=True)

trainer.tune(model)
```

If your model is using an arbitrary value instead of `self.lr` or `self.learning_rate`, set that value as `auto_lr_find`:

```
model = LitModel()

# to set to your own hparams.my_value
trainer = Trainer(auto_lr_find='my_value')

trainer.tune(model)
```

If you want to inspect the results of the learning rate finder or just play around with the parameters of the algorithm, this can be done by invoking the `lr_find` method of the trainer. A typical example of this would look like

```
model = MyModelClass(hparams)
trainer = Trainer()

# Run learning rate finder
lr_finder = trainer.tuner.lr_find(model)

# Results can be found in
lr_finder.results

# Plot with
fig = lr_finder.plot(suggest=True)
fig.show()

# Pick point based on plot, or get suggestion
new_lr = lr_finder.suggestion()

# update hparams of the model
model.hparams.lr = new_lr

# Fit model
trainer.fit(model)
```

The figure produced by `lr_finder.plot()` should look something like the figure below. It is recommended to not pick the learning rate that achieves the lowest loss, but instead something in the middle of the sharpest downward slope (red point). This is the point returned by `lr_finder.suggestion()`.

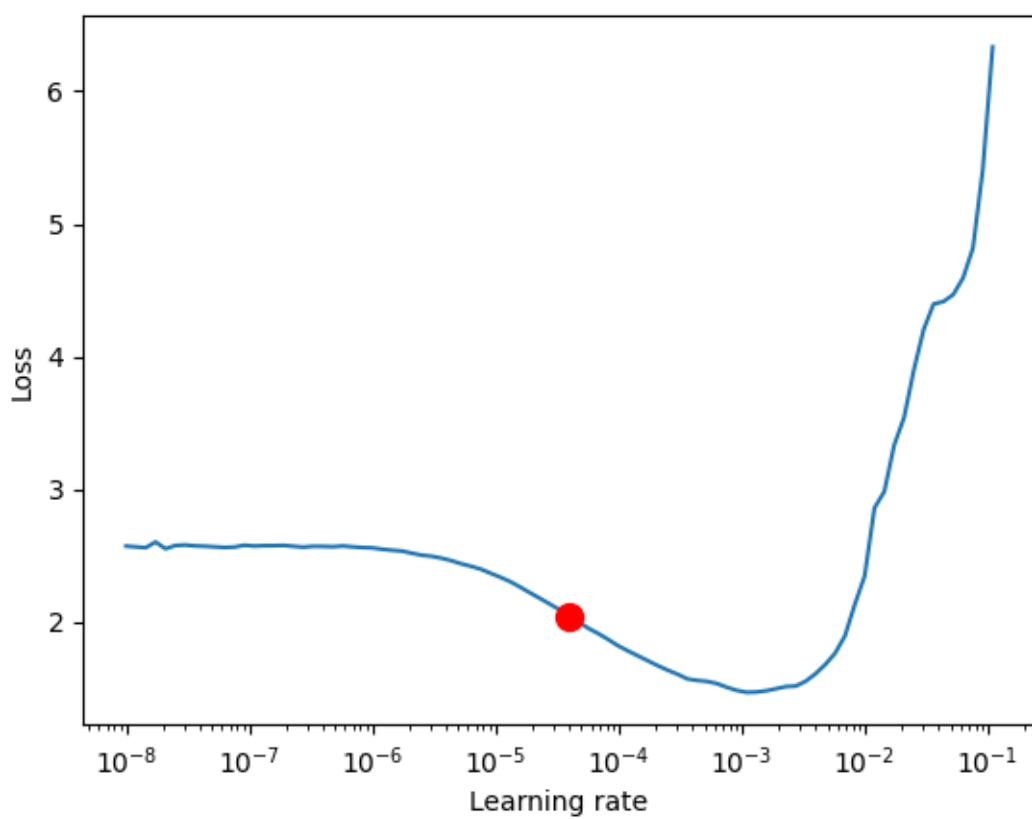
The parameters of the algorithm can be seen below.

```
pytorch_lightning.tuner.lr_finder.lr_find(trainer, model, train_dataloader=None,
                                           val_dataloaders=None, min_lr=1e-08,
                                           max_lr=1, num_training=100,
                                           mode='exponential', early_stop_threshold=4.0,
                                           datamodule=None)
```

`lr_find` enables the user to do a range test of good initial learning rates, to reduce the amount of guesswork in picking a good starting learning rate.

Parameters

- **model** `(LightningModule)` – Model to do range testing for
- **train_dataloader** `(Optional[DataLoader])` – A PyTorch *DataLoader* with training samples. If the model has a predefined `train_dataloader` method, this will be skipped.
- **min_lr** `(float)` – minimum learning rate to investigate
- **max_lr** `(float)` – maximum learning rate to investigate



- **num_training** (int) – number of learning rates to test
- **mode** (str) – search strategy, either ‘linear’ or ‘exponential’. If set to ‘linear’ the learning rate will be searched by linearly increasing after each batch. If set to ‘exponential’, will increase learning rate exponentially.
- **early_stop_threshold** (float) – threshold for stopping the search. If the loss at any point is larger than `early_stop_threshold*best_loss` then the search is stopped. To disable, set to None.
- **datamodule** (Optional[LightningDataModule]) – An optional *LightningDataModule* which holds the training and validation dataloader(s). Note that the `train_dataloader` and `val_dataloaders` parameters cannot be used at the same time as this parameter, or a *MisconfigurationException* will be raised.

Example:

```
# Setup model and trainer
model = MyModelClass(hparams)
trainer = pl.Trainer()

# Run lr finder
lr_finder = trainer.tuner.lr_find(model, ...)

# Inspect results
fig = lr_finder.plot(); fig.show()
suggested_lr = lr_finder.suggestion()

# Overwrite lr and create new model
hparams.lr = suggested_lr
model = MyModelClass(hparams)

# Ready to train with new learning rate
trainer.fit(model)
```

MULTI-GPU TRAINING

Lightning supports multiple ways of doing distributed training.

25.1 Preparing your code

To train on CPU/GPU/TPU without changing your code, we need to build a few good habits :)

25.1.1 Delete `.cuda()` or `.to()` calls

Delete any calls to `.cuda()` or `.to(device)`.

```
# before lightning
def forward(self, x):
    x = x.cuda(0)
    layer_1.cuda(0)
    x_hat = layer_1(x)

# after lightning
def forward(self, x):
    x_hat = layer_1(x)
```

25.1.2 Init tensors using `type_as` and `register_buffer`

When you need to create a new tensor, use `type_as`. This will make your code scale to any arbitrary number of GPUs or TPUs with Lightning.

```
# before lightning
def forward(self, x):
    z = torch.Tensor(2, 3)
    z = z.cuda(0)

# with lightning
def forward(self, x):
```

(continues on next page)

(continued from previous page)

```
z = torch.Tensor(2, 3)
z = z.type_as(x)
```

The `LightningModule` knows what device it is on. You can access the reference via `self.device`. Sometimes it is necessary to store tensors as module attributes. However, if they are not parameters they will remain on the CPU even if the module gets moved to a new device. To prevent that and remain device agnostic, register the tensor as a buffer in your modules's `__init__` method with `register_buffer()`.

```
class LitModel(LightningModule):

    def __init__(self):
        ...
        self.register_buffer("sigma", torch.eye(3))
        # you can now access self.sigma anywhere in your module
```

25.1.3 Remove samplers

In PyTorch, you must use `torch.nn.DistributedSampler` for multi-node or TPU training. The sampler makes sure each GPU sees the appropriate part of your data.

```
# without lightning
def train_dataloader(self):
    dataset = MNIST(...)
    sampler = None

    if self.on_tpu:
        sampler = DistributedSampler(dataset)

    return DataLoader(dataset, sampler=sampler)
```

Lightning adds the correct samplers when needed, so no need to explicitly add samplers.

```
# with lightning
def train_dataloader(self):
    dataset = MNIST(...)
    return DataLoader(dataset)
```

Note: You can disable this behavior with `Trainer(replace_sampler_ddp=False)`

Note: For iterable datasets, we don't do this automatically.

25.1.4 Make models pickleable

It's very likely your code is already `pickleable`, in that case no change is necessary. However, if you run a distributed model and get the following error:

```
self._launch(process_obj)
File "/net/software/local/python/3.6.5/lib/python3.6/multiprocessing/popen_spawn_
↳posix.py", line 47,
in _launch reduction.dump(process_obj, fp)
File "/net/software/local/python/3.6.5/lib/python3.6/multiprocessing/reduction.py",
↳line 60, in dump
ForkingPickler(file, protocol).dump(obj)
_pickle.PicklingError: Can't pickle <function <lambda> at 0x2b599e088ae8>:
attribute lookup <lambda> on __main__ failed
```

This means something in your model definition, transforms, optimizer, dataloader or callbacks cannot be pickled, and the following code will fail:

```
import pickle
pickle.dump(some_object)
```

This is a limitation of using multiple processes for distributed training within PyTorch. To fix this issue, find your piece of code that cannot be pickled. The end of the stacktrace is usually helpful. ie: in the stacktrace example here, there seems to be a lambda function somewhere in the code which cannot be pickled.

```
self._launch(process_obj)
File "/net/software/local/python/3.6.5/lib/python3.6/multiprocessing/popen_spawn_
↳posix.py", line 47,
in _launch reduction.dump(process_obj, fp)
File "/net/software/local/python/3.6.5/lib/python3.6/multiprocessing/reduction.py",
↳line 60, in dump
ForkingPickler(file, protocol).dump(obj)
_pickle.PicklingError: Can't pickle [THIS IS THE THING TO FIND AND DELETE]:
attribute lookup <lambda> on __main__ failed
```

25.2 Select GPU devices

You can select the GPU devices using ranges, a list of indices or a string containing a comma separated list of GPU ids:

```
# DEFAULT (int) specifies how many GPUs to use per node
Trainer(gpus=k)

# Above is equivalent to
Trainer(gpus=list(range(k)))

# Specify which GPUs to use (don't use when running on cluster)
Trainer(gpus=[0, 1])

# Equivalent using a string
Trainer(gpus='0, 1')

# To use all available GPUs put -1 or '-1'
```

(continues on next page)

(continued from previous page)

```
# equivalent to list(range(torch.cuda.available_devices()))
Trainer(gpus=-1)
```

The table below lists examples of possible input formats and how they are interpreted by Lightning. Note in particular the difference between `gpus=0`, `gpus=[0]` and `gpus="0"`.

<i>gpus</i>	Type	Parsed	Meaning
None	NoneType	None	CPU
0	int	None	CPU
3	int	[0, 1, 2]	first 3 GPUs
-1	int	[0, 1, 2, ...]	all available GPUs
[0]	list	[0]	GPU 0
[1, 3]	list	[1, 3]	GPUs 1 and 3
"0"	str	[0]	GPU 0
"3"	str	[3]	GPU 3
"1, 3"	str	[1, 3]	GPUs 1 and 3
"-1"	str	[0, 1, 2, ...]	all available GPUs

Note: When specifying number of gpus as an integer `gpus=k`, setting the trainer flag `auto_select_gpus=True` will automatically help you find `k` gpus that are not occupied by other processes. This is especially useful when GPUs are configured to be in “exclusive mode”, such that only one process at a time can access them.

25.2.1 Remove CUDA flags

CUDA flags make certain GPUs visible to your script. Lightning sets these for you automatically, there’s NO NEED to do this yourself.

```
# lightning will set according to what you give the trainer
os.environ["CUDA_DEVICE_ORDER"] = "PCI_BUS_ID"
os.environ["CUDA_VISIBLE_DEVICES"] = "0"
```

However, when using a cluster, Lightning will NOT set these flags (and you should not either). SLURM will set these for you. For more details see the [SLURM cluster guide](#).

25.3 Distributed modes

Lightning allows multiple ways of training

- Data Parallel (`distributed_backend='dp'`) (multiple-gpus, 1 machine)
- DistributedDataParallel (`distributed_backend='ddp'`) (multiple-gpus across many machines (python script based)).
- DistributedDataParallel (`distributed_backend='ddp_spawn'`) (multiple-gpus across many machines (spawn based)).
- DistributedDataParallel 2 (`distributed_backend='ddp2'`) (DP in a machine, DDP across machines).
- Horovod (`distributed_backend='horovod'`) (multi-machine, multi-gpu, configured at runtime)

- TPUs (*tpu_cores=8|x*) (tpu or TPU pod)

Note: If you request multiple GPUs or nodes without setting a mode, DDP will be automatically used.

For a deeper understanding of what Lightning is doing, feel free to read this [guide](#).

25.3.1 Data Parallel

`DataParallel` (DP) splits a batch across k GPUs. That is, if you have a batch of 32 and use DP with 2 gpus, each GPU will process 16 samples, after which the root node will aggregate the results.

Warning: DP use is discouraged by PyTorch and Lightning. Use DDP which is more stable and at least 3x faster

```
# train on 2 GPUs (using DP mode)
trainer = Trainer(gpus=2, distributed_backend='dp')
```

25.3.2 Distributed Data Parallel

`DistributedDataParallel` (DDP) works as follows:

1. Each GPU across each node gets its own process.
2. Each GPU gets visibility into a subset of the overall dataset. It will only ever see that subset.
3. Each process inits the model.

Note: Make sure to set the random seed before the instantiation of a `Trainer()` so that each model initializes with the same weights.

4. Each process performs a full forward and backward pass in parallel.
5. The gradients are synced and averaged across all processes.
6. Each process updates its optimizer.

```
# train on 8 GPUs (same machine (ie: node))
trainer = Trainer(gpus=8, distributed_backend='ddp')

# train on 32 GPUs (4 nodes)
trainer = Trainer(gpus=8, distributed_backend='ddp', num_nodes=4)
```

This Lightning implementation of DDP calls your script under the hood multiple times with the correct environment variables:

```
# example for 3 GPUs DDP
MASTER_ADDR=localhost MASTER_PORT=random() WORLD_SIZE=3 NODE_RANK=0 LOCAL_RANK=0
↪python my_file.py --gpus 3 --etc
MASTER_ADDR=localhost MASTER_PORT=random() WORLD_SIZE=3 NODE_RANK=1 LOCAL_RANK=0
↪python my_file.py --gpus 3 --etc
MASTER_ADDR=localhost MASTER_PORT=random() WORLD_SIZE=3 NODE_RANK=2 LOCAL_RANK=0
↪python my_file.py --gpus 3 --etc
```

We use DDP this way because `ddp_spawn` has a few limitations (due to Python and PyTorch):

1. Since `.spawn()` trains the model in subprocesses, the model on the main process does not get updated.
2. `Dataloader(num_workers=N)`, where `N` is large, bottlenecks training with DDP... ie: it will be VERY slow or won't work at all. This is a PyTorch limitation.
3. Forces everything to be picklable.

There are cases in which it is NOT possible to use DDP. Examples are:

- Jupyter Notebook, Google COLAB, Kaggle, etc.
- You have a nested script without a root package
- Your script needs to invoke both `.fit` and `.test`, or one of them multiple times

In these situations you should use `dp` or `ddp_spawn` instead.

25.3.3 Distributed Data Parallel 2

In certain cases, it's advantageous to use all batches on the same machine instead of a subset. For instance, you might want to compute a NCE loss where it pays to have more negative samples.

In this case, we can use DDP2 which behaves like DP in a machine and DDP across nodes. DDP2 does the following:

1. Copies a subset of the data to each node.
2. Inits a model on each node.
3. Runs a forward and backward pass using DP.
4. Syncs gradients across nodes.
5. Applies the optimizer updates.

```
# train on 32 GPUs (4 nodes)
trainer = Trainer(gpus=8, distributed_backend='ddp2', num_nodes=4)
```

25.3.4 Distributed Data Parallel Spawn

`ddp_spawn` is exactly like `ddp` except that it uses `.spawn` to start the training processes.

Warning: It is STRONGLY recommended to use *DDP* for speed and performance.

```
mp.spawn(self.ddp_train, nprocs=self.num_processes, args=(model, ))
```

If your script does not support being called from the command line (ie: it is nested without a root project module) you can use the following method:

```
# train on 8 GPUs (same machine (ie: node))
trainer = Trainer(gpus=8, distributed_backend='ddp')
```

We STRONGLY discourage this use because it has limitations (due to Python and PyTorch):

1. The model you pass in will not update. Please save a checkpoint and restore from there.
2. Set `Dataloader(num_workers=0)` or it will bottleneck training.

`ddp` is MUCH faster than `ddp_spawn`. We recommend you

1. Install a top-level module for your project using `setup.py`

```
# setup.py
#!/usr/bin/env python

from setuptools import setup, find_packages

setup(name='src',
      version='0.0.1',
      description='Describe Your Cool Project',
      author='',
      author_email='',
      url='https://github.com/YourSeed', # REPLACE WITH YOUR OWN GITHUB PROJECT LINK
      install_requires=[
          'pytorch-lightning'
      ],
      packages=find_packages()
    )
```

2. Setup your project like so:

```
/project
  /src
    some_file.py
    /or_a_folder
  setup.py
```

3. Install as a root-level package

```
cd /project
pip install -e .
```

You can then call your scripts anywhere

```
cd /project/src
python some_file.py --distributed_backend 'ddp' --gpus 8
```

25.3.5 Horovod

Horovod allows the same training script to be used for single-GPU, multi-GPU, and multi-node training.

Like Distributed Data Parallel, every process in Horovod operates on a single GPU with a fixed subset of the data. Gradients are averaged across all GPUs in parallel during the backward pass, then synchronously applied before beginning the next step.

The number of worker processes is configured by a driver application (*horovodrun* or *mpirun*). In the training script, Horovod will detect the number of workers from the environment, and automatically scale the learning rate to compensate for the increased total batch size.

Horovod can be configured in the training script to run with any number of GPUs / processes as follows:

```
# train Horovod on GPU (number of GPUs / machines provided on command-line)
trainer = Trainer(distributed_backend='horovod', gpus=1)

# train Horovod on CPU (number of processes / machines provided on command-line)
trainer = Trainer(distributed_backend='horovod')
```

When starting the training job, the driver application will then be used to specify the total number of worker processes:

```
# run training with 4 GPUs on a single machine
horovodrun -np 4 python train.py

# run training with 8 GPUs on two machines (4 GPUs each)
horovodrun -np 8 -H hostname1:4,hostname2:4 python train.py
```

See the official [Horovod documentation](#) for details on installation and performance tuning.

25.3.6 DP/DDP2 caveats

In DP and DDP2 each GPU within a machine sees a portion of a batch. DP and ddp2 roughly do the following:

```
def distributed_forward(batch, model):
    batch = torch.Tensor(32, 8)
    gpu_0_batch = batch[:8]
    gpu_1_batch = batch[8:16]
    gpu_2_batch = batch[16:24]
    gpu_3_batch = batch[24:]

    y_0 = model_copy_gpu_0(gpu_0_batch)
    y_1 = model_copy_gpu_1(gpu_1_batch)
    y_2 = model_copy_gpu_2(gpu_2_batch)
    y_3 = model_copy_gpu_3(gpu_3_batch)

    return [y_0, y_1, y_2, y_3]
```

So, when Lightning calls any of the *training_step*, *validation_step*, *test_step* you will only be operating on one of those pieces.

```
# the batch here is a portion of the FULL batch
def training_step(self, batch, batch_idx):
    y_0 = batch
```

For most metrics, this doesn't really matter. However, if you want to add something to your computational graph (like softmax) using all batch parts you can use the *training_step_end* step.

```
def training_step_end(self, outputs):
    # only use when on dp
    outputs = torch.cat(outputs, dim=1)
    softmax = softmax(outputs, dim=1)
    out = softmax.mean()
    return out
```

In pseudocode, the full sequence is:

```
# get data
batch = next(dataloader)

# copy model and data to each gpu
batch_splits = split_batch(batch, num_gpus)
models = copy_model_to_gpus(model)

# in parallel, operate on each batch chunk
all_results = []
```

(continues on next page)

(continued from previous page)

```

for gpu_num in gpus:
    batch_split = batch_splits[gpu_num]
    gpu_model = models[gpu_num]
    out = gpu_model(batch_split)
    all_results.append(out)

# use the full batch for something like softmax
full_out = model.training_step_end(all_results)

```

To illustrate why this is needed, let's look at DataParallel

```

def training_step(self, batch, batch_idx):
    x, y = batch
    y_hat = self(batch)

    # on dp or ddp2 if we did softmax now it would be wrong
    # because batch is actually a piece of the full batch
    return y_hat

def training_step_end(self, batch_parts_outputs):
    # batch_parts_outputs has outputs of each part of the batch

    # do softmax here
    outputs = torch.cat(outputs, dim=1)
    softmax = softmax(outputs, dim=1)
    out = softmax.mean()

    return out

```

If `training_step_end` is defined it will be called regardless of TPU, DP, DDP, etc... which means it will behave the same regardless of the backend.

Validation and test step have the same option when using DP.

```

def validation_step_end(self, batch_parts_outputs):
    ...

def test_step_end(self, batch_parts_outputs):
    ...

```

25.3.7 Distributed and 16-bit precision

Due to an issue with Apex and DataParallel (PyTorch and NVIDIA issue), Lightning does not allow 16-bit and DP training. We tried to get this to work, but it's an issue on their end.

Below are the possible configurations we support.

1 GPU	1+ GPUs	DP	DDP	16-bit	command
Y					<code>Trainer(gpus=1)</code>
Y				Y	<code>Trainer(gpus=1, use_amp=True)</code>
	Y	Y			<code>Trainer(gpus=k, distributed_backend='dp')</code>
	Y		Y		<code>Trainer(gpus=k, distributed_backend='ddp')</code>
	Y		Y	Y	<code>Trainer(gpus=k, distributed_backend='ddp', use_amp=True)</code>

25.3.8 Implement Your Own Distributed (DDP) training

If you need your own way to init PyTorch DDP you can override `pytorch_lightning.core.LightningModule()`.

If you also need to use your own DDP implementation, override: `pytorch_lightning.core.LightningModule.configure_ddp()`.

25.4 Batch size

When using distributed training make sure to modify your learning rate according to your effective batch size.

Let's say you have a batch size of 7 in your dataloader.

```
class LitModel(LightningModule):  
  
    def train_dataloader(self):  
        return Dataset(..., batch_size=7)
```

In (DDP, Horovod) your effective batch size will be $7 * \text{gpus} * \text{num_nodes}$.

```
# effective batch size = 7 * 8  
Trainer(gpus=8, distributed_backend='ddp|horovod')  
  
# effective batch size = 7 * 8 * 10  
Trainer(gpus=8, num_nodes=10, distributed_backend='ddp|horovod')
```

In DDP2, your effective batch size will be $7 * \text{num_nodes}$. The reason is that the full batch is visible to all GPUs on the node when using DDP2.

```
# effective batch size = 7  
Trainer(gpus=8, distributed_backend='ddp2')  
  
# effective batch size = 7 * 10  
Trainer(gpus=8, num_nodes=10, distributed_backend='ddp2')
```

Note: Huge batch sizes are actually really bad for convergence. Check out: [Accurate, Large Minibatch SGD: Training ImageNet in 1 Hour](#)

25.5 PytorchElastic

Lightning supports the use of PytorchElastic to enable fault-tolerant and elastic distributed job scheduling. To use it, specify the 'ddp' or 'ddp2' backend and the number of gpus you want to use in the trainer.

```
Trainer(gpus=8, distributed_backend='ddp')
```

Following the [PytorchElastic Quickstart documentation](#), you then need to start a single-node etcd server on one of the hosts:

```
etcd --enable-v2
--listen-client-urls http://0.0.0.0:2379,http://127.0.0.1:4001
--advertise-client-urls PUBLIC_HOSTNAME:2379
```

And then launch the elastic job with:

```
python -m torchelastic.distributed.launch
--nnodes=MIN_SIZE:MAX_SIZE
--nproc_per_node=TRAINERS_PER_NODE
--rdzv_id=JOB_ID
--rdzv_backend=etcd
--rdzv_endpoint=ETCD_HOST:ETCD_PORT
YOUR_LIGHTNING_TRAINING_SCRIPT.py (--arg1 ... train script args...)
```

See the official [PytorchElastic documentation](#) for details on installation and more use cases.

25.6 Jupyter Notebooks

Unfortunately any *ddp_* is not supported in jupyter notebooks. Please use *dp* for multiple GPUs. This is a known Jupyter issue. If you feel like taking a stab at adding this support, feel free to submit a PR!

25.7 Pickle Errors

Multi-GPU training sometimes requires your model to be pickled. If you run into an issue with pickling try the following to figure out the issue

```
import pickle

model = YourModel()
pickle.dumps(model)
```

However, if you use *ddp* the pickling requirement is not there and you should be fine. If you use *ddp_spawn* the pickling requirement remains. This is a limitation of Python.

MULTIPLE DATASETS

Lightning supports multiple dataloaders in a few ways.

1. Create a dataloader that iterates multiple datasets under the hood.
 2. In the validation and test loop you also have the option to return multiple dataloaders which lightning will call sequentially.
-

26.1 Multiple training dataloaders

For training, the best way to use multiple dataloaders is to create a `DataLoader` class which wraps your multiple dataloaders (this of course also works for testing and validation dataloaders).

(reference)

```
class ConcatDataset(torch.utils.data.Dataset):
    def __init__(self, *datasets):
        self.datasets = datasets

    def __getitem__(self, i):
        return tuple(d[i] for d in self.datasets)

    def __len__(self):
        return min(len(d) for d in self.datasets)

class LitModel(LightningModule):

    def train_dataloader(self):
        concat_dataset = ConcatDataset(
            datasets.ImageFolder(trainindir_A),
            datasets.ImageFolder(trainindir_B)
        )

        loader = torch.utils.data.DataLoader(
            concat_dataset,
            batch_size=args.batch_size,
            shuffle=True,
            num_workers=args.workers,
            pin_memory=True
        )
        return loader
```

(continues on next page)

(continued from previous page)

```
def val_dataloader(self):  
    # SAME  
    ...  
  
def test_dataloader(self):  
    # SAME  
    ...
```

26.2 Test/Val dataloaders

For validation and test dataloaders, lightning also gives you the additional option of passing multiple dataloaders back from each call.

See the following for more details:

- `val_dataloader()`
- `test_dataloader()`

```
def val_dataloader(self):  
    loader_1 = Dataloader()  
    loader_2 = Dataloader()  
    return [loader_1, loader_2]
```


SAVING AND LOADING WEIGHTS

Lightning automates saving and loading checkpoints. Checkpoints capture the exact value of all parameters used by a model.

Checkpointing your training allows you to resume a training process in case it was interrupted, fine-tune a model or use a pre-trained model for inference without having to retrain the model.

27.1 Checkpoint saving

A Lightning checkpoint has everything needed to restore a training session including:

- 16-bit scaling factor (apex)
- Current epoch
- Global step
- Model state_dict
- State of all optimizers
- State of all learningRate schedulers
- State of all callbacks
- The hyperparameters used for that model if passed in as hparams (Argparse.Namespace)

27.1.1 Automatic saving

Lightning automatically saves a checkpoint for you in your current working directory, with the state of your last training epoch. This makes sure you can resume training in case it was interrupted.

To change the checkpoint path pass in:

```
# saves checkpoints to '/your/path/to/save/checkpoints' at every epoch end
trainer = Trainer(default_root_dir='/your/path/to/save/checkpoints')
```

You can customize the checkpointing behavior to monitor any quantity of your training or validation steps. For example, if you want to update your checkpoints based on your validation loss:

1. Calculate any metric or other quantity you wish to monitor, such as validation loss.
2. Log the quantity using `log()` method, with a key such as `val_loss`.
3. Initializing the `ModelCheckpoint` callback, and set `monitor` to be the key of your quantity.
4. Pass the callback to `checkpoint_callback` Trainer flag.

```
from pytorch_lightning.callbacks import ModelCheckpoint

class LitAutoEncoder(pl.LightningModule):
    def validation_step(self, batch, batch_idx):
        x, y = batch
        y_hat = self.backbone(x)

        # 1. calculate loss
        loss = F.cross_entropy(y_hat, y)

        # 2. log `val_loss`
        self.log('val_loss', loss)

# 3. Init ModelCheckpoint callback, monitoring 'val_loss'
checkpoint_callback = ModelCheckpoint(monitor='val_loss')

# 4. Pass your callback to checkpoint_callback trainer flag
trainer = Trainer(checkpoint_callback=checkpoint_callback)
```

You can also control more advanced options, like `save_top_k`, to save the best k models and the mode of the monitored quantity (min/max/auto, where the mode is automatically inferred from the name of the monitored quantity), `save_weights_only` or `period` to set the interval of epochs between checkpoints, to avoid slowdowns.

```
from pytorch_lightning.callbacks import ModelCheckpoint

class LitAutoEncoder(pl.LightningModule):
    def validation_step(self, batch, batch_idx):
        x, y = batch
        y_hat = self.backbone(x)
        loss = F.cross_entropy(y_hat, y)
        self.log('val_loss', loss)

# saves a file like: my/path/sample-mnist-epoch=02-val_loss=0.32.ckpt
checkpoint_callback = ModelCheckpoint(
    monitor='val_loss',
    filepath='my/path/sample-mnist-{epoch:02d}-{val_loss:.2f}' ,
    save_top_k=3,
    mode='min')

trainer = Trainer(checkpoint_callback=checkpoint_callback)
```

You can retrieve the checkpoint after training by calling

```
checkpoint_callback = ModelCheckpoint(filepath='my/path/')
trainer = Trainer(checkpoint_callback=checkpoint_callback)
trainer.fit(model)
checkpoint_callback.best_model_path
```

Disabling checkpoints

You can disable checkpointing by passing

```
trainer = Trainer(checkpoint_callback=False)
```

The Lightning checkpoint also saves the arguments passed into the LightningModule init under the *module_arguments* key in the checkpoint.

```
class MyLightningModule(LightningModule):

    def __init__(self, learning_rate, *args, **kwargs):
        super().__init__()

    # all init args were saved to the checkpoint
    checkpoint = torch.load(CKPT_PATH)
    print(checkpoint['module_arguments'])
    # {'learning_rate': the_value}
```

27.1.2 Manual saving

You can manually save checkpoints and restore your model from the checkpointed state.

```
model = MyLightningModule(hparams)
trainer.fit(model)
trainer.save_checkpoint("example.ckpt")
new_model = MyModel.load_from_checkpoint(checkpoint_path="example.ckpt")
```

27.2 Checkpoint loading

To load a model along with its weights, biases and *module_arguments* use the following method:

```
model = MyLightningModule.load_from_checkpoint(PATH)

print(model.learning_rate)
# prints the learning_rate you used in this checkpoint

model.eval()
y_hat = model(x)
```

But if you don't want to use the values saved in the checkpoint, pass in your own here

```
class LitModel(LightningModule):

    def __init__(self, in_dim, out_dim):
        super().__init__()
        self.save_hyperparameters()
        self.ll = nn.Linear(self.hparams.in_dim, self.hparams.out_dim)
```

you can restore the model like this

```
# if you train and save the model like this it will use these values when loading
# the weights. But you can overwrite this
```

(continues on next page)

(continued from previous page)

```
LitModel(in_dim=32, out_dim=10)

# uses in_dim=32, out_dim=10
model = LitModel.load_from_checkpoint(PATH)

# uses in_dim=128, out_dim=10
model = LitModel.load_from_checkpoint(PATH, in_dim=128, out_dim=10)
```

27.2.1 Restoring Training State

If you don't just want to load weights, but instead restore the full training, do the following:

```
model = LitModel()
trainer = Trainer(resume_from_checkpoint='some/path/to/my_checkpoint.ckpt')

# automatically restores model, epoch, step, LR schedulers, apex, etc...
trainer.fit(model)
```

OPTIMIZATION

Lightning offers two modes for managing the optimization process:

- automatic optimization (AutoOpt)
- manual optimization

For the majority of research cases, **automatic optimization** will do the right thing for you and it is what most users should use.

For advanced/expert users who want to do esoteric optimization schedules or techniques, use **manual optimization**.

28.1 Manual optimization

For advanced research topics like reinforcement learning, sparse coding, or GAN research, it may be desirable to manually manage the optimization process. To do so, do the following:

- Disable automatic optimization in Trainer: `Trainer(automatic_optimization=False)`
- Drop or ignore the `optimizer_idx` argument
- Use `self.manual_backward(loss)` instead of `loss.backward()` to automatically scale your loss

```
def training_step(self, batch, batch_idx, optimizer_idx):
    # ignore optimizer_idx
    (opt_g, opt_d) = self.optimizers()

    # do anything you want
    loss_a = ...

    # use self.backward which will also handle scaling the loss when using amp
    self.manual_backward(loss_a, opt_g)
    opt_g.step()
    opt_g.zero_grad()

    # do anything you want
    loss_b = ...

    # pass in any args that loss.backward() normally takes
    self.manual_backward(loss_b, opt_d, retain_graph=True)
    self.manual_backward(loss_b, opt_d)
    opt_d.step()
    opt_d.zero_grad()
```

Note: This is only recommended for experts who need ultimate flexibility

Manual optimization does not yet support accumulated gradients but will be live in 1.1.0

28.2 Automatic optimization

With Lightning most users don't have to think about when to call `.backward()`, `.step()`, `.zero_grad()`, since Lightning automates that for you.

Under the hood Lightning does the following:

```
for epoch in epochs:
    for batch_id data:
        loss = model.training_step(batch, batch_idx, ...)
        loss.backward()
        optimizer.step()
        optimizer.zero_grad()

    for scheduler in scheduler:
        scheduler.step()
```

In the case of multiple optimizers, Lightning does the following:

```
for epoch in epochs:
    for batch in data:
        for opt in optimizers:
            disable_grads_for_other_optimizers()
            train_step(opt)
            opt.step()

    for scheduler in scheduler:
        scheduler.step()
```

28.2.1 Learning rate scheduling

Every optimizer you use can be paired with any [LearningRateScheduler](#).

```
# no LR scheduler
def configure_optimizers(self):
    return Adam(...)

# Adam + LR scheduler
def configure_optimizers(self):
    optimizer = Adam(...)
    scheduler = ReduceLROnPlateau(optimizer, ...)
    return [optimizer], [scheduler]

# Two optimizers each with a scheduler
def configure_optimizers(self):
    optimizer1 = Adam(...)
    optimizer2 = SGD(...)
```

(continues on next page)

(continued from previous page)

```

scheduler1 = ReduceLROnPlateau(optimizer1, ...)
scheduler2 = LambdaLR(optimizer2, ...)
return [optimizer1, optimizer2], [scheduler1, scheduler2]

# Same as above with additional params passed to the first scheduler
def configure_optimizers(self):
    optimizers = [Adam(...), SGD(...)]
    schedulers = [
        {
            'scheduler': ReduceLROnPlateau(optimizers[0], ...),
            'monitor': 'val_recall', # Default: val_loss
            'interval': 'epoch',
            'frequency': 1
        },
        LambdaLR(optimizers[1], ...)
    ]
    return optimizers, schedulers

```

28.2.2 Use multiple optimizers (like GANs)

To use multiple optimizers return > 1 optimizers from `pytorch_lightning.core.LightningModule.configure_optimizers()`

```

# one optimizer
def configure_optimizers(self):
    return Adam(...)

# two optimizers, no schedulers
def configure_optimizers(self):
    return Adam(...), SGD(...)

# Two optimizers, one scheduler for adam only
def configure_optimizers(self):
    return [Adam(...), SGD(...)], [ReduceLROnPlateau()]

```

Lightning will call each optimizer sequentially:

```

for epoch in epochs:
    for batch in data:
        for opt in optimizers:
            train_step(opt)
            opt.step()

    for scheduler in scheduler:
        scheduler.step()

```

28.2.3 Step optimizers at arbitrary intervals

To do more interesting things with your optimizers such as learning rate warm-up or odd scheduling, override the `optimizer_step()` function.

For example, here step optimizer A every 2 batches and optimizer B every 4 batches

```
def optimizer_step(self, current_epoch, batch_nb, optimizer, optimizer_idx, second_
    ↳order_closure=None, on_tpu=False, using_native_amp=False, using_lbfgs=False):
    optimizer.step()

def optimizer_zero_grad(self, current_epoch, batch_idx, optimizer, opt_idx):
    optimizer.zero_grad()

# Alternating schedule for optimizer steps (ie: GANs)
def optimizer_step(self, current_epoch, batch_nb, optimizer, optimizer_idx, second_
    ↳order_closure=None, on_tpu=False, using_native_amp=False, using_lbfgs=False):
    # update generator opt every 2 steps
    if optimizer_i == 0:
        if batch_nb % 2 == 0 :
            optimizer.step()
            optimizer.zero_grad()

    # update discriminator opt every 4 steps
    if optimizer_i == 1:
        if batch_nb % 4 == 0 :
            optimizer.step()
            optimizer.zero_grad()

    # ...
    # add as many optimizers as you want
```

Here we add a learning-rate warm up

```
# learning rate warm-up
def optimizer_step(self, current_epoch, batch_nb, optimizer, optimizer_idx, second_
    ↳order_closure=None, on_tpu=False, using_native_amp=False, using_lbfgs=False):
    # warm up lr
    if self.trainer.global_step < 500:
        lr_scale = min(1., float(self.trainer.global_step + 1) / 500.)
        for pg in optimizer.param_groups:
            pg['lr'] = lr_scale * self.hparams.learning_rate

    # update params
    optimizer.step()
    optimizer.zero_grad()
```


28.2.4 Using the closure functions for optimization

When using optimization schemes such as LBFGS, the `second_order_closure` needs to be enabled. By default, this function is defined by wrapping the `training_step` and the backward steps as follows

```
def second_order_closure(pl_module, split_batch, batch_idx, opt_idx, optimizer,
    ↪hidden):
    # Model training step on a given batch
    result = pl_module.training_step(split_batch, batch_idx, opt_idx, hidden)

    # Model backward pass
    pl_module.backward(result, optimizer, opt_idx)

    # on_after_backward callback
    pl_module.on_after_backward(result.training_step_output, batch_idx, result.loss)

    return result

# This default `second_order_closure` function can be enabled by passing it directly
    ↪into the `optimizer.step`
def optimizer_step(self, current_epoch, batch_nb, optimizer, optimizer_idx, second_
    ↪order_closure, on_tpu=False, using_native_amp=False, using_lbfgs=False):
    # update params
    optimizer.step(second_order_closure)
```


PERFORMANCE AND BOTTLENECK PROFILER

Profiling your training run can help you understand if there are any bottlenecks in your code.

29.1 Built-in checks

PyTorch Lightning supports profiling standard actions in the training loop out of the box, including:

- `on_epoch_start`
- `on_epoch_end`
- `on_batch_start`
- `tbptt_split_batch`
- `model_forward`
- `model_backward`
- `on_after_backward`
- `optimizer_step`
- `on_batch_end`
- `training_step_end`
- `on_training_end`

29.2 Enable simple profiling

If you only wish to profile the standard actions, you can set `profiler=True` when constructing your *Trainer* object.

```
trainer = Trainer(..., profiler=True)
```

The profiler's results will be printed at the completion of a training *fit()*.

Profiler Report

Action	Mean duration (s)	Total time (s)
on_epoch_start	5.993e-06	5.993e-06
get_train_batch	0.0087412	16.398
on_batch_start	5.0865e-06	0.0095372

(continues on next page)

(continued from previous page)

model_forward	0.0017818	3.3408
model_backward	0.0018283	3.4282
on_after_backward	4.2862e-06	0.0080366
optimizer_step	0.0011072	2.0759
on_batch_end	4.5202e-06	0.0084753
on_epoch_end	3.919e-06	3.919e-06
on_train_end	5.449e-06	5.449e-06

29.3 Advanced Profiling

If you want more information on the functions called during each event, you can use the *AdvancedProfiler*. This option uses Python's *cProfiler* to provide a report of time spent on *each* function called within your code.

```
profiler = AdvancedProfiler()
trainer = Trainer(..., profiler=profiler)
```

The profiler's results will be printed at the completion of a training *fit()*. This profiler report can be quite long, so you can also specify an *output_filename* to save the report instead of logging it to the output in your terminal. The output below shows the profiling for the action *get_train_batch*.

```
Profiler Report

Profile stats for: get_train_batch
    4869394 function calls (4863767 primitive calls) in 18.893 seconds
Ordered by: cumulative time
List reduced from 76 to 10 due to restriction <10>
ncalls  tottime  percall  cumtime  percall filename:lineno(function)
3752/1876    0.011    0.000    18.887    0.010 {built-in method builtins.next}
    1876    0.008    0.000    18.877    0.010 dataloader.py:344(__next__)
    1876    0.074    0.000    18.869    0.010 dataloader.py:383(_next_data)
    1875    0.012    0.000    18.721    0.010 fetch.py:42(fetch)
    1875    0.084    0.000    18.290    0.010 fetch.py:44(<listcomp>)
60000    1.759    0.000    18.206    0.000 mnist.py:80(__getitem__)
60000    0.267    0.000    13.022    0.000 transforms.py:68(__call__)
60000    0.182    0.000    7.020    0.000 transforms.py:93(__call__)
60000    1.651    0.000    6.839    0.000 functional.py:42(to_tensor)
60000    0.260    0.000    5.734    0.000 transforms.py:167(__call__)
```

You can also reference this profiler in your *LightningModule* to profile specific actions of interest. If you don't want to always have the profiler turned on, you can optionally pass a *PassThroughProfiler* which will allow you to skip profiling without having to make any code changes. Each profiler has a method *profile()* which returns a context handler. Simply pass in the name of your action that you want to track and the profiler will record performance for code executed within this context.

```
from pytorch_lightning.profiler import Profiler, PassThroughProfiler

class MyModel(LightningModule):
    def __init__(self, profiler=None):
        self.profiler = profiler or PassThroughProfiler()

    def custom_processing_step(self, data):
        with profiler.profile('my_custom_action'):
            # custom processing step
```

(continues on next page)

(continued from previous page)

```

    return data

profiler = Profiler()
model = MyModel(profiler)
trainer = Trainer(profiler=profiler, max_epochs=1)

```

class pytorch_lightning.profiler.**BaseProfiler** (*output_streams=None*)
 Bases: `abc.ABC`

If you wish to write a custom profiler, you should inherit from this class.

Params: `stream_out`: callable

describe ()

Logs a profile report after the conclusion of the training run.

Return type `None`

profile (*action_name*)

Yields a context manager to encapsulate the scope of a profiled action.

Example:

```

with self.profile('load training data'):
    # load training data code

```

The profiler will start once you've entered the context and will automatically stop once you exit the code block.

Return type `None`

abstract start (*action_name*)

Defines how to start recording an action.

Return type `None`

abstract stop (*action_name*)

Defines how to record the duration once an action is complete.

Return type `None`

abstract summary ()

Create profiler summary in text format.

Return type `str`

class pytorch_lightning.profiler.**SimpleProfiler** (*output_filename=None*)

Bases: `pytorch_lightning.profiler.profilers.BaseProfiler`

This profiler simply records the duration of actions (in seconds) and reports the mean duration of each action and the total time spent over the entire training run.

Params:

output_filename (`str`): optionally save profile results to file instead of printing to std out when training is finished.

describe ()

Logs a profile report after the conclusion of the training run.

start (*action_name*)

Defines how to start recording an action.

Return type `None`

stop (*action_name*)

Defines how to record the duration once an action is complete.

Return type `None`

summary ()

Create profiler summary in text format.

Return type `str`

```
class pytorch_lightning.profiler.AdvancedProfiler (output_filename=None,  
                                              line_count_restriction=1.0)
```

Bases: `pytorch_lightning.profiler.profilers.BaseProfiler`

This profiler uses Python's cProfiler to record more detailed information about time spent in each function call recorded during a given action. The output is quite verbose and you should only use this if you want very detailed reports.

Parameters

- **output_filename** `// (Optional[str])` – optionally save profile results to file instead of printing to std out when training is finished.
- **line_count_restriction** `// (float)` – this can be used to limit the number of functions reported for each action. either an integer (to select a count of lines), or a decimal fraction between 0.0 and 1.0 inclusive (to select a percentage of lines)

describe ()

Logs a profile report after the conclusion of the training run.

start (*action_name*)

Defines how to start recording an action.

Return type `None`

stop (*action_name*)

Defines how to record the duration once an action is complete.

Return type `None`

summary ()

Create profiler summary in text format.

Return type `str`

```
class pytorch_lightning.profiler.PassThroughProfiler
```

Bases: `pytorch_lightning.profiler.profilers.BaseProfiler`

This class should be used when you don't want the (small) overhead of profiling. The Trainer uses this class by default.

Params: `stream_out`: callable

start (*action_name*)

Defines how to start recording an action.

Return type `None`

stop (*action_name*)

Defines how to record the duration once an action is complete.

Return type `None`

summary ()

Create profiler summary in text format.

Return type `str`

SINGLE GPU TRAINING

Make sure you are running on a machine that has at least one GPU. Lightning handles all the NVIDIA flags for you, there's no need to set them yourself.

```
# train on 1 GPU (using dp mode)
trainer = Trainer(gpus=1)
```


SEQUENTIAL DATA

Lightning has built in support for dealing with sequential data.

31.1 Packed sequences as inputs

When using PackedSequence, do 2 things:

1. Return either a padded tensor in dataset or a list of variable length tensors in the dataloader collate_fn (example shows the list implementation).
2. Pack the sequence in forward or training and validation steps depending on use case.

```
# For use in dataloader
def collate_fn(batch):
    x = [item[0] for item in batch]
    y = [item[1] for item in batch]
    return x, y

# In module
def training_step(self, batch, batch_nb):
    x = rnn.pack_sequence(batch[0], enforce_sorted=False)
    y = rnn.pack_sequence(batch[1], enforce_sorted=False)
```

31.2 Truncated Backpropagation Through Time

There are times when multiple backwards passes are needed for each batch. For example, it may save memory to use Truncated Backpropagation Through Time when training RNNs.

Lightning can handle TBTT automatically via this flag.

```
# DEFAULT (single backwards pass per batch)
trainer = Trainer(truncated_bptt_steps=None)

# (split batch into sequences of size 2)
trainer = Trainer(truncated_bptt_steps=2)
```

Note: If you need to modify how the batch is split, override `pytorch_lightning.core.LightningModule.tbptt_split_batch()`.

Note: Using this feature requires updating your `LightningModule`'s `pytorch_lightning.core.LightningModule.training_step()` to include a *hidens* arg.

31.3 Iterable Datasets

Lightning supports using `IterableDatasets` as well as map-style `Datasets`. `IterableDatasets` provide a more natural option when using sequential data.

Note: When using an `IterableDataset` you must set the `val_check_interval` to 1.0 (the default) or an int (specifying the number of training batches to run before validation) when initializing the `Trainer`. This is because the `IterableDataset` does not have a `__len__` and Lightning requires this to calculate the validation interval when `val_check_interval` is less than one. Similarly, you can set `limit_{mode}_batches` to a float or an int. If it is set to 0.0 or 0 it will set `num_{mode}_batches` to 0, if it is an int it will set `num_{mode}_batches` to `limit_{mode}_batches`, if it is set to 1.0 it will run for the whole dataset, otherwise it will throw an exception. Here mode can be train/val/test.

```
# IterableDataset
class CustomDataset(IterableDataset):

    def __init__(self, data):
        self.data_source

    def __iter__(self):
        return iter(self.data_source)

# Setup DataLoader
def train_dataloader(self):
    seq_data = ['A', 'long', 'time', 'ago', 'in', 'a', 'galaxy', 'far', 'far', 'away']
    iterable_dataset = CustomDataset(seq_data)

    dataloader = DataLoader(dataset=iterable_dataset, batch_size=5)
    return dataloader
```

```
# Set val_check_interval
trainer = Trainer(val_check_interval=100)

# Set limit_val_batches to 0.0 or 0
trainer = Trainer(limit_val_batches=0.0)

# Set limit_val_batches as an int
trainer = Trainer(limit_val_batches=100)
```

TRAINING TRICKS

Lightning implements various tricks to help during training

32.1 Accumulate gradients

Accumulated gradients runs K small batches of size N before doing a backwards pass. The effect is a large effective batch size of size KxN.

See also:

Trainer

```
# DEFAULT (ie: no accumulated grads)
trainer = Trainer(accumulate_grad_batches=1)
```

32.2 Gradient Clipping

Gradient clipping may be enabled to avoid exploding gradients. Specifically, this will [clip the gradient norm](#) computed over all model parameters together.

See also:

Trainer

```
# DEFAULT (ie: don't clip)
trainer = Trainer(gradient_clip_val=0)

# clip gradients with norm above 0.5
trainer = Trainer(gradient_clip_val=0.5)
```

32.3 Auto scaling of batch size

Auto scaling of batch size may be enabled to find the largest batch size that fits into memory. Larger batch size often yields better estimates of gradients, but may also result in longer training time. Inspired by <https://github.com/BlackHC/toma>.

See also:

Trainer

```
# DEFAULT (ie: don't scale batch size automatically)
trainer = Trainer(auto_scale_batch_size=None)

# Autoscale batch size
trainer = Trainer(auto_scale_batch_size=None|'power'|'binsearch')

# find the batch size
trainer.tune(model)
```

Currently, this feature supports two modes *'power'* scaling and *'binsearch'* scaling. In *'power'* scaling, starting from a batch size of 1 keeps doubling the batch size until an out-of-memory (OOM) error is encountered. Setting the argument to *'binsearch'* will initially also try doubling the batch size until it encounters an OOM, after which it will do a binary search that will finetune the batch size. Additionally, it should be noted that the batch size scaler cannot search for batch sizes larger than the size of the training dataset.

Note: This feature expects that a *batch_size* field is either located as a model attribute i.e. *model.batch_size* or as a field in your *hparams* i.e. *model.hparams.batch_size*. The field should exist and will be overridden by the results of this algorithm. Additionally, your *train_dataloader()* method should depend on this field for this feature to work i.e.

```
def train_dataloader(self):
    return DataLoader(train_dataset, batch_size=self.batch_size|self.hparams.batch_
↪size)
```

Warning: Due to these constraints, this features does *NOT* work when passing dataloaders directly to *.fit()*.

The scaling algorithm has a number of parameters that the user can control by invoking the trainer method *.scale_batch_size* themselves (see description below).

```
# Use default in trainer construction
trainer = Trainer()
tuner = Tuner(trainer)

# Invoke method
new_batch_size = tuner.scale_batch_size(model, *extra_parameters_here)

# Override old batch size
model.hparams.batch_size = new_batch_size

# Fit as normal
trainer.fit(model)
```

The algorithm in short works by:

1. Dumping the current state of the model and trainer

2. Iteratively until convergence or maximum number of tries *max_trials* (default 25) has been reached:

- Call *fit()* method of trainer. This evaluates *steps_per_trial* (default 3) number of training steps. Each training step can trigger an OOM error if the tensors (training batch, weights, gradients ect.) allocated during the steps have a too large memory footprint.
 - If an OOM error is encountered, decrease batch size else increase it. How much the batch size is increased/decreased is determined by the choosen stratgy.
3. The found batch size is saved to either *model.batch_size* or *model.hparams.batch_size*
 4. Restore the initial state of model and trainer

```
class pytorch_lightning.tuner.tuning.Tuner(trainer)
```

Bases: `object`

```
scale_batch_size(model, mode='power', steps_per_trial=3, init_val=2, max_trials=25,
                  batch_arg_name='batch_size', **fit_kwargs)
```

Will iteratively try to find the largest batch size for a given model that does not give an out of memory (OOM) error.

Parameters

- **model** – Model to fit.
- **mode** (`str`) – string setting the search mode. Either *power* or *binsearch*. If mode is *power* we keep multiplying the batch size by 2, until we get an OOM error. If mode is 'binsearch', we will initially also keep multiplying by 2 and after encountering an OOM error do a binary search between the last successful batch size and the batch size that failed.
- **steps_per_trial** (`int`) – number of steps to run with a given batch size. Ideally 1 should be enough to test if a OOM error occurs, however in practise a few are needed
- **init_val** (`int`) – initial batch size to start the search with
- **max_trials** (`int`) – max number of increase in batch size done before algorithm is terminated
- **batch_arg_name** (`str`) – name of the attribute that stores the batch size. It is expected that the user has provided a model or datamodule that has a hyperparameter with that name. We will look for this attribute name in the following places
 - *model*
 - *model.hparams*
 - *model.datamodule*
 - *trainer.datamodule* (the datamodule passed to the tune method)
- ****fit_kwargs** – remaining arguments to be passed to *.fit()*, e.g., *dataloader* or *data-module*.

Warning: Batch size finder is not supported for DDP yet, it is coming soon.

TRANSFER LEARNING

33.1 Using Pretrained Models

Sometimes we want to use a `LightningModule` as a pretrained model. This is fine because a `LightningModule` is just a `torch.nn.Module`!

Note: Remember that a `LightningModule` is EXACTLY a `torch.nn.Module` but with more capabilities.

Let's use the *AutoEncoder* as a feature extractor in a separate model.

```
class Encoder(torch.nn.Module):
    ...

class AutoEncoder(LightningModule):
    def __init__(self):
        self.encoder = Encoder()
        self.decoder = Decoder()

class CIFAR10Classifier(LightningModule):
    def __init__(self):
        # init the pretrained LightningModule
        self.feature_extractor = AutoEncoder.load_from_checkpoint(PATH)
        self.feature_extractor.freeze()

        # the autoencoder outputs a 100-dim representation and CIFAR-10 has 10 classes
        self.classifier = nn.Linear(100, 10)

    def forward(self, x):
        representations = self.feature_extractor(x)
        x = self.classifier(representations)
        ...
```

We used our pretrained Autoencoder (a `LightningModule`) for transfer learning!

33.2 Example: Imagenet (computer Vision)

```
import torchvision.models as models

class ImagenetTransferLearning(LightningModule):
    def __init__(self):
        # init a pretrained resnet
        num_target_classes = 10
        self.feature_extractor = models.resnet50(pretrained=True)
        self.feature_extractor.eval()

        # use the pretrained model to classify cifar-10 (10 image classes)
        self.classifier = nn.Linear(2048, num_target_classes)

    def forward(self, x):
        representations = self.feature_extractor(x)
        x = self.classifier(representations)
        ...
```

Finetune

```
model = ImagenetTransferLearning()
trainer = Trainer()
trainer.fit(model)
```

And use it to predict your data of interest

```
model = ImagenetTransferLearning.load_from_checkpoint(PATH)
model.freeze()

x = some_images_from_cifar10()
predictions = model(x)
```

We used a pretrained model on imagenet, finetuned on CIFAR-10 to predict on CIFAR-10. In the non-academic world we would finetune on a tiny dataset you have and predict on your dataset.

33.3 Example: BERT (NLP)

Lightning is completely agnostic to what's used for transfer learning so long as it is a `torch.nn.Module` subclass.

Here's a model that uses [Huggingface transformers](#).

```
class BertMNLIFinetuner(LightningModule):

    def __init__(self):
        super().__init__()

        self.bert = BertModel.from_pretrained('bert-base-cased', output_
→attentions=True)
        self.W = nn.Linear(bert.config.hidden_size, 3)
        self.num_classes = 3

    def forward(self, input_ids, attention_mask, token_type_ids):
```

(continues on next page)

(continued from previous page)

```
h, _, attn = self.bert(input_ids=input_ids,
                        attention_mask=attention_mask,
                        token_type_ids=token_type_ids)

h_cls = h[:, 0]
logits = self.W(h_cls)
return logits, attn
```


TPU SUPPORT

Lightning supports running on TPUs. At this moment, TPUs are available on Google Cloud (GCP), Google Colab and Kaggle Environments. For more information on TPUs [watch this video](#).

34.1 TPU Terminology

A TPU is a Tensor processing unit. Each TPU has 8 cores where each core is optimized for 128x128 matrix multiplies. In general, a single TPU is about as fast as 5 V100 GPUs!

A TPU pod hosts many TPUs on it. Currently, TPU pod v2 has 2048 cores! You can request a full pod from Google cloud or a “slice” which gives you some subset of those 2048 cores.

34.2 How to access TPUs

To access TPUs, there are three main ways.

1. Using Google Colab.
 2. Using Google Cloud (GCP).
 3. Using Kaggle.
-

34.3 Colab TPUs

Colab is like a jupyter notebook with a free GPU or TPU hosted on GCP.

To get a TPU on colab, follow these steps:

1. Go to <https://colab.research.google.com/>.
 2. Click “new notebook” (bottom right of pop-up).
-

3. Click runtime > change runtime settings. Select Python 3, and hardware accelerator “TPU”. This will give you a TPU with 8 cores.
4. Next, insert this code into the first cell and execute. This will install the xla library that interfaces between PyTorch and the TPU.

```
!curl https://raw.githubusercontent.com/pytorch/xla/master/contrib/scripts/env-  
↪setup.py -o pytorch-xla-env-setup.py  
!python pytorch-xla-env-setup.py --version nightly --apt-packages libomp5_  
↪libopenblas-dev
```

5. Once the above is done, install PyTorch Lightning (v 0.7.0+).

```
!pip install pytorch-lightning
```

6. Then set up your LightningModule as normal.
-

34.4 DistributedSamplers

Lightning automatically inserts the correct samplers - no need to do this yourself!

Usually, with TPUs (and DDP), you would need to define a DistributedSampler to move the right chunk of data to the appropriate TPU. As mentioned, this is not needed in Lightning

Note: Don’t add distributedSamplers. Lightning does this automatically

If for some reason you still need to, this is how to construct the sampler for TPU use

```
import torch_xla.core.xla_model as xm  
  
def train_dataloader(self):  
    dataset = MNIST(  
        os.getcwd(),  
        train=True,  
        download=True,  
        transform=transforms.ToTensor()  
    )  
  
    # required for TPU support  
    sampler = None  
    if use_tpu:  
        sampler = torch.utils.data.distributed.DistributedSampler(  
            dataset,  
            num_replicas=xm.xrt_world_size(),  
            rank=xm.get_ordinal(),  
            shuffle=True  
        )  
  
    loader = DataLoader(  
        dataset,  
        sampler=sampler,  
        batch_size=32  
    )
```

(continues on next page)

(continued from previous page)

```
return loader
```

Configure the number of TPU cores in the trainer. You can only choose 1 or 8. To use a full TPU pod skip to the TPU pod section.

```
import pytorch_lightning as pl

my_model = MyLightningModule()
trainer = pl.Trainer(tpu_cores=8)
trainer.fit(my_model)
```

That's it! Your model will train on all 8 TPU cores.

34.5 Single TPU core training

Lightning supports training on a single TPU core. Just pass the TPU core ID [1-8] in a list.

```
trainer = pl.Trainer(tpu_cores=[1])
```

34.6 Distributed Backend with TPU

The `distributed_backend` option used for GPUs does not apply to TPUs. TPUs work in DDP mode by default (distributing over each core)

34.7 TPU Pod

To train on more than 8 cores, your code actually doesn't change! All you need to do is submit the following command:

```
$ python -m torch_xla.distributed.xla_dist
--tpu=$TPU_POD_NAME
--conda-env=torch-xla-nightly
-- python /usr/share/torch-xla-0.5/pytorch/xla/test/test_train_imagenet.py --fake_data
```

See [this guide](#) on how to set up the instance groups and VMs needed to run TPU Pods.

34.8 16 bit precision

Lightning also supports training in 16-bit precision with TPUs. By default, TPU training will use 32-bit precision. To enable 16-bit, set the 16-bit flag.

```
import pytorch_lightning as pl

my_model = MyLightningModule()
trainer = pl.Trainer(tpu_cores=8, precision=16)
trainer.fit(my_model)
```

Under the hood the xla library will use the `bfloat16` type.

34.9 About XLA

XLA is the library that interfaces PyTorch with the TPUs. For more information check out [XLA](#).

Guide for [troubleshooting XLA](#)

TEST SET

Lightning forces the user to run the test set separately to make sure it isn't evaluated by mistake.

35.1 Test after fit

To run the test set after training completes, use this method.

```
# run full training
trainer.fit(model)

# (1) load the best checkpoint automatically (lightning tracks this for you)
trainer.test()

# (2) don't load a checkpoint, instead use the model with the latest weights
trainer.test(ckpt_path=None)

# (3) test using a specific checkpoint
trainer.test(ckpt_path='/path/to/my_checkpoint.ckpt')

# (4) test with an explicit model (will use this model and not load a checkpoint)
trainer.test(model)
```

35.2 Test multiple models

You can run the test set on multiple models using the same trainer instance.

```
model1 = LitModel()
model2 = GANModel()

trainer = Trainer()
trainer.test(model1)
trainer.test(model2)
```

35.3 Test pre-trained model

To run the test set on a pre-trained model, use this method.

```
model = MyLightningModule.load_from_checkpoint(  
    checkpoint_path='/path/to/pytorch_checkpoint.ckpt',  
    hparams_file='/path/to/test_tube/experiment/version/hparams.yaml',  
    map_location=None  
)  
  
# init trainer with whatever options  
trainer = Trainer(...)  
  
# test (pass in the model)  
trainer.test(model)
```

In this case, the options you pass to trainer will be used when running the test set (ie: 16-bit, dp, ddp, etc...)

35.4 Test with additional data loaders

You can still run inference on a test set even if the `test_dataloader` method hasn't been defined within your *Lightning-Module* instance. This would be the case when your test data is not available at the time your model was declared.

```
# setup your data loader  
test = DataLoader(...)  
  
# test (pass in the loader)  
trainer.test(test_dataloaders=test)
```

You can either pass in a single dataloader or a list of them. This optional named parameter can be used in conjunction with any of the above use cases.

INFERENCE IN PRODUCTION

PyTorch Lightning eases the process of deploying models into production.

36.1 Exporting to ONNX

PyTorch Lightning provides a handy function to quickly export your model to ONNX format, which allows the model to be independent of PyTorch and run on an ONNX Runtime.

To export your model to ONNX format call the `to_onnx` function on your Lightning Module with the filepath and `input_sample`.

```
filepath = 'model.onnx'
model = SimpleModel()
input_sample = torch.randn(1, 64)
model.to_onnx(filepath, input_sample, export_params=True)
```

You can also skip passing the input sample if the `example_input_array` property is specified in your LightningModule.

Once you have the exported model, you can run it on your ONNX runtime in the following way:

```
ort_session = onnxruntime.InferenceSession(filepath)
input_name = ort_session.get_inputs()[0].name
ort_inputs = {input_name: np.random.randn(1, 64).astype(np.float32)}
ort_outs = ort_session.run(None, ort_inputs)
```

36.2 Exporting to TorchScript

TorchScript allows you to serialize your models in a way that it can be loaded in non-Python environments. The LightningModule has a handy method `to_torchscript()` that returns a scripted module which you can save or directly use.

```
model = SimpleModel()
script = model.to_torchscript()

# save for use in production environment
torch.jit.save(script, "model.pt")
```

It is recommended that you install the latest supported version of PyTorch to use this feature without limitations.

ASR & TTS

These are amazing ecosystems to help with Automatic Speech Recognition (ASR) and Text to speech (TTS).

37.1 NeMo

NVIDIA NeMo is a toolkit for building new State-of-the-Art Conversational AI models. NeMo has separate collections for Automatic Speech Recognition (ASR), Natural Language Processing (NLP), and Text-to-Speech (TTS) models. Each collection consists of prebuilt modules that include everything needed to train on your data. Every module can easily be customized, extended, and composed to create new Conversational AI model architectures.

Conversational AI architectures are typically very large and require a lot of data and compute for training. NeMo uses PyTorch Lightning for easy and performant multi-GPU/multi-node mixed-precision training.

Note: Every NeMo model is a LightningModule that comes equipped with all supporting infrastructure for training and reproducibility.

37.1.1 NeMo Models

NeMo Models contain everything needed to train and reproduce state of the art Conversational AI research and applications, including:

- neural network architectures
- datasets/data loaders
- data preprocessing/postprocessing
- data augmentors
- optimizers and schedulers
- tokenizers
- language models

NeMo uses Hydra for configuring both NeMo models and the PyTorch Lightning Trainer. Depending on the domain and application, many different AI libraries will have to be configured to build the application. Hydra makes it easy to bring all of these libraries together so that each can be configured from .yaml or the Hydra CLI.

Note: Every NeMo model has an example configuration file and a corresponding script that contains all configurations needed for training.

The end result of using NeMo, Pytorch Lightning, and Hydra is that NeMo models all have the same look and feel. This makes it easy to do Conversational AI research across multiple domains. NeMo models are also fully compatible with the PyTorch ecosystem.

Installing NeMo

Before installing NeMo, please install Cython first.

```
pip install Cython
```

For ASR and TTS models, also install these linux utilities.

```
apt-get update && apt-get install -y libsndfile1 ffmpeg
```

Then installing the latest NeMo release is a simple pip install.

```
pip install nemo_toolkit[all]==1.0.0b1
```

To install the main branch from GitHub:

```
python -m pip install git+https://github.com/NVIDIA/NeMo.git@main#egg=nemo_
↳ toolkit[all]
```

To install from a local clone of NeMo:

```
./reinstall.sh # from cloned NeMo's git root
```

For Docker users, the NeMo container is available on [NGC](#).

```
docker pull nvcr.io/nvidia/nemo:v1.0.0b1
```

```
docker run --runtime=nvidia -it --rm -v --shm-size=8g -p 8888:8888 -p 6006:6006 --
↳ ulimit memlock=-1 --ulimit stack=67108864 nvcr.io/nvidia/nemo:1.0.0b1
```

Experiment Manager

NeMo's Experiment Manager leverages PyTorch Lightning for model checkpointing, TensorBoard Logging, and Weights and Biases logging. The Experiment Manager is included by default in all NeMo example scripts.

```
exp_manager(trainer, cfg.get("exp_manager", None))
```

And is configurable via .yaml with Hydra.

```
exp_manager:
  exp_dir: null
  name: *name
  create_tensorboard_logger: True
  create_checkpoint_callback: True
```

Optionally launch Tensorboard to view training results in `./nemo_experiments` (by default).

```
tensorboard --bind_all --logdir nemo_experiments
```

37.1.2 Automatic Speech Recognition (ASR)

Everything needed to train Convolutional ASR models is included with NeMo. NeMo supports multiple Speech Recognition architectures, including Jasper and QuartzNet. [NeMo Speech Models](#) can be trained from scratch on custom datasets or fine-tuned using pre-trained checkpoints trained on thousands of hours of audio that can be restored for immediate use.

Some typical ASR tasks are included with NeMo:

- [Audio transcription](#)
- [Byte Pair/Word Piece Training](#)
- [Speech Commands](#)
- [Voice Activity Detection](#)
- [Speaker Recognition](#)

See this [asr notebook](#) for a full tutorial on doing ASR with NeMo, PyTorch Lightning, and Hydra.

Specify ASR Model Configurations with YAML File

NeMo Models and the PyTorch Lightning Trainer can be fully configured from .yaml files using Hydra.

See this [asr config](#) for the entire speech to text .yaml file.

```
# configure the PyTorch Lightning Trainer
trainer:
  gpus: 0 # number of gpus
  max_epochs: 5
  max_steps: null # computed at runtime if not set
  num_nodes: 1
  distributed_backend: ddp
  ...
# configure the ASR model
model:
  ...
  encoder:
    _target_: nemo.collections.asr.modules.ConvASREncoder
    params:
      feat_in: *n_mels
      activation: relu
      conv_mask: true

    jasper:
      - filters: 128
        repeat: 1
        kernel: [11]
        stride: [1]
        dilation: [1]
        dropout: *dropout
      ...
```

(continues on next page)

(continued from previous page)

```
# all other configuration, data, optimizer, preprocessor, etc
...
```

Developing ASR Model From Scratch

speech_to_text.py

```
# hydra_runner calls hydra.main and is useful for multi-node experiments
@hydra_runner(config_path="conf", config_name="config")
def main(cfg):
    trainer = Trainer(**cfg.trainer)
    asr_model = EncDecCTCModel(cfg.model, trainer)
    trainer.fit(asr_model)
```

Hydra makes every aspect of the NeMo model, including the PyTorch Lightning Trainer, customizable from the command line.

```
python NeMo/examples/asr/speech_to_text.py --config-name=quartznet_15x5 \
    trainer.gpus=4 \
    trainer.max_epochs=128 \
    +trainer.precision=16 \
    model.train_ds.manifest_filepath=<PATH_TO_DATA>/librispeech-train-all.json \
    model.validation_ds.manifest_filepath=<PATH_TO_DATA>/librispeech-dev-other.json \
    model.train_ds.batch_size=64 \
    +model.validation_ds.num_workers=16 \
    +model.train_ds.num_workers=16
```

Note: Training NeMo ASR models can take days/weeks so it is highly recommended to use multiple GPUs and multiple nodes with the PyTorch Lightning Trainer.

Using State-Of-The-Art Pre-trained ASR Model

Transcribe audio with QuartzNet model pretrained on ~3300 hours of audio.

```
quartznet = EncDecCTCModel.from_pretrained('QuartzNet15x5Base-En')

files = ['path/to/my.wav'] # file duration should be less than 25 seconds

for fname, transcription in zip(files, quartznet.transcribe(paths2audio_files=files)):
    print(f"Audio in {fname} was recognized as: {transcription}")
```

To see the available pretrained checkpoints:

```
EncDecCTCModel.list_available_models()
```


NeMo ASR Model Under the Hood

Any aspect of ASR training or model architecture design can easily be customized with PyTorch Lightning since every NeMo model is a Lightning Module.

```
class EncDecCTCModel(ASRModel):
    """Base class for encoder decoder CTC-based models."""
    ...
    @typecheck()
    def forward(self, input_signal, input_signal_length):
        processed_signal, processed_signal_len = self.preprocessor(
            input_signal=input_signal, length=input_signal_length,
        )
        # Spec augment is not applied during evaluation/testing
        if self.spec_augmentation is not None and self.training:
            processed_signal = self.spec_augmentation(input_spec=processed_signal)
        encoded, encoded_len = self.encoder(audio_signal=processed_signal,
        →length=processed_signal_len)
        log_probs = self.decoder(encoder_output=encoded)
        greedy_predictions = log_probs.argmax(dim=-1, keepdim=False)
        return log_probs, encoded_len, greedy_predictions

    # PTL-specific methods
    def training_step(self, batch, batch_nb):
        audio_signal, audio_signal_len, transcript, transcript_len = batch
        log_probs, encoded_len, predictions = self.forward(
            input_signal=audio_signal, input_signal_length=audio_signal_len
        )
        loss_value = self.loss(
            log_probs=log_probs, targets=transcript, input_lengths=encoded_len,
        →target_lengths=transcript_len
        )
        wer_num, wer_denom = self._wer(predictions, transcript, transcript_len)
        tensorboard_logs = {
            'train_loss': loss_value,
            'training_batch_wer': wer_num / wer_denom,
            'learning_rate': self._optimizer.param_groups[0]['lr'],
        }
        return {'loss': loss_value, 'log': tensorboard_logs}
```

Neural Types in NeMo ASR

NeMo Models and Neural Modules come with Neural Type checking. Neural type checking is extremely useful when combining many different neural network architectures for a production-grade application.

```
@property
def input_types(self) -> Optional[Dict[str, NeuralType]]:
    if hasattr(self.preprocessor, '_sample_rate'):
        audio_eltype = AudioSignal(freq=self.preprocessor._sample_rate)
    else:
        audio_eltype = AudioSignal()
    return {
        "input_signal": NeuralType(('B', 'T'), audio_eltype),
        "input_signal_length": NeuralType(tuple('B'), LengthsType()),
    }
```

(continues on next page)

(continued from previous page)

```

@property
def output_types(self) -> Optional[Dict[str, NeuralType]]:
    return {
        "outputs": NeuralType(('B', 'T', 'D'), LogprobsType()),
        "encoded_lengths": NeuralType(tuple('B'), LengthsType()),
        "greedy_predictions": NeuralType(('B', 'T'), LabelsType()),
    }

```

37.1.3 Natural Language Processing (NLP)

Everything needed to finetune BERT-like language models for NLP tasks is included with NeMo. [NeMo NLP Models](#) include [HuggingFace Transformers](#) and [NVIDIA Megatron-LM](#) BERT and Bio-Megatron models. NeMo can also be used for pretraining BERT-based language models from HuggingFace.

Any of the HuggingFace encoders or Megatron-LM encoders can easily be used for the NLP tasks that are included with NeMo:

- [Glue Benchmark \(All tasks\)](#)
- [Intent Slot Classification](#)
- [Language Modeling \(BERT Pretraining\)](#)
- [Question Answering](#)
- [Text Classification \(including Sentiment Analysis\)](#)
- [Token Classification \(including Named Entity Recognition\)](#)
- [Punctuation and Capitalization](#)

Named Entity Recognition (NER)

NER (or more generally token classification) is the NLP task of detecting and classifying key information (entities) in text. This task is very popular in Healthcare and Finance. In finance, for example, it can be important to identify geographical, geopolitical, organizational, persons, events, and natural phenomenon entities. See this [NER notebook](#) for a full tutorial on doing NER with NeMo, PyTorch Lightning, and Hydra.

Specify NER Model Configurations with YAML File

..note NeMo Models and the PyTorch Lightning Trainer can be fully configured from .yaml files using Hydra.

See this [token classification config](#) for the entire NER (token classification) .yaml file.

```

# configure any argument of the PyTorch Lightning Trainer
trainer:
    gpus: 1 # the number of gpus, 0 for CPU
    num_nodes: 1
    max_epochs: 5
    ...
# configure any aspect of the token classification model here
model:
    dataset:
        data_dir: ??? # /path/to/data

```

(continues on next page)

(continued from previous page)

```

    class_balancing: null # choose from [null, weighted_loss]. Weighted_loss_
↳enables the weighted class balancing of the loss, may be used for handling_
↳unbalanced classes
    max_seq_length: 128
    ...
tokenizer:
    tokenizer_name: ${model.language_model.pretrained_model_name} # or sentencepiece
    vocab_file: null # path to vocab file
    ...
# the language model can be from HuggingFace or Megatron-LM
language_model:
    pretrained_model_name: bert-base-uncased
    lm_checkpoint: null
    ...
# the classifier for the downstream task
head:
    num_fc_layers: 2
    fc_dropout: 0.5
    activation: 'relu'
    ...
# all other configuration: train/val/test/ data, optimizer, experiment manager, etc
...

```

Developing NER Model From Scratch

token_classification.py

```

# hydra_runner calls hydra.main and is useful for multi-node experiments
@hydra_runner(config_path="conf", config_name="token_classification_config")
def main(cfg: DictConfig) -> None:
    trainer = pl.Trainer(**cfg.trainer)
    model = TokenClassificationModel(cfg.model, trainer=trainer)
    trainer.fit(model)

```

After training, we can do inference with the saved NER model using PyTorch Lightning.

Inference from file:

```

gpu = 1 if cfg.trainer.gpus != 0 else 0
trainer = pl.Trainer(gpus=gpu)
model.set_trainer(trainer)
model.evaluate_from_file(
    text_file=os.path.join(cfg.model.dataset.data_dir, cfg.model.validation_ds.text_
↳file),
    labels_file=os.path.join(cfg.model.dataset.data_dir, cfg.model.validation_ds.
↳labels_file),
    output_dir=exp_dir,
    add_confusion_matrix=True,
    normalize_confusion_matrix=True,
)

```

Or we can run inference on a few examples:

```

queries = ['we bought four shirts from the nvidia gear store in santa clara.',
↳'Nvidia is a company in Santa Clara.']

```

(continues on next page)

(continued from previous page)

```
results = model.add_predictions(queries)

for query, result in zip(queries, results):
    logging.info(f'Query : {query}')
    logging.info(f'Result: {result.strip()}\n')
```

Hydra makes every aspect of the NeMo model, including the PyTorch Lightning Trainer, customizable from the command line.

```
python token_classification.py \
    model.language_model.pretrained_model_name=bert-base-cased \
    model.head.num_fc_layers=2 \
    model.dataset.data_dir=/path/to/my/data \
    trainer.max_epochs=5 \
    trainer.gpus=[0,1]
```

37.1.4 Tokenizers

Tokenization is the process of converting natural language text into integer arrays which can be used for machine learning. For NLP tasks, tokenization is an essential part of data preprocessing. NeMo supports all BERT-like model tokenizers from [HuggingFace's AutoTokenizer](#) and also supports [Google's SentencePieceTokenizer](#) which can be trained on custom data.

To see the list of supported tokenizers:

```
from nemo.collections import nlp as nemo_nlp

nemo_nlp.modules.get_tokenizer_list()
```

See this [tokenizer notebook](#) for a full tutorial on using tokenizers in NeMo.

Language Models

Language models are used to extract information from (tokenized) text. Much of the state-of-the-art in natural language processing is achieved by fine-tuning pretrained language models on the downstream task.

With NeMo, you can either [pretrain](#) a BERT model on your data or use a pretrained language model from [HuggingFace Transformers](#) or [NVIDIA Megatron-LM](#).

To see the list of language models available in NeMo:

```
nemo_nlp.modules.get_pretrained_lm_models_list(include_external=True)
```

Easily switch between any language model in the above list by using `.get_lm_model`.

```
nemo_nlp.modules.get_lm_model(pretrained_model_name='distilbert-base-uncased')
```

See this [language model notebook](#) for a full tutorial on using pretrained language models in NeMo.

Using a Pre-trained NER Model

NeMo has pre-trained NER models that can be used to get started with Token Classification right away. Models are automatically downloaded from NGC, cached locally to disk, and loaded into GPU memory using the `.from_pretrained` method.

```
# load pre-trained NER model
pretrained_ner_model = TokenClassificationModel.from_pretrained(model_name="NERModel")

# define the list of queries for inference
queries = [
    'we bought four shirts from the nvidia gear store in santa clara.',
    'Nvidia is a company.',
    'The Adventures of Tom Sawyer by Mark Twain is an 1876 novel about a young boy_
↳growing '
    + 'up along the Mississippi River.',
]
results = pretrained_ner_model.add_predictions(queries)

for query, result in zip(queries, results):
    print()
    print(f'Query : {query}')
    print(f'Result: {result.strip()}\\n')
```

NeMo NER Model Under the Hood

Any aspect of NLP training or model architecture design can easily be customized with PyTorch Lightning since every NeMo model is a Lightning Module.

```
class TokenClassificationModel(ModelPT):
    """
    Token Classification Model with BERT, applicable for tasks such as Named Entity_
↳Recognition
    """
    ...
    @typecheck()
    def forward(self, input_ids, token_type_ids, attention_mask):
        hidden_states = self.bert_model(
            input_ids=input_ids, token_type_ids=token_type_ids, attention_
↳mask=attention_mask
        )
        logits = self.classifier(hidden_states=hidden_states)
        return logits

    # PTL-specific methods
    def training_step(self, batch, batch_idx):
        """
        Lightning calls this inside the training loop with the data from the training_
↳dataloader
        passed in as `batch`.
        """
        input_ids, input_type_ids, input_mask, subtokens_mask, loss_mask, labels = _
↳batch
        logits = self(input_ids=input_ids, token_type_ids=input_type_ids, attention_
↳mask=input_mask)
```

(continues on next page)

(continued from previous page)

```

        loss = self.loss(logits=logits, labels=labels, loss_mask=loss_mask)
        tensorboard_logs = {'train_loss': loss, 'lr': self._optimizer.param_groups[0][
↪ 'lr']}
        return {'loss': loss, 'log': tensorboard_logs}
    ...

```

Neural Types in NeMo NLP

NeMo Models and Neural Modules come with Neural Type checking. Neural type checking is extremely useful when combining many different neural network architectures for a production-grade application.

```

@property
def input_types(self) -> Optional[Dict[str, NeuralType]]:
    return self.bert_model.input_types

@property
def output_types(self) -> Optional[Dict[str, NeuralType]]:
    return self.classifier.output_types

```

37.1.5 Text-To-Speech (TTS)

Everything needed to train TTS models and generate audio is included with NeMo. [NeMo TTS Models](#) can be trained from scratch on your own data or pretrained models can be downloaded automatically. NeMo currently supports a two step inference procedure. First, a model is used to generate a mel spectrogram from text. Second, a model is used to generate audio from a mel spectrogram.

Mel Spectrogram Generators:

- [Tacotron 2](#)
- [Glow-TTS](#)

Audio Generators:

- [Griffin-Lim](#)
- [WaveGlow](#)
- [SqueezeWave](#)

Specify TTS Model Configurations with YAML File

..note NeMo Models and PyTorch Lightning Trainer can be fully configured from .yaml files using Hydra.

tts/conf/glow_tts.yaml

```

# configure the PyTorch Lightning Trainer
trainer:
    gpus: -1 # number of gpus
    max_epochs: 350
    num_nodes: 1
    distributed_backend: ddp
    ...

```

(continues on next page)

(continued from previous page)

```
# configure the TTS model
model:
    ...
    encoder:
        _target_: nemo.collections.tts.modules.glow_tts.TextEncoder
        params:
            n_vocab: 148
            out_channels: *n_mels
            hidden_channels: 192
            filter_channels: 768
            filter_channels_dp: 256
        ...
# all other configuration, data, optimizer, parser, preprocessor, etc
...
```

Developing TTS Model From Scratch

tts/glow_tts.py

```
# hydra_runner calls hydra.main and is useful for multi-node experiments
@hydra_runner(config_path="conf", config_name="glow_tts")
def main(cfg):
    trainer = pl.Trainer(**cfg.trainer)
    model = GlowTTSModel(cfg=cfg.model, trainer=trainer)
    trainer.fit(model)
```

Hydra makes every aspect of the NeMo model, including the PyTorch Lightning Trainer, customizable from the command line.

```
python NeMo/examples/tts/glow_tts.py \
    trainer.gpus=4 \
    trainer.max_epochs=400 \
    ...
    train_dataset=/path/to/train/data \
    validation_datasets=/path/to/val/data \
    model.train_ds.batch_size = 64 \
```

..note Training NeMo TTS models from scratch take days/weeks so it is highly recommended to use multiple GPUs and multiple nodes with the PyTorch Lightning Trainer.

Using State-Of-The-Art Pre-trained TTS Model

Generate speech using models trained on *LJSpeech* <<https://keithito.com/LJ-Speech-Dataset/>>, around 24 hours of single speaker data.

See this [TTS notebook](#) for a full tutorial on generating speech with NeMo, PyTorch Lightning, and Hydra.

```
# load pretrained spectrogram model
spec_gen = SpecModel.from_pretrained('GlowTTS-22050Hz').cuda()

# load pretrained Generators
vocoder = WaveGlowModel.from_pretrained('WaveGlow-22050Hz').cuda()
```

(continues on next page)

(continued from previous page)

```
def infer(spec_gen_model, vocoder_model, str_input):
    with torch.no_grad():
        parsed = spec_gen.parse(text_to_generate)
        spectrogram = spec_gen.generate_spectrogram(tokens=parsed)
        audio = vocoder.convert_spectrogram_to_audio(spec=spectrogram)
        if isinstance(spectrogram, torch.Tensor):
            spectrogram = spectrogram.to('cpu').numpy()
        if len(spectrogram.shape) == 3:
            spectrogram = spectrogram[0]
        if isinstance(audio, torch.Tensor):
            audio = audio.to('cpu').numpy()
        return spectrogram, audio
```

```
text_to_generate = input("Input what you want the model to say: ")
spec, audio = infer(spec_gen, vocoder, text_to_generate)
```

To see the available pretrained checkpoints:

```
# spec generator
GlowTTSModel.list_available_models()

# vocoder
WaveGlowModel.list_available_models()
```

NeMo TTS Model Under the Hood

Any aspect of TTS training or model architecture design can easily be customized with PyTorch Lightning since every NeMo model is a LightningModule.

glow_tts.py

```
class GlowTTSModel(SpectrogramGenerator):
    """
    GlowTTS model used to generate spectrograms from text
    Consists of a text encoder and an invertible spectrogram decoder
    """
    ...
    # NeMo models come with neural type checking
    @typecheck(
        input_types={
            "x": NeuralType(('B', 'T'), TokenIndex()),
            "x_lengths": NeuralType(('B'), LengthsType()),
            "y": NeuralType(('B', 'D', 'T'), MelSpectrogramType(), optional=True),
            "y_lengths": NeuralType(('B'), LengthsType(), optional=True),
            "gen": NeuralType(optional=True),
            "noise_scale": NeuralType(optional=True),
            "length_scale": NeuralType(optional=True),
        }
    )
    def forward(self, *, x, x_lengths, y=None, y_lengths=None, gen=False, noise_
↳ scale=0.3, length_scale=1.0):
        if gen:
            return self.glow_tts.generate_spect(
                text=x, text_lengths=x_lengths, noise_scale=noise_scale, length_
↳ scale=length_scale
```

(continues on next page)

(continued from previous page)

```

        )
    else:
        return self.glow_tts(text=x, text_lengths=x_lengths, spect=y, spect_
↪lengths=y_lengths)
    ...
    def step(self, y, y_lengths, x, x_lengths):
        z, y_m, y_logs, logdet, logw, logw_, y_lengths, attn = self(
            x=x, x_lengths=x_lengths, y=y, y_lengths=y_lengths, gen=False
        )

        l_mle, l_length, logdet = self.loss(
            z=z,
            y_m=y_m,
            y_logs=y_logs,
            logdet=logdet,
            logw=logw,
            logw_=logw_,
            x_lengths=x_lengths,
            y_lengths=y_lengths,
        )

        loss = sum([l_mle, l_length])

        return l_mle, l_length, logdet, loss, attn

    # PTL-specific methods
    def training_step(self, batch, batch_idx):
        y, y_lengths, x, x_lengths = batch

        y, y_lengths = self.preprocessor(input_signal=y, length=y_lengths)

        l_mle, l_length, logdet, loss, _ = self.step(y, y_lengths, x, x_lengths)

        output = {
            "loss": loss, # required
            "progress_bar": {"l_mle": l_mle, "l_length": l_length, "logdet": logdet},
            "log": {"loss": loss, "l_mle": l_mle, "l_length": l_length, "logdet": ↪
↪logdet},
        }

        return output
    ...

```

Neural Types in NeMo TTS

NeMo Models and Neural Modules come with Neural Type checking. Neural type checking is extremely useful when combining many different neural network architectures for a production-grade application.

```

@typecheck(
    input_types={
        "x": NeuralType(('B', 'T'), TokenIndex()),
        "x_lengths": NeuralType(('B'), LengthsType()),
        "y": NeuralType(('B', 'D', 'T'), MelSpectrogramType(), optional=True),
        "y_lengths": NeuralType(('B'), LengthsType(), optional=True),
        "gen": NeuralType(optional=True),
    }
)

```

(continues on next page)

(continued from previous page)

```
        "noise_scale": NeuralType(optional=True),
        "length_scale": NeuralType(optional=True),
    }
)
def forward(self, *, x, x_lengths, y=None, y_lengths=None, gen=False, noise_scale=0.3,
    ↪ length_scale=1.0):
    ...
```

37.1.6 Learn More

Download pre-trained [ASR](#), [NLP](#), and [TTS](#) models on [NVIDIA NGC](#) to quickly get started with NeMo.

Become an expert on Building Conversational AI applications with our [tutorials](#), and [example scripts](#),

Note: Most NeMo tutorial notebooks can be run on [Google Colab](#).

[NVIDIA NeMo](#) is actively being developed on GitHub. [Contributions](#) are welcome!

See our [developer guide](#) for more information on core NeMo concepts, ASR/NLP/TTS collections, and the NeMo API.

CONTRIBUTOR COVENANT CODE OF CONDUCT

38.1 Our Pledge

In the interest of fostering an open and welcoming environment, we as contributors and maintainers pledge to making participation in our project and our community a harassment-free experience for everyone, regardless of age, body size, disability, ethnicity, sex characteristics, gender identity and expression, level of experience, education, socio-economic status, nationality, personal appearance, race, religion, or sexual identity and orientation.

38.2 Our Standards

Examples of behavior that contributes to creating a positive environment include:

- Using welcoming and inclusive language
- Being respectful of differing viewpoints and experiences
- Gracefully accepting constructive criticism
- Focusing on what is best for the community
- Showing empathy towards other community members

Examples of unacceptable behavior by participants include:

- The use of sexualized language or imagery and unwelcome sexual attention or advances
- Trolling, insulting/derogatory comments, and personal or political attacks
- Public or private harassment
- Publishing others' private information, such as a physical or electronic address, without explicit permission
- Other conduct which could reasonably be considered inappropriate in a professional setting

38.3 Our Responsibilities

Project maintainers are responsible for clarifying the standards of acceptable behavior and are expected to take appropriate and fair corrective action in response to any instances of unacceptable behavior.

Project maintainers have the right and responsibility to remove, edit, or reject comments, commits, code, wiki edits, issues, and other contributions that are not aligned to this Code of Conduct, or to ban temporarily or permanently any contributor for other behaviors that they deem inappropriate, threatening, offensive, or harmful.

38.4 Scope

This Code of Conduct applies both within project spaces and in public spaces when an individual is representing the project or its community. Examples of representing a project or community include using an official project e-mail address, posting via an official social media account, or acting as an appointed representative at an online or offline event. Representation of a project may be further defined and clarified by project maintainers.

38.5 Enforcement

Instances of abusive, harassing, or otherwise unacceptable behavior may be reported by contacting the project team at waf2107@columbia.edu. All complaints will be reviewed and investigated and will result in a response that is deemed necessary and appropriate to the circumstances. The project team is obligated to maintain confidentiality with regard to the reporter of an incident. Further details of specific enforcement policies may be posted separately.

Project maintainers who do not follow or enforce the Code of Conduct in good faith may face temporary or permanent repercussions as determined by other members of the project's leadership.

38.6 Attribution

This Code of Conduct is adapted from the [Contributor Covenant](https://www.contributor-covenant.org/version/1/4/code-of-conduct.html), version 1.4, available at <https://www.contributor-covenant.org/version/1/4/code-of-conduct.html>

For answers to common questions about this code of conduct, see <https://www.contributor-covenant.org/faq>

CONTRIBUTING

Welcome to the PyTorch Lightning community! We're building the most advanced research platform on the planet to implement the latest, best practices that the amazing PyTorch team rolls out!

39.1 Main Core Value: One less thing to remember

Simplify the API as much as possible from the user perspective. Any additions or improvements should minimize the things the user needs to remember.

For example: One benefit of the `validation_step` is that the user doesn't have to remember to set the model to `.eval()`. This helps users avoid all sorts of subtle errors.

39.2 Lightning Design Principles

We encourage all sorts of contributions you're interested in adding! When coding for lightning, please follow these principles.

39.2.1 No PyTorch Interference

We don't want to add any abstractions on top of pure PyTorch. This gives researchers all the control they need without having to learn yet another framework.

39.2.2 Simple Internal Code

It's useful for users to look at the code and understand very quickly what's happening. Many users won't be engineers. Thus we need to value clear, simple code over condensed ninja moves. While that's super cool, this isn't the project for that :)

39.2.3 Force User Decisions To Best Practices

There are 1,000 ways to do something. However, eventually one popular solution becomes standard practice, and everyone follows. We try to find the best way to solve a particular problem, and then force our users to use it for readability and simplicity. A good example is accumulated gradients. There are many different ways to implement it, we just pick one and force users to use it. A bad forced decision would be to make users use a specific library to do something.

When something becomes a best practice, we add it to the framework. This is usually something like bits of code in utils or in the model file that everyone keeps adding over and over again across projects. When this happens, bring that code inside the trainer and add a flag for it.

39.2.4 Simple External API

What makes sense to you may not make sense to others. When creating an issue with an API change suggestion, please validate that it makes sense for others. Treat code changes the way you treat a startup: validate that it's a needed feature, then add if it makes sense for many people.

39.2.5 Backward-compatible API

We all hate updating our deep learning packages because we don't want to refactor a bunch of stuff. In Lightning, we make sure every change we make which could break an API is backward compatible with good deprecation warnings.

You shouldn't be afraid to upgrade Lightning :)

39.2.6 Gain User Trust

As a researcher, you can't have any part of your code going wrong. So, make thorough tests to ensure that every implementation of a new trick or subtle change is correct.

39.2.7 Interoperability

Have a favorite feature from other libraries like fast.ai or transformers? Those should just work with lightning as well. Grab your favorite model or learning rate scheduler from your favorite library and run it in Lightning.

39.3 Contribution Types

We are always looking for help implementing new features or fixing bugs.

A lot of good work has already been done in project mechanics (requirements/base.txt, setup.py, pep8, badges, ci, etc. . .) so we're in a good state there thanks to all the early contributors (even pre-beta release)!

39.3.1 Bug Fixes:

1. If you find a bug please submit a github issue.
 - Make sure the title explains the issue.
 - Describe your setup, what you are trying to do, expected vs. actual behaviour. Please add configs and code samples.
 - Add details on how to reproduce the issue - a minimal test case is always best, colab is also great. Note, that the sample code shall be minimal and if needed with publicly available data.
2. Try to fix it or recommend a solution. We highly recommend to use test-driven approach:
 - Convert your minimal code example to a unit/integration test with assert on expected results.
 - Start by debugging the issue... You can run just this particular test in your IDE and draft a fix.
 - Verify that your test case fails on the master branch and only passes with the fix applied.
3. Submit a PR!

Note, even if you do not find the solution, sending a PR with a test covering the issue is a valid contribution and we can help you or finish it with you :]

39.3.2 New Features:

1. Submit a github issue - describe what is the motivation of such feature (adding the use case or an example is helpful).
2. Let's discuss to determine the feature scope.
3. Submit a PR! We recommend test driven approach to adding new features as well:
 - Write a test for the functionality you want to add.
 - Write the functional code until the test passes.
4. Add/update the relevant tests!
 - [This PR](#) is a good example for adding a new metric, and [this one](#) for a new logger.

39.3.3 Test cases:

Want to keep Lightning healthy? Love seeing those green tests? So do we! How to we keep it that way? We write tests! We value tests contribution even more than new features.

Most of the tests in PyTorch Lightning train a trial MNIST model under various trainer conditions (ddp, ddp2+amp, etc...). The tests expect the model to perform to a reasonable degree of testing accuracy to pass. Want to add a new test case and not sure how? [Talk to us!](#)

39.4 Guidelines

39.4.1 Original code

All added or edited code shall be the own original work of the particular contributor. If you use some third-party implementation, all such blocks/functions/modules shall be properly referred and if possible also agreed by code's author. For example - This code is inspired from `http://...`. In case you adding new dependencies, make sure that they are compatible with the actual PyTorch Lightning license (ie. dependencies should be *at least* as permissive as the PyTorch Lightning license).

39.4.2 Coding Style

1. Use f-strings for output formation (except logging when we stay with lazy `logging.info("Hello %s!", name)`);
2. Black code formatter is used using `pre-commit` hook.

39.4.3 Documentation

We are using Sphinx with Napoleon extension. Moreover, we set Google style to follow with type convention.

- Napoleon formatting with Google style
- ReStructured Text (reST)
- Paragraph-level markup

See following short example of a sample function taking one position string and optional

```
from typing import Optional

def my_func(param_a: int, param_b: Optional[float] = None) -> str:
    """Sample function.

    Args:
        param_a: first parameter
        param_b: second parameter

    Return:
        sum of both numbers

    Example:
        Sample doctest example...
        >>> my_func(1, 2)
        3

    .. note:: If you want to add something.
    """
    p = param_b if param_b else 0
    return str(param_a + p)
```

When updating the docs make sure to build them first locally and visually inspect the html files (in the browser) for formatting errors. In certain cases, a missing blank line or a wrong indent can lead to a broken layout. Run these commands


```
pip install -r requirements/docs.txt
cd docs
make html
```

and open `docs/build/html/index.html` in your browser.

Notes:

- You need to have LaTeX installed for rendering math equations. You can for example install TeXLive by doing one of the following:
 - on Ubuntu (Linux) run `apt-get install texlive` or otherwise follow the instructions on the TeXLive website
 - use the [RTD docker image](#)
- with PL used class meta you need to use python 3.7 or higher

When you send a PR the continuous integration will run tests and build the docs. You can access a preview of the html pages in the *Artifacts* tab in CircleCI when you click on the task named *ci/circleci: Build-Docs* at the bottom of the PR page.

39.4.4 Testing

Local: Testing your work locally will help you speed up the process since it allows you to focus on particular (failing) test-cases. To setup a local development environment, install both local and test dependencies:

```
python -m pip install ".[dev, examples]"
python -m pip install pre-commit
```

You can run the full test-case in your terminal via this bash script:

```
bash .run_local_tests.sh
```

Note: if your computer does not have multi-GPU nor TPU these tests are skipped.

GitHub Actions: For convenience, you can also use your own GHActions building which will be triggered with each commit. This is useful if you do not test against all required dependency versions.

Docker: Another option is utilize the [pytorch lightning cuda base docker image](#). You can then run:

```
python -m pytest pytorch_lightning tests pl_examples -v --flake8
```

39.4.5 Pull Request

We welcome any useful contribution! For your convenience here's a recommended workflow:

1. Think about what you want to do - fix a bug, repair docs, etc. If you want to implement a new feature or enhance an existing one, start by opening a GitHub issue to explain the feature and the motivation. Members from core-contributors will take a look (it might take some time - we are often overloaded with issues!) and discuss it. Once an agreement was reached - start coding.
2. Start your work locally (usually until you need our CI testing).
 - Create a branch and prepare your changes.
 - Tip: do not work with your master directly, it may become complicated when you need to rebase.
 - Tip: give your PR a good name! It will be useful later when you may work on multiple tasks/PRs.

3. Test your code!
 - It is always good practice to start coding by creating a test case, verifying it breaks with current behaviour, and passes with your new changes.
 - Make sure your new tests cover all different edge cases.
 - Make sure all exceptions are handled.
4. Create a “Draft PR” which is clearly marked, to let us know you don’t need feedback yet.
5. When you feel ready for integrating your work, mark your PR “Ready for review”.
 - Your code should be readable and follow the project’s design principles.
 - Make sure all tests are passing.
 - Make sure you add a GitHub issue to your PR.
6. Use tags in PR name for following cases:
 - **[blocked by #]** if your work is depending on others changes.
 - **[wip]** when you start to re-edit your work, mark it so no one will accidentally merge it in meantime.

39.4.6 Question & Answer

1. How can I help/contribute?

All help is extremely welcome - reporting bugs, fixing documentation, adding test cases, solving issues and preparing bug fixes. To solve some issues you can start with label [good first issue](#) or chose something close to your domain with label [help wanted](#). Before you start to implement anything check that the issue description that it is clear and self-assign the task to you (if it is not possible, just comment that you take it and we assign it to you...).

2. Is there a recommendation for branch names?

We do not rely on the name convention so far you are working with your own fork. Anyway it would be nice to follow this convention `<type>/<issue-id>_<short-name>` where the types are: `bugfix`, `feature`, `docs`, `tests`, ...

3. How to rebase my PR?

We recommend creating a PR in separate branch other than `master`, especially if you plan submitting several changes and do not want to wait until the first one is resolved (we can work on them in parallel).

First, make sure you have set [upstream](#) by running:

```
git remote add upstream https://github.com/PyTorchLightning/pytorch-lightning.git
```

You’ll know its set up right if you run `git remote -v` and see something similar to this:

```
origin  https://github.com/{YOUR_USERNAME}/pytorch-lightning.git (fetch)
origin  https://github.com/{YOUR_USERNAME}/pytorch-lightning.git (push)
upstream      https://github.com/PyTorchLightning/pytorch-lightning.git (fetch)
upstream      https://github.com/PyTorchLightning/pytorch-lightning.git (push)
```

Now you can update your master with upstream’s master by running:

```
git fetch --all --prune
git checkout master
git merge upstream/master
```

Finally, checkout your feature branch and rebase it with master before pushing up your feature branch:

```
git checkout my-PR-branch
git rebase master
# follow git instructions to resolve conflicts
git push -f
```

Eventually, you can perform the rebasing directly from upstream after setting it up:

```
git fetch --all --prune
git rebase upstream/master
# follow git instructions to resolve conflicts
git push -f
```

39.4.7 Bonus Workflow Tip

If you don't want to remember all the commands above every time you want to push some code/setup a Lightning Dev environment on a new VM, you can set up bash aliases for some common commands. You can add these to one of your `~/.bashrc`, `~/.zshrc`, or `~/.bash_aliases` files.

NOTE: Once you edit one of these files, remember to source it or restart your shell. (ex. `source ~/.bashrc` if you added these to your `~/.bashrc` file).

```
plclone (){
    git clone https://github.com/{YOUR_USERNAME}/pytorch-lightning.git
    cd pytorch-lightning
    git remote add upstream https://github.com/PyTorchLightning/pytorch-lightning.git
    # This is just here to print out info about your remote upstream/origin
    git remote -v
}

plfetch (){
    git fetch --all --prune
    git checkout master
    git merge upstream/master
}

# Rebase your branch with upstream's master
# plrebase <your-branch-name>
plrebase (){
    git checkout $@
    git rebase master
}
```

Now, you can:

- clone your fork and set up upstream by running `plclone` from your terminal
- fetch upstream and update your local master branch with it by running `plfetch`
- rebase your feature branch (after running `plfetch`) by running `plrebase your-branch-name`

HOW TO BECOME A CORE CONTRIBUTOR

Thanks for your interest in joining the Lightning team! We're a rapidly growing project which is poised to become the go-to framework for DL researchers! We're currently recruiting for a team of 5 core maintainers.

As a core maintainer you will have a strong say in the direction of the project. Big changes will require a majority of maintainers to agree.

40.1 Code of conduct

First and foremost, you'll be evaluated against [these core values](#). Any code we commit or feature we add needs to align with those core values.

40.2 The bar for joining the team

Lightning is being used to solve really hard problems at the top AI labs in the world. As such, the bar for adding team members is extremely high. Candidates must have solid engineering skills, have a good eye for user experience, and must be a power user of Lightning and PyTorch.

With that said, the Lightning team will be diverse and a reflection of an inclusive AI community. You don't have to be an engineer to contribute! Scientists with great usability intuition and PyTorch ninja skills are welcomed!

40.3 Responsibilities:

The responsibilities mainly revolve around 3 things.

40.3.1 Github issues

- Here we want to help users have an amazing experience. These range from questions from new people getting into DL to questions from researchers about doing something esoteric with Lightning. Often, these issues require some sort of bug fix, document clarification or new functionality to be scoped out.
- To become a core member you must resolve at least 10 Github issues which align with the API design goals for Lightning. By the end of these 10 issues I should feel comfortable in the way you answer user questions. Pleasant/helpful tone.
- Can abstract from that issue or bug into functionality that might solve other related issues or makes the platform more flexible.

- Don't make users feel like they don't know what they're doing. We're here to help and to make everyone's experience delightful.

40.3.2 Pull requests

- Here we need to ensure the code that enters Lightning is high quality. For each PR we need to:
- Make sure code coverage does not decrease
- Documents are updated
- Code is elegant and simple
- Code is NOT overly engineered or hard to read
- Ask yourself, could a non-engineer understand what's happening here?
- Make sure new tests are written
- Is this NECESSARY for Lightning? There are some PRs which are just purely about adding engineering complexity which have no place in Lightning. Guidance
- Some other PRs are for people who are wanting to get involved and add something unnecessary. We do want their help though! So don't approve the PR, but direct them to a Github issue that they might be interested in helping with instead!
- To be considered for core contributor, please review 10 PRs and help the authors land it on master. Once you've finished the review, ping me for a sanity check. At the end of 10 PRs if your PR reviews are inline with expectations described above, then you can merge PRs on your own going forward, otherwise we'll do a few more until we're both comfortable :)

40.3.3 Project directions

There are some big decisions which the project must make. For these I expect core contributors to have something meaningful to add if it's their area of expertise.

40.3.4 Diversity

Lightning should reflect the broader community it serves. As such we should have scientists/researchers from different fields contributing!

The first 5 core contributors will fit this profile. Thus if you overlap strongly with experiences and expertise as someone else on the team, you might have to wait until the next set of contributors are added.

40.3.5 Summary: Requirements to apply

- Solve 10 Github issues. The goal is to be inline with expectations for solving issues by the last one so you can do them on your own. If not, I might ask you to solve a few more specific ones.
- Do 10 PR reviews. The goal is to be inline with expectations for solving issues by the last one so you can do them on your own. If not, I might ask you to solve a few more specific ones.

If you want to be considered, ping me on gitter and start [tracking your progress here](#).

PYTORCH LIGHTNING GOVERNANCE | PERSONS OF INTEREST

41.1 Leads

- William Falcon ([williamFalcon](#)) (Lightning founder)
- Jirka Borovec ([Borda](#))
- Ethan Harris ([ethanwharris](#)) (Torchbearer founder)
- Matthew Painter ([MattPainter01](#)) (Torchbearer founder)
- Justus Schock ([justusschock](#)) (Former Core Member PyTorch Ignite)

41.2 Core Maintainers

- Nic Eggert ([neggert](#))
- Jeff Ling ([jeffling](#))
- Jeremy Jordan ([jeremyjordan](#))
- Tullie Murrell ([tullie](#))
- Adrian Wälchli ([awaelchli](#))
- Nicki Skafte ([skaftenicki](#))
- Peter Yu ([yukw777](#))
- Rohit Gupta ([rohitgr7](#))
- Lezwon Castelino ([lezwon](#))
- Jeff Yang ([ydcjeff](#))
- Roger Shieh ([s-rog](#))

CHANGELOG

All notable changes to this project will be documented in this file.

The format is based on [Keep a Changelog](#).

42.1 [unreleased] - YYYY-MM-DD

42.1.1 [unreleased] - Added

42.1.2 [unreleased] - Changed

42.1.3 [unreleased] - Deprecated

42.1.4 [unreleased] - Removed

42.1.5 [unreleased] - Fixed

- Fixed hparams saving - save the state when `save_hyperparameters()` is called [in `__init__`] (#4163)
- Fixed runtime failure while exporting hparams to yaml (#4158)

42.2 [1.0.1] - 2020-10-14

42.2.1 [1.0.1] - Added

- Added `getstate/setstate` method for `torch.save` serialization (#4127)

42.3 [1.0.0] - 2020-10-13

42.3.1 [1.0.0] - Added

- Added Explained Variance Metric + metric fix (#4013)
- Added Metric <-> Lightning Module integration tests (#4008)
- Added parsing OS env vars in Trainer (#4022)
- Added classification metrics (#4043)
- Updated explained variance metric (#4024)
- Enabled plugins (#4041)
- Enabled custom clusters (#4048)
- Enabled passing in custom accelerators (#4050)
- Added `LightningModule.toggle_optimizer` (#4058)
- Added `LightningModule.manual_backward` (#4063)

42.3.2 [1.0.0] - Changed

- Integrated metrics API with `self.log` (#3961)
- Decoupled Apex (#4052, #4054, #4055, #4056, #4058, #4060, #4061, #4062, #4063, #4064, #4065)
- Renamed all backends to `Accelerator` (#4066)
- Enabled manual returns (#4089)

42.3.3 [1.0.0] - Removed

- Removed output argument from `*_batch_end` hooks (#3965, #3966)
- Removed output argument from `*_epoch_end` hooks (#3967)
- Removed support for `EvalResult` and `TrainResult` (#3968)
- Removed deprecated trainer flags: `overfit_pct`, `log_save_interval`, `row_log_interval` (#3969)
- Removed deprecated `early_stop_callback` (#3982)
- Removed deprecated model hooks (#3980)
- Removed deprecated callbacks (#3979)
- Removed `trainer` argument in `LightningModule.backward` #4056)

42.3.4 [1.0.0] - Fixed

- Fixed `current_epoch` property update to reflect true epoch number inside `LightningDataModule`, when `reload_dataloaders_every_epoch=True`. (#3974)
- Fixed to print scaler value in progress bar (#4053)
- Fixed mismatch between docstring and code regarding when `on_load_checkpoint` hook is called (#3996)

42.4 [0.10.0] - 2020-10-07

42.4.1 [0.10.0] - Added

- Added new Metrics API. (#3868, #3921)
- Enable PyTorch 1.7 compatibility (#3541)
- Added `LightningModule.to_torchscript` to support exporting as `ScriptModule` (#3258)
- Added warning when dropping unpicklable hparams (#2874)
- Added EMB similarity (#3349)
- Added `ModelCheckpoint.to_yaml` method (#3048)
- Allow `ModelCheckpoint` monitor to be `None`, meaning it will always save (#3630)
- Disabled optimizers setup during testing (#3059)
- Added support for datamodules to save and load checkpoints when training (#3563)
- Added support for datamodule in learning rate finder (#3425)
- Added gradient clip test for native AMP (#3754)
- Added dist lib to enable syncing anything across devices (#3762)
- Added broadcast to `TPUBackend` (#3814)
- Added `XLADeviceUtils` class to check XLA device type (#3274)

42.4.2 [0.10.0] - Changed

- Refactored accelerator backends:
 - moved TPU `xxx_step` to backend (#3118)
 - refactored DDP backend `forward` (#3119)
 - refactored GPU backend `__step` (#3120)
 - refactored Horovod backend (#3121, #3122)
 - remove obscure forward call in eval + CPU backend `__step` (#3123)
 - reduced all simplified forward (#3126)
 - added hook base method (#3127)
 - refactor eval loop to use hooks - use `test_mode` for if so we can split later (#3129)
 - moved `__step_end` hooks (#3130)
 - training forward refactor (#3134)

- training AMP scaling refactor (#3135)
- eval step scaling factor (#3136)
- add eval loop object to streamline eval loop (#3138)
- refactored dataloader process hook (#3139)
- refactored inner eval loop (#3141)
- final inner eval loop hooks (#3154)
- clean up hooks in `run_evaluation` (#3156)
- clean up data reset (#3161)
- expand eval loop out (#3165)
- moved hooks around in eval loop (#3195)
- remove `_evaluate_fx` (#3197)
- `Trainer.fit` hook clean up (#3198)
- DDPs train hooks (#3203)
- refactor DDP backend (#3204, #3207, #3208, #3209, #3210)
- reduced accelerator selection (#3211)
- group prepare data hook (#3212)
- added data connector (#3285)
- modular `is_overridden` (#3290)
- adding `Trainer.tune()` (#3293)
- move `run_pretrain_routine` -> `setup_training` (#3294)
- move train outside of setup training (#3297)
- move `prepare_data` to data connector (#3307)
- moved accelerator router (#3309)
- train loop refactor - moving train loop to own object (#3310, #3312, #3313, #3314)
- duplicate data interface definition up into `DataHooks` class (#3344)
- inner train loop (#3359, #3361, #3362, #3363, #3365, #3366, #3367, #3368, #3369, #3370, #3371, #3372, #3373, #3374, #3375, #3376, #3385, #3388, #3397)
- all logging related calls in a connector (#3395)
- device parser (#3400, #3405)
- added model connector (#3407)
- moved eval loop logging to loggers (#3408)
- moved eval loop (#3412#3408)
- trainer/separate argparse (#3421, #3428, #3432)
- move `lr_finder` (#3434)
- organize args (##3435, #3442, #3447, #3448, #3449, #3456)
- move specific accelerator code (#3457)

- group connectors (#3472)
- accelerator connector methods x/n (#3469, #3470, #3474)
- merge backends x/n (#3476, #3477, #3478, #3480, #3482)
- apex plugin (#3502)
- precision plugins (#3504)
- Result - make monitor default to checkpoint_on to simplify (#3571)
- reference to the Trainer on the LightningDataModule (#3684)
- add .log to lightning module (#3686, #3699, #3701, #3704, #3715)
- enable tracking original metric when step and epoch are both true (#3685)
- deprecated results obj, added support for simpler comms (#3681)
- move backends back to individual files (#3712)
- fixes logging for eval steps (#3763)
- decoupled DDP, DDP spawn (#3733, #3766, #3767, #3774, #3802, #3806)
- remove weight loading hack for ddp_cpu (#3808)
- separate torchelastic from DDP (#3810)
- separate SLURM from DDP (#3809)
- decoupled DDP2 (#3816)
- bug fix with logging val epoch end + monitor (#3812)
- decoupled DDP, DDP spawn (#3733, #3817, #3819, #3927)
- callback system and init DDP (#3836)
- adding compute environments (#3837, #3842)
- epoch can now log independently (#3843)
- test selecting the correct backend. temp backends while slurm and TorchElastic are decoupled (#3848)
- fixed init_slurm_connection causing hostname errors (#3856)
- moves init apex from LM to apex connector (#3923)
- moves sync bn to each backend (#3925)
- moves configure ddp to each backend (#3924)
- Deprecation warning (#3844)
- Changed LearningRateLogger to LearningRateMonitor (#3251)
- Used fsspec instead of gfile for all IO (#3320)
 - Swaped torch.load for fsspec load in DDP spawn backend (#3787)
 - Swaped torch.load for fsspec load in cloud_io loading (#3692)
 - Added support for to_disk() to use remote filepaths with fsspec (#3930)
 - Updated model_checkpoint's to_yaml to use fsspec open (#3801)
 - Fixed fsspec is inconsistant when doing fs.ls (#3805)
- Refactor GPUStatsMonitor to improve training speed (#3257)

- Changed IoU score behavior for classes absent in target and pred (#3098)
- Changed IoU `remove_bg` bool to `ignore_index` optional int (#3098)
- Changed defaults of `save_top_k` and `save_last` to `None` in `ModelCheckpoint` (#3680)
- `row_log_interval` and `log_save_interval` are now based on training loop's `global_step` instead of epoch-internal batch index (#3667)
- Silenced some warnings. verified ddp refactors (#3483)
- Cleaning up stale logger tests (#3490)
- Allow `ModelCheckpoint` `monitor` to be `None` (#3633)
- Enable `None` model checkpoint default (#3669)
- Skipped `best_model_path` if `checkpoint_callback` is `None` (#2962)
- Used `raise .. from ..` to explicitly chain exceptions (#3750)
- Mocking loggers (#3596, #3617, #3851, #3859, #3884, #3853, #3910, #3889, #3926)
- Write predictions in `LightningModule` instead of `EvalResult` #3882

42.4.3 [0.10.0] - Deprecated

- Deprecated `TrainResult` and `EvalResult`, use `self.log` and `self.write` from the `LightningModule` to log metrics and write predictions. `training_step` can now only return a scalar (for the loss) or a dictionary with anything you want. (#3681)
- Deprecate `early_stop_callback` `Trainer` argument (#3845)
- Rename `Trainer` arguments `row_log_interval` >> `log_every_n_steps` and `log_save_interval` >> `flush_logs_every_n_steps` (#3748)

42.4.4 [0.10.0] - Removed

- Removed experimental Metric API (#3868, #3943, #3949, #3946), listed changes before final removal:
 - Added `EmbeddingSimilarity` metric (#3349, #3358)
 - Added hooks to metric module interface (#2528)
 - Added error when AUROC metric is used for multiclass problems (#3350)
 - Fixed `ModelCheckpoint` with `save_top_k=-1` option not tracking the best models when a monitor metric is available (#3735)
 - Fixed counter-intuitive error being thrown in `Accuracy` metric for zero target tensor (#3764)
 - Fixed aggregation of metrics (#3517)
 - Fixed Metric aggregation (#3321)
 - Fixed RMSLE metric (#3188)
 - Renamed `reduction` to `class_reduction` in classification metrics (#3322)
 - Changed `class_reduction` similar to `sklearn` for classification metrics (#3322)
 - Renaming of precision recall metric (#3308)

42.4.5 [0.10.0] - Fixed

- Fixed `on_train_batch_start` hook to end epoch early (#3700)
- Fixed `num_sanity_val_steps` is clipped to `limit_val_batches` (#2917)
- Fixed ONNX model save on GPU (#3145)
- Fixed `GpuUsageLogger` to work on different platforms (#3008)
- Fixed auto-scale batch size not dumping `auto_lr_find` parameter (#3151)
- Fixed `batch_outputs` with optimizer frequencies (#3229)
- Fixed setting batch size in `LightningModule.datamodule` when using `auto_scale_batch_size` (#3266)
- Fixed Horovod distributed backend compatibility with native AMP (#3404)
- Fixed batch size auto scaling exceeding the size of the dataset (#3271)
- Fixed getting `experiment_id` from MLFlow only once instead of each training loop (#3394)
- Fixed `overfit_batches` which now correctly disables shuffling for the training loader. (#3501)
- Fixed gradient norm tracking for `row_log_interval > 1` (#3489)
- Fixed `ModelCheckpoint` name formatting (3164)
- Fixed auto-scale batch size (#3151)
- Fixed example implementation of `AutoEncoder` (#3190)
- Fixed invalid paths when remote logging with `TensorBoard` (#3236)
- Fixed change `t()` to `transpose()` as XLA devices do not support `.t()` on 1-dim tensor (#3252)
- Fixed (weights only) checkpoints loading without PL (#3287)
- Fixed `gather_all_tensors` cross GPUs in DDP (#3319)
- Fixed CometML save dir (#3419)
- Fixed forward key metrics (#3467)
- Fixed normalize mode at confusion matrix (replace NaNs with zeros) (#3465)
- Fixed global step increment in training loop when `training_epoch_end` hook is used (#3673)
- Fixed `dataloader shuffling` not getting turned off with `overfit_batches > 0` and `distributed_backend = "ddp"` (#3534)
- Fixed determinism in `DDPSpawnBackend` when using `seed_everything` in main process (#3335)
- Fixed `ModelCheckpoint` period to actually save every period epochs (#3630)
- Fixed `val_progress_bar` total with `num_sanity_val_steps` (#3751)
- Fixed Tuner dump: add `current_epoch` to `dumped_params` (#3261)
- Fixed `current_epoch` and `global_step` properties mismatch between `Trainer` and `LightningModule` (#3785)
- Fixed learning rate scheduler for optimizers with internal state (#3897)
- Fixed `tbptt_reduce_fx` when non-floating tensors are logged (#3796)
- Fixed model checkpoint frequency (#3852)
- Fixed logging non-tensor scalar with result breaks subsequent epoch aggregation (#3855)

- Fixed `TrainerEvaluationLoopMixin` activates `model.train()` at the end (#3858)
- Fixed `overfit_batches` when using with multiple `val/test_dataloaders` (#3857)
- Fixed enables `training_step` to return `None` (#3862)
- Fixed init nan for checkpointing (#3863)
- Fixed for `load_from_checkpoint` (#2776)
- Fixes incorrect `batch_sizes` when `Dataloader` returns a dict with multiple tensors (#3668)
- Fixed unexpected signature for `validation_step` (#3947)

42.5 [0.9.0] - 2020-08-20

42.5.1 [0.9.0] - Added

- Added `SyncBN` for DDP (#2801, #2838)
- Added basic `CSVLogger` (#2721)
- Added SSIM metrics (#2671)
- Added BLEU metrics (#2535)
- Added support to export a model to ONNX format (#2596)
- Added support for `Trainer(num_sanity_val_steps=-1)` to check all validation data before training (#2246)
- Added struct. output:
 - tests for val loop flow (#2605)
 - `EvalResult` support for train and val. loop (#2615, #2651)
 - weighted average in results obj (#2930)
 - fix result obj DP auto reduce (#3013)
- Added class `LightningDataModule` (#2668)
- Added support for PyTorch 1.6 (#2745)
- Added call `DataModule` hooks implicitly in trainer (#2755)
- Added support for Mean in DDP Sync (#2568)
- Added remaining sklearn metrics: `AveragePrecision`, `BalancedAccuracy`, `CohenKappaScore`, `DCG`, `Hamming`, `Hinge`, `Jaccard`, `MeanAbsoluteError`, `MeanSquaredError`, `MeanSquaredLogError`, `MedianAbsoluteError`, `R2Score`, `MeanPoissonDeviance`, `MeanGammaDeviance`, `MeanTweedieDeviance`, `ExplainedVariance` (#2562)
- Added support for `limit_{mode}_batches (int)` to work with infinite dataloader (`IterableDataset`) (#2840)
- Added support returning python scalars in DP (#1935)
- Added support to Tensorboard logger for `OmegaConf hparams` (#2846)
- Added tracking of basic states in `Trainer` (#2541)
- Tracks all outputs including TBPTT and multiple optimizers (#2890)

- Added GPU Usage Logger (#2932)
- Added `strict=False` for `load_from_checkpoint` (#2819)
- Added saving test predictions on multiple GPUs (#2926)
- Auto log the computational graph for loggers that support this (#3003)
- Added warning when changing monitor and using results obj (#3014)
- Added a hook `transfer_batch_to_device` to the `LightningDataModule` (#3038)

42.5.2 [0.9.0] - Changed

- Truncated long version numbers in progress bar (#2594)
- Enabling val/test loop disabling (#2692)
- Refactored into `accelerator` module:
 - GPU training (#2704)
 - TPU training (#2708)
 - DDP(2) backend (#2796)
 - Retrieve last logged val from result by key (#3049)
- Using `.comet.config` file for `CometLogger` (#1913)
- Updated hooks arguments - breaking for setup and teardown (#2850)
- Using `gfile` to support remote directories (#2164)
- Moved optimizer creation after device placement for DDP backends (#2904)
- Support `**DictConfig` for `hparam` serialization (#2519)
- Removed callback metrics from test results obj (#2994)
- Re-enabled naming metrics in ckpt name (#3060)
- Changed progress bar epoch counting to start from 0 (#3061)

42.5.3 [0.9.0] - Deprecated

- Deprecated Trainer attribute `ckpt_path`, which will now be set by `weights_save_path` (#2681)

42.5.4 [0.9.0] - Removed

- Removed deprecated: (#2760)
 - core decorator `data_loader`
 - Module hook `on_sanity_check_start` and loading `load_from_metrics`
 - package `pytorch_lightning.logging`
 - Trainer arguments: `show_progress_bar`, `num_tpu_cores`, `use_amp`, `print_nan_grads`
 - LR Finder argument `num_accumulation_steps`

42.5.5 [0.9.0] - Fixed

- Fixed `accumulate_grad_batches` for last batch (#2853)
- Fixed setup call while testing (#2624)
- Fixed local rank zero casting (#2640)
- Fixed single scalar return from training (#2587)
- Fixed Horovod backend to scale LR schedulers with the optimizer (#2626)
- Fixed `dtype` and `device` properties not getting updated in submodules (#2657)
- Fixed `fast_dev_run` to run for all dataloaders (#2581)
- Fixed `save_dir` in loggers getting ignored by default value of `weights_save_path` when user did not specify `weights_save_path` (#2681)
- Fixed `weights_save_path` getting ignored when `logger=False` is passed to `Trainer` (#2681)
- Fixed TPU multi-core and `Float16` (#2632)
- Fixed test metrics not being logged with `LoggerCollection` (#2723)
- Fixed data transfer to device when using `torchtext.data.Field` and `include_lengths` is `True` (#2689)
- Fixed shuffle argument for distributed sampler (#2789)
- Fixed logging interval (#2694)
- Fixed loss value in the progress bar is wrong when `accumulate_grad_batches > 1` (#2738)
- Fixed correct CWD for ddp sub-processes when using Hydra (#2719)
- Fixed selecting GPUs using `CUDA_VISIBLE_DEVICES` (#2739, #2796)
- Fixed false `num_classes` warning in metrics (#2781)
- Fixed shell injection vulnerability in subprocess call (#2786)
- Fixed LR finder and `hparams` compatibility (#2821)
- Fixed `ModelCheckpoint` not saving the latest information when `save_last=True` (#2881)
- Fixed ImageNet example: learning rate scheduler, number of workers and batch size when using DDP (#2889)
- Fixed apex gradient clipping (#2829)
- Fixed save apex scaler states (#2828)
- Fixed a model loading issue with inheritance and variable positional arguments (#2911)
- Fixed passing `non_blocking=True` when transferring a batch object that does not support it (#2910)
- Fixed checkpointing to remote file paths (#2925)
- Fixed adding val step argument to metrics (#2986)
- Fixed an issue that caused `Trainer.test()` to stall in ddp mode (#2997)
- Fixed gathering of results with tensors of varying shape (#3020)
- Fixed batch size auto-scaling feature to set the new value on the correct model attribute (#3043)
- Fixed automatic batch scaling not working with half precision (#3045)
- Fixed setting device to root gpu (#3042)

42.6 [0.8.5] - 2020-07-09

42.6.1 [0.8.5] - Added

- Added a PSNR metric: peak signal-to-noise ratio (#2483)
- Added functional regression metrics (#2492)

42.6.2 [0.8.5] - Removed

- Removed auto val reduce (#2462)

42.6.3 [0.8.5] - Fixed

- Flattening Wandb Hyperparameters (#2459)
- Fixed using the same DDP python interpreter and actually running (#2482)
- Fixed model summary input type conversion for models that have input dtype different from model parameters (#2510)
- Made `TensorBoardLogger` and `CometLogger` pickleable (#2518)
- Fixed a problem with `MLflowLogger` creating multiple run folders (#2502)
- Fixed `global_step` increment (#2455)
- Fixed TPU hanging example (#2488)
- Fixed `argparse` default value bug (#2526)
- Fixed Dice and IoU to avoid NaN by adding small eps (#2545)
- Fixed accumulate gradients schedule at epoch 0 (continued) (#2513)
- Fixed `Trainer .fit()` returning last not best weights in “ddp_spawn” (#2565)
- Fixed passing (do not pass) TPU weights back on test (#2566)
- Fixed DDP tests and `.test()` (#2512, #2570)

42.7 [0.8.4] - 2020-07-01

42.7.1 [0.8.4] - Added

- Added reduce ddp results on eval (#2434)
- Added a warning when an `IterableDataset` has `__len__` defined (#2437)

42.7.2 [0.8.4] - Changed

- Enabled no returns from eval ([#2446](#))

42.7.3 [0.8.4] - Fixed

- Fixes train outputs ([#2428](#))
- Fixes Conda dependencies ([#2412](#))
- Fixed Apex scaling with decoupled backward ([#2433](#))
- Fixed crashing or wrong displaying progressbar because of missing ipywidgets ([#2417](#))
- Fixed TPU saving dir ([fc26078e](#), [04e68f02](#))
- Fixed logging on rank 0 only ([#2425](#))

42.8 [0.8.3] - 2020-06-29

42.8.1 [0.8.3] - Fixed

- Fixed AMP wrong call ([593837e](#))
- Fixed batch typo ([92d1e75](#))

42.9 [0.8.2] - 2020-06-28

42.9.1 [0.8.2] - Added

- Added TorchText support for moving data to GPU ([#2379](#))

42.9.2 [0.8.2] - Changed

- Changed epoch indexing from 0 instead of 1 ([#2289](#))
- Refactor Model backward ([#2276](#))
- Refactored `training_batch` + tests to verify correctness ([#2327](#), [#2328](#))
- Refactored training loop ([#2336](#))
- Made optimization steps for hooks ([#2363](#))
- Changed default apex level to 'O2' ([#2362](#))

42.9.3 [0.8.2] - Removed

- Moved `TrainsLogger` to `Bolts` (#2384)

42.9.4 [0.8.2] - Fixed

- Fixed parsing TPU arguments and TPU tests (#2094)
- Fixed number batches in case of multiple dataloaders and `limit_{*}_batches` (#1920, #2226)
- Fixed an issue with forward hooks not being removed after model summary (#2298)
- Fix for `load_from_checkpoint()` not working with absolute path on Windows (#2294)
- Fixed an issue how `_has_len` handles `NotImplementedError` e.g. raised by `torchtext.data.Iterator` (#2293), (#2307)
- Fixed `average_precision` metric (#2319)
- Fixed ROC metric for CUDA tensors (#2304)
- Fixed `average_precision` metric (#2319)
- Fixed lost compatibility with custom datatypes implementing `.to` (#2335)
- Fixed loading model with kwargs (#2387)
- Fixed `sum(0)` for `trainer.num_val_batches` (#2268)
- Fixed checking if the parameters are a `DictConfig` Object (#2216)
- Fixed SLURM weights saving (#2341)
- Fixed swaps LR scheduler order (#2356)
- Fixed adding tensorboard hparams logging test (#2342)
- Fixed use model ref for tear down (#2360)
- Fixed logger crash on DDP (#2388)
- Fixed several issues with early stopping and checkpoint callbacks (#1504, #2391)
- Fixed loading past checkpoints from v0.7.x (#2405)
- Fixed loading model without arguments (#2403)
- Fixed Windows compatibility issue (#2358)

42.10 [0.8.1] - 2020-06-19

42.10.1 [0.8.1] - Fixed

- Fixed the `load_from_checkpoint` path detected as URL bug (#2244)
- Fixed hooks - added barrier (#2245, #2257, #2260)
- Fixed hparams - remove frame inspection on `self.hparams` (#2253)
- Fixed setup and on fit calls (#2252)
- Fixed GPU template (#2255)

42.11 [0.8.0] - 2020-06-18

42.11.1 [0.8.0] - Added

- Added `overfit_batches`, `limit_{val|test}_batches` flags (overfit now uses training set for all three) (#2213)
- Added metrics
 - Base classes (#1326, #1877)
 - Sklearn metrics classes (#1327)
 - Native torch metrics (#1488, #2062)
 - docs for all Metrics (#2184, #2209)
 - Regression metrics (#2221)
- Added type hints in `Trainer.fit()` and `Trainer.test()` to reflect that also a list of dataloaders can be passed in (#1723)
- Allow dataloaders without sampler field present (#1907)
- Added option `save_last` to save the model at the end of every epoch in `ModelCheckpoint` (#1908)
- Early stopping checks on `_validation_end` (#1458)
- Attribute `best_model_path` to `ModelCheckpoint` for storing and later retrieving the path to the best saved model file (#1799)
- Speed up single-core TPU training by loading data using `ParallelLoader` (#2033)
- Added a model hook `transfer_batch_to_device` that enables moving custom data structures to the target device (1756)
- Added `black` formatter for the code with code-checker on pull (1610)
- Added back the slow spawn ddp implementation as `ddp_spawn` (#2115)
- Added loading checkpoints from URLs (#1667)
- Added a callback method `on_keyboard_interrupt` for handling `KeyboardInterrupt` events during training (#2134)
- Added a decorator `auto_move_data` that moves data to the correct device when using the `LightningModule` for inference (#1905)
- Added `ckpt_path` option to `LightningModule.test(...)` to load particular checkpoint (#2190)
- Added `setup` and `teardown` hooks for model (#2229)

42.11.2 [0.8.0] - Changed

- Allow user to select individual TPU core to train on (#1729)
- Removed non-finite values from loss in `LRFinder` (#1862)
- Allow passing model hyperparameters as complete kwarg list (#1896)
- Renamed `ModelCheckpoint`'s attributes `best` to `best_model_score` and `kth_best_model` to `kth_best_model_path` (#1799)
- Re-Enable `Logger`'s `ImportErrors` (#1938)

- Changed the default value of the Trainer argument `weights_summary` from `full` to `top` (#2029)
- Raise an error when lightning replaces an existing sampler (#2020)
- Enabled `prepare_data` from correct processes - clarify local vs global rank (#2166)
- Remove explicit flush from tensorboard logger (#2126)
- Changed epoch indexing from 1 instead of 0 (#2206)

42.11.3 [0.8.0] - Deprecated

- Deprecated flags: (#2213)
 - `overfit_pct` in favour of `overfit_batches`
 - `val_percent_check` in favour of `limit_val_batches`
 - `test_percent_check` in favour of `limit_test_batches`
- Deprecated `ModelCheckpoint`'s attributes `best` and `kth_best_model` (#1799)
- Dropped official support/testing for older PyTorch versions <1.3 (#1917)
- Deprecated Trainer `proc_rank` in favour of `global_rank` (#2166, #2269)

42.11.4 [0.8.0] - Removed

- Removed unintended Trainer argument `progress_bar_callback`, the callback should be passed in by `Trainer(callbacks=[...])` instead (#1855)
- Removed obsolete `self._device` in Trainer (#1849)
- Removed deprecated API (#2073)
 - Packages: `pytorch_lightning.pt_overrides`, `pytorch_lightning.root_module`
 - Modules: `pytorch_lightning.logging.comet_logger`, `pytorch_lightning.logging.mlflow_logger`, `pytorch_lightning.logging.test_tube_logger`, `pytorch_lightning.overrides.override_data_parallel`, `pytorch_lightning.core.model_saving`, `pytorch_lightning.core.root_module`
 - Trainer arguments: `add_row_log_interval`, `default_save_path`, `gradient_clip`, `nb_gpu_nodes`, `max_nb_epochs`, `min_nb_epochs`, `nb_sanity_val_steps`
 - Trainer attributes: `nb_gpu_nodes`, `num_gpu_nodes`, `gradient_clip`, `max_nb_epochs`, `min_nb_epochs`, `nb_sanity_val_steps`, `default_save_path`, `tng_tqdm_dic`

42.11.5 [0.8.0] - Fixed

- Run graceful training teardown on interpreter exit (#1631)
- Fixed user warning when apex was used together with learning rate schedulers (#1873)
- Fixed multiple calls of `EarlyStopping` callback (#1863)
- Fixed an issue with `Trainer.from_argparse_args` when passing in unknown Trainer args (#1932)
- Fixed bug related to logger not being reset correctly for model after tuner algorithms (#1933)
- Fixed root node resolution for SLURM cluster with dash in host name (#1954)

- Fixed `LearningRateLogger` in multi-scheduler setting (#1944)
- Fixed test configuration check and testing (#1804)
- Fixed an issue with `Trainer` constructor silently ignoring unknown/misspelled arguments (#1820)
- Fixed `save_weights_only` in `ModelCheckpoint` (#1780)
- Allow use of same `WandbLogger` instance for multiple training loops (#2055)
- Fixed an issue with `_auto_collect_arguments` collecting local variables that are not constructor arguments and not working for signatures that have the instance not named `self` (#2048)
- Fixed mistake in parameters' grad norm tracking (#2012)
- Fixed CPU and hanging GPU crash (#2118)
- Fixed an issue with the model summary and `example_input_array` depending on a specific ordering of the submodules in a `LightningModule` (#1773)
- Fixed Tpu logging (#2230)
- Fixed Pid port + duplicate `rank_zero` logging (#2140, #2231)

42.12 [0.7.6] - 2020-05-16

42.12.1 [0.7.6] - Added

- Added callback for logging learning rates (#1498)
- Added transfer learning example (for a binary classification task in computer vision) (#1564)
- Added type hints in `Trainer.fit()` and `Trainer.test()` to reflect that also a list of dataloaders can be passed in (#1723).
- Added auto scaling of batch size (#1638)
- The progress bar metrics now also get updated in `training_epoch_end` (#1724)
- Enable `NeptuneLogger` to work with `distributed_backend=ddp` (#1753)
- Added option to provide seed to random generators to ensure reproducibility (#1572)
- Added override for hparams in `load_from_ckpt` (#1797)
- Added support multi-node distributed execution under `torchelastic` (#1811, #1818)
- Added using `store_true` for bool args (#1822, #1842)
- Added dummy logger for internally disabling logging for some features (#1836)

42.12.2 [0.7.6] - Changed

- Enable `non-blocking` for device transfers to GPU (#1843)
- Replace `meta_tags.csv` with `hparams.yaml` (#1271)
- Reduction when `batch_size < num_gpus` (#1609)
- Updated `LightningTemplateModel` to look more like Colab example (#1577)
- Don't convert `namedtuple` to `tuple` when transferring the batch to target device (#1589)
- Allow passing hparams as keyword argument to `LightningModule` when loading from checkpoint (#1639)

- Args should come after the last positional argument (#1807)
- Made ddp the default if no backend specified with multiple GPUs (#1789)

42.12.3 [0.7.6] - Deprecated

- Deprecated `tags_csv` in favor of `hparams_file` (#1271)

42.12.4 [0.7.6] - Fixed

- Fixed broken link in PR template (#1675)
- Fixed ModelCheckpoint not None checking filepath (#1654)
- Trainer now calls `on_load_checkpoint()` when resuming from a checkpoint (#1666)
- Fixed sampler logic for ddp with iterable dataset (#1734)
- Fixed `_reset_eval_dataloader()` for IterableDataset (#1560)
- Fixed Horovod distributed backend to set the `root_gpu` property (#1669)
- Fixed wandb logger `global_step` affects other loggers (#1492)
- Fixed disabling progress bar on non-zero ranks using Horovod backend (#1709)
- Fixed bugs that prevent lr finder to be used together with early stopping and validation dataloaders (#1676)
- Fixed a bug in Trainer that prepended the checkpoint path with `version_` when it shouldn't (#1748)
- Fixed lr key name in case of param groups in LearningRateLogger (#1719)
- Fixed saving native AMP scaler state (introduced in #1561)
- Fixed accumulation parameter and suggestion method for learning rate finder (#1801)
- Fixed num processes wasn't being set properly and auto sampler was ddp failing (#1819)
- Fixed bugs in semantic segmentation example (#1824)
- Fixed saving native AMP scaler state (#1561, #1777)
- Fixed native amp + ddp (#1788)
- Fixed hparam logging with metrics (#1647)

42.13 [0.7.5] - 2020-04-27

42.13.1 [0.7.5] - Changed

- Allow logging of metrics together with hparams (#1630)
- Allow metrics logged together with hparams (#1630)

42.13.2 [0.7.5] - Removed

- Removed Warning from trainer loop (#1634)

42.13.3 [0.7.5] - Fixed

- Fixed ModelCheckpoint not being fixable (#1632)
- Fixed CPU DDP breaking change and DDP change (#1635)
- Tested pickling (#1636)

42.14 [0.7.4] - 2020-04-26

42.14.1 [0.7.4] - Added

- Added flag `replace_sampler_ddp` to manually disable sampler replacement in DDP (#1513)
- Added speed parity tests (max 1 sec difference per epoch)(#1482)
- Added `auto_select_gpus` flag to trainer that enables automatic selection of available GPUs on exclusive mode systems.
- Added learning rate finder (#1347)
- Added support for ddp mode in clusters without SLURM (#1387)
- Added `test_dataloaders` parameter to `Trainer.test()` (#1434)
- Added `terminate_on_nan` flag to trainer that performs a NaN check with each training iteration when set to True (#1475)
- Added speed parity tests (max 1 sec difference per epoch)(#1482)
- Added `terminate_on_nan` flag to trainer that performs a NaN check with each training iteration when set to True. (#1475)
- Added `ddp_cpu` backend for testing ddp without GPUs (#1158)
- Added `Horovod` support as a distributed backend `Trainer(distributed_backend='horovod')` (#1529)
- Added support for 8 core distributed training on Kaggle TPU's (#1568)
- Added support for native AMP (#1561, #1580)

42.14.2 [0.7.4] - Changed

- Changed the default behaviour to no longer include a NaN check with each training iteration. (#1475)
- Decoupled the progress bar from trainer` it is a callback now and can be customized or even be replaced entirely (#1450).
- Changed lr schedule step interval behavior to update every backwards pass instead of every forwards pass (#1477)
- Defines shared proc. rank, remove rank from instances (e.g. loggers) (#1408)
- Updated semantic segmentation example with custom U-Net and logging (#1371)

- Disabled val and test shuffling (#1600)

42.14.3 [0.7.4] - Deprecated

- Deprecated `training_tqdm_dict` in favor of `progress_bar_dict` (#1450).

42.14.4 [0.7.4] - Removed

- Removed `test_dataloaders` parameter from `Trainer.fit()` (#1434)

42.14.5 [0.7.4] - Fixed

- Added the possibility to pass nested metrics dictionaries to loggers (#1582)
- Fixed memory leak from `opt` return (#1528)
- Fixed saving checkpoint before deleting old ones (#1453)
- Fixed loggers - flushing last logged metrics even before continue, e.g. `trainer.test()` results (#1459)
- Fixed optimizer configuration when `configure_optimizers` returns dict without `lr_scheduler` (#1443)
- Fixed `LightningModule` - mixing hparams and arguments in `LightningModule.__init__()` crashes `load_from_checkpoint()` (#1505)
- Added a missing call to the `on_before_zero_grad` model hook (#1493).
- Allow use of sweeps with `WandbLogger` (#1512)
- Fixed a bug that caused the `callbacks` `Trainer` argument to reference a global variable (#1534).
- Fixed a bug that set all boolean CLI arguments from `Trainer.add_argparse_args` always to `True` (#1571)
- Fixed do not copy the batch when training on a single GPU (#1576, #1579)
- Fixed soft checkpoint removing on DDP (#1408)
- Fixed automatic parser bug (#1585)
- Fixed bool conversion from string (#1606)

42.15 [0.7.3] - 2020-04-09

42.15.1 [0.7.3] - Added

- Added `rank_zero_warn` for warning only in rank 0 (#1428)

42.15.2 [0.7.3] - Fixed

- Fixed default `DistributedSampler` for DDP training (#1425)
- Fixed workers warning not on windows (#1430)
- Fixed returning tuple from `run_training_batch` (#1431)
- Fixed gradient clipping (#1438)
- Fixed pretty print (#1441)

42.16 [0.7.2] - 2020-04-07

42.16.1 [0.7.2] - Added

- Added same step loggers' metrics aggregation (#1278)
- Added parity test between a vanilla MNIST model and lightning model (#1284)
- Added parity test between a vanilla RNN model and lightning model (#1351)
- Added Reinforcement Learning - Deep Q-network (DQN) lightning example (#1232)
- Added support for hierarchical dict (#1152)
- Added `TrainsLogger` class (#1122)
- Added type hints to `pytorch_lightning.core` (#946)
- Added support for `IterableDataset` in validation and testing (#1104)
- Added support for non-primitive types in hparams for `TensorboardLogger` (#1130)
- Added a check that stops the training when loss or weights contain NaN or inf values. (#1097)
- Added support for `IterableDataset` when `val_check_interval=1.0` (default), this will trigger validation at the end of each epoch. (#1283)
- Added `summary` method to `Profilers`. (#1259)
- Added informative errors if user defined dataloader has zero length (#1280)
- Added testing for python 3.8 (#915)
- Added a `training_epoch_end` method which is the mirror of `validation_epoch_end`. (#1357)
- Added model configuration checking (#1199)
- Added support for optimizer frequencies through `LightningModule.configure_optimizers()` (#1269)
- Added option to run without an optimizer by returning `None` from `configure_optimizers`. (#1279)
- Added a warning when the number of data loader workers is small. (#1378)

42.16.2 [0.7.2] - Changed

- Changed (renamed and refactored) `TensorRunningMean` -> `TensorRunningAccum`: running accumulations were generalized. (#1278)
- Changed `progress_bar_refresh_rate` trainer flag to disable progress bar when set to 0. (#1108)
- Enhanced `load_from_checkpoint` to also forward params to the model (#1307)
- Updated references to `self.forward()` to instead use the `__call__` interface. (#1211)
- Changed default behaviour of `configure_optimizers` to use no optimizer rather than Adam. (#1279)
- Allow to upload models on W&B (#1339)
- On DP and DDP2 unsqueeze is automated now (#1319)
- Did not always create a `DataLoader` during reinstantiation, but the same type as before (if subclass of `DataLoader`) (#1346)
- Did not interfere with a default sampler (#1318)
- Remove default Adam optimizer (#1317)
- Give warnings for unimplemented required lightning methods (#1317)
- Made `evaluate` method private >> `Trainer._evaluate(...)`. (#1260)
- Simplify the PL examples structure (shallower and more readable) (#1247)
- Changed min max gpu memory to be on their own plots (#1358)
- Remove `.item` which causes sync issues (#1254)
- Changed smoothing in TQDM to decrease variability of time remaining between training / eval (#1194)
- Change default logger to dedicated one (#1064)

42.16.3 [0.7.2] - Deprecated

- Deprecated Trainer argument `print_nan_grads` (#1097)
- Deprecated Trainer argument `show_progress_bar` (#1108)

42.16.4 [0.7.2] - Removed

- Removed test for no test dataloader in `.fit` (#1495)
- Removed duplicated module `pytorch_lightning.utilities.arg_parse` for loading CLI arguments (#1167)
- Removed wandb logger's `finalize` method (#1193)
- Dropped `torchvision` dependency in tests and added own MNIST dataset class instead (#986)

42.16.5 [0.7.2] - Fixed

- Fixed `model_checkpoint` when saving all models (#1359)
- `Trainer.add_argparse_args` classmethod fixed. Now it adds a type for the arguments (#1147)
- Fixed bug related to type checking of `ReduceLROnPlateau` lr schedulers (#1126)
- Fixed a bug to ensure lightning checkpoints to be backward compatible (#1132)
- Fixed a bug that created an extra dataloader with active `reload_dataloaders_every_epoch` (#1196)
- Fixed all warnings and errors in the docs build process (#1191)
- Fixed an issue where `val_percent_check=0` would not disable validation (#1251)
- Fixed average of incomplete `TensorRunningMean` (#1309)
- Fixed `WandbLogger.watch` with `wandb.init()` (#1311)
- Fixed an issue with early stopping that would prevent it from monitoring training metrics when validation is disabled / not implemented (#1235).
- Fixed a bug that would cause `trainer.test()` to run on the validation set when overloading `validation_epoch_end` and `test_end` (#1353)
- Fixed `WandbLogger.watch` - use of the `watch` method without importing `wandb` (#1311)
- Fixed `WandbLogger` to be used with 'ddp' - allow reinit in sub-processes (#1149, #1360)
- Made `training_epoch_end` behave like `validation_epoch_end` (#1357)
- Fixed `fast_dev_run` running validation twice (#1365)
- Fixed pickle error from quick patch `__code__` (#1352)
- Fixed memory leak on GPU0 (#1094, #1349)
- Fixed checkpointing interval (#1272)
- Fixed validation and training loops run the partial dataset (#1192)
- Fixed running `on_validation_end` only on main process in DDP (#1125)
- Fixed `load_spawn_weights` only in proc rank 0 (#1385)
- Fixes `use_amp` issue (#1145)
- Fixes using deprecated `use_amp` attribute (#1145)
- Fixed Tensorboard logger error: `lightning_logs` directory not exists in multi-node DDP on nodes with rank `!= 0` (#1377)
- Fixed `Unimplemented backend XLA` error on TPU (#1387)

42.17 [0.7.1] - 2020-03-07

42.17.1 [0.7.1] - Fixed

- Fixes print issues and `data_loader` (#1080)

42.18 [0.7.0] - 2020-03-06

42.18.1 [0.7.0] - Added

- Added automatic sampler setup. Depending on DDP or TPU, lightning configures the sampler correctly (user needs to do nothing) (#926)
- Added `reload_dataloaders_every_epoch=False` flag for trainer. Some users require reloading data every epoch (#926)
- Added `progress_bar_refresh_rate=50` flag for trainer. Throttle refresh rate on notebooks (#926)
- Updated governance docs
- Added a check to ensure that the metric used for early stopping exists before training commences (#542)
- Added `optimizer_idx` argument to backward hook (#733)
- Added `entity` argument to `WandbLogger` to be passed to `wandb.init` (#783)
- Added a tool for profiling training runs (#782)
- Improved flexibility for naming of TensorBoard logs, can now set `version` to a `str` to just save to that directory, and use `name=' '` to prevent experiment-name directory (#804)
- Added option to specify `step` key when logging metrics (#808)
- Added `train_dataloader`, `val_dataloader` and `test_dataloader` arguments to `Trainer.fit()`, for alternative data parsing (#759)
- Added Tensor Processing Unit (TPU) support (#868)
- Added semantic segmentation example (#751, #876, #881)
- Split callbacks in multiple files (#849)
- Support for user defined callbacks (#889 and #950)
- Added support for multiple loggers to be passed to `Trainer` as an iterable (e.g. list, tuple, etc.) (#903)
- Added support for step-based learning rate scheduling (#941)
- Added support for logging `hparams` as `dict` (#1029)
- Checkpoint and early stopping now work without `val. step` (#1041)
- Support graceful training cleanup after Keyboard Interrupt (#856, #1019)
- Added type hints for function arguments (#912,)
- Added default `argparser` for `Trainer` (#952, #1023)
- Added TPU gradient clipping (#963)
- Added max/min number of steps in `Trainer` (#728)

42.18.2 [0.7.0] - Changed

- Improved NeptuneLogger by adding `close_after_fit` argument to allow logging after training(#908)
- Changed default TQDM to use `tqdm.auto` for prettier outputs in IPython notebooks (#752)
- Changed `pytorch_lightning.logging` to `pytorch_lightning.loggers` (#767)
- Moved the default `tqdm_dict` definition from Trainer to LightningModule, so it can be overridden by the user (#749)
- Moved functionality of `LightningModule.load_from_metrics` into `LightningModule.load_from_checkpoint` (#995)
- Changed Checkpoint path parameter from `filepath` to `dirpath` (#1016)
- Freezed models hparams as Namespace property (#1029)
- Dropped logging config in package init (#1015)
- Renames model steps (#1051)
 - `training_end` >> `training_epoch_end`
 - `validation_end` >> `validation_epoch_end`
 - `test_end` >> `test_epoch_end`
- Refactor dataloading, supports infinite dataloader (#955)
- Create single file in TensorBoardLogger (#777)

42.18.3 [0.7.0] - Deprecated

- Deprecated `pytorch_lightning.logging` (#767)
- Deprecated `LightningModule.load_from_metrics` in favour of `LightningModule.load_from_checkpoint` (#995, #1079)
- Deprecated `@data_loader` decorator (#926)
- Deprecated model steps `training_end`, `validation_end` and `test_end` (#1051, #1056)

42.18.4 [0.7.0] - Removed

- Removed dependency on pandas (#736)
- Removed dependency on torchvision (#797)
- Removed dependency on scikit-learn (#801)

42.18.5 [0.7.0] - Fixed

- Fixed a bug where early stopping `on_end_epoch` would be called inconsistently when `check_val_every_n_epoch == 0` (#743)
- Fixed a bug where the model checkpointer didn't write to the same directory as the logger (#771)
- Fixed a bug where the `TensorBoardLogger` class would create an additional empty log file during fitting (#777)
- Fixed a bug where `global_step` was advanced incorrectly when using `accumulate_grad_batches > 1` (#832)
- Fixed a bug when calling `self.logger.experiment` with multiple loggers (#1009)
- Fixed a bug when calling `logger.append_tags` on a `NeptuneLogger` with a single tag (#1009)
- Fixed sending back data from `.spawn` by saving and loading the trained model in/out of the process (#1017)
- Fixed port collision on DDP (#1010)
- Fixed/tested pass overrides (#918)
- Fixed comet logger to log after train (#892)
- Remove deprecated args to learning rate step function (#890)

42.19 [0.6.0] - 2020-01-21

42.19.1 [0.6.0] - Added

- Added support for resuming from a specific checkpoint via `resume_from_checkpoint` argument (#516)
- Added support for `ReduceLROnPlateau` scheduler (#320)
- Added support for Apex mode O2 in conjunction with Data Parallel (#493)
- Added option (`save_top_k`) to save the top k models in the `ModelCheckpoint` class (#128)
- Added `on_train_start` and `on_train_end` hooks to `ModelHooks` (#598)
- Added `TensorBoardLogger` (#607)
- Added support for weight summary of model with multiple inputs (#543)
- Added `map_location` argument to `load_from_metrics` and `load_from_checkpoint` (#625)
- Added option to disable validation by setting `val_percent_check=0` (#649)
- Added `NeptuneLogger` class (#648)
- Added `WandbLogger` class (#627)

42.19.2 [0.6.0] - Changed

- Changed the default progress bar to print to stdout instead of stderr (#531)
- Renamed `step_idx` to `step`, `epoch_idx` to `epoch`, `max_num_epochs` to `max_epochs` and `min_num_epochs` to `min_epochs` (#589)
- Renamed `total_batch_nb` to `total_batches`, `nb_val_batches` to `num_val_batches`, `nb_training_batches` to `num_training_batches`, `max_nb_epochs` to `max_epochs`, `min_nb_epochs` to `min_epochs`, `nb_test_batches` to `num_test_batches`, and `nb_val_batches` to `num_val_batches` (#567)
- Changed gradient logging to use parameter names instead of indexes (#660)
- Changed the default logger to `TensorBoardLogger` (#609)
- Changed the directory for tensorboard logging to be the same as model checkpointing (#706)

42.19.3 [0.6.0] - Deprecated

- Deprecated `max_nb_epochs` and `min_nb_epochs` (#567)
- Deprecated the `on_sanity_check_start` hook in `ModelHooks` (#598)

42.19.4 [0.6.0] - Removed

- Removed the `save_best_only` argument from `ModelCheckpoint`, use `save_top_k=1` instead (#128)

42.19.5 [0.6.0] - Fixed

- Fixed a bug which occurred when using Adagrad with cuda (#554)
- Fixed a bug where training would be on the GPU despite setting `gpus=0` or `gpus=[]` (#561)
- Fixed an error with `print_nan_gradients` when some parameters do not require gradient (#579)
- Fixed a bug where the progress bar would show an incorrect number of total steps during the validation sanity check when using multiple validation data loaders (#597)
- Fixed support for PyTorch 1.1.0 (#552)
- Fixed an issue with early stopping when using a `val_check_interval < 1.0` in `Trainer` (#492)
- Fixed bugs relating to the `CometLogger` object that would cause it to not work properly (#481)
- Fixed a bug that would occur when returning `-1` from `on_batch_start` following an early exit or when the batch was `None` (#509)
- Fixed a potential race condition with several processes trying to create checkpoint directories (#530)
- Fixed a bug where batch 'segments' would remain on the GPU when using `truncated_bptt > 1` (#532)
- Fixed a bug when using `IterableDataset` (#547)
- Fixed a bug where `.item` was called on non-tensor objects (#602)
- Fixed a bug where `Trainer.train` would crash on an uninitialized variable if the trainer was run after resuming from a checkpoint that was already at `max_epochs` (#608)
- Fixed a bug where early stopping would begin two epochs early (#617)

- Fixed a bug where `num_training_batches` and `num_test_batches` would sometimes be rounded down to zero (#649)
- Fixed a bug where an additional batch would be processed when manually setting `num_training_batches` (#653)
- Fixed a bug when batches did not have a `.copy` method (#701)
- Fixed a bug when using `log_gpu_memory=True` in Python 3.6 (#715)
- Fixed a bug where checkpoint writing could exit before completion, giving incomplete checkpoints (#689)
- Fixed a bug where `on_train_end` was not called when early stopping (#723)

42.20 [0.5.3] - 2019-11-06

42.20.1 [0.5.3] - Added

- Added option to disable default logger, checkpointer, and early stopping by passing `logger=False`, `checkpoint_callback=False` and `early_stop_callback=False` respectively
- Added `CometLogger` for use with Comet.ml
- Added `val_check_interval` argument to `Trainer` allowing validation to be performed at every given number of batches
- Added functionality to save and load hyperparameters using the standard checkpoint mechanism
- Added call to `torch.cuda.empty_cache` before training starts
- Added option for user to override the call to `backward`
- Added support for truncated backprop through time via the `truncated_bptt_steps` argument in `Trainer`
- Added option to operate on all outputs from `training_step` in DDP2
- Added a hook for modifying DDP init
- Added a hook for modifying Apex

42.20.2 [0.5.3] - Changed

- Changed experiment version to be padded with zeros (e.g. `/dir/version_9` becomes `/dir/version_0009`)
- Changed callback metrics to include any metrics given in logs or progress bar
- Changed the default for `save_best_only` in `ModelCheckpoint` to `True`
- Added `tng_data_loader` for backwards compatibility
- Renamed `MLFlowLogger.client` to `MLFlowLogger.experiment` for consistency
- Moved `global_step` increment to happen after the batch has been processed
- Changed weights restore to first attempt HPC weights before restoring normally, preventing both weights being restored and running out of memory
- Changed progress bar functionality to add multiple progress bars for train/val/test
- Changed calls to `print` to use `logging` instead

42.20.3 [0.5.3] - Deprecated

- Deprecated `tng_dataloader`

42.20.4 [0.5.3] - Fixed

- Fixed an issue where the number of batches was off by one during training
- Fixed a bug that occurred when setting a checkpoint callback and `early_stop_callback=False`
- Fixed an error when importing `CometLogger`
- Fixed a bug where the `gpus` argument had some unexpected behaviour
- Fixed a bug where the computed total number of batches was sometimes incorrect
- Fixed a bug where the progress bar would sometimes not show the total number of batches in test mode
- Fixed a bug when using the `log_gpu_memory='min_max'` option in `Trainer`
- Fixed a bug where checkpointing would sometimes erase the current directory

42.21 [0.5.2] - 2019-10-10

42.21.1 [0.5.2] - Added

- Added `weights_summary` argument to `Trainer` to be set to `full` (full summary), `top` (just top level modules) or `other`
- Added `tags` argument to `MLFlowLogger`

42.21.2 [0.5.2] - Changed

- Changed default for `amp_level` to `O1`

42.21.3 [0.5.2] - Removed

- Removed the `print_weights_summary` argument from `Trainer`

42.21.4 [0.5.2] - Fixed

- Fixed a bug where logs were not written properly
- Fixed a bug where `logger.finalize` wasn't called after training is complete
- Fixed callback metric errors in DDP
- Fixed a bug where `TestTubeLogger` didn't log to the correct directory

42.22 [0.5.1] - 2019-10-05

42.22.1 [0.5.1] - Added

- Added the `LightningLoggerBase` class for experiment loggers
- Added `MLFlowLogger` for logging with `mlflow`
- Added `TestTubeLogger` for logging with `test_tube`
- Added a different implementation of DDP (`distributed_backend='ddp2'`) where every node has one model using all GPUs
- Added support for optimisers which require a closure (e.g. LBFGS)
- Added automatic `MASTER_PORT` default for DDP when not set manually
- Added new GPU memory logging options `'min_max'` (log only the min/max utilization) and `'all'` (log all the GPU memory)

42.22.2 [0.5.1] - Changed

- Changed schedulers to always be called with the current epoch
- Changed `test_tube` to an optional dependency
- Changed data loaders to internally use a getter instead of a python property
- Disabled auto GPU loading when restoring weights to prevent out of memory errors
- Changed logging, early stopping and checkpointing to occur by default

42.22.3 [0.5.1] - Fixed

- Fixed a bug with samplers that do not specify `set_epoch`
- Fixed a bug when using the `MLFlowLogger` with unsupported data types, this will now raise a warning
- Fixed a bug where gradient norms were always zero using `track_grad_norm`
- Fixed a bug which causes a crash when logging memory

42.23 [0.5.0] - 2019-09-26

42.23.1 [0.5.0] - Changed

- Changed `data_batch` argument to `batch` throughout
- Changed `batch_i` argument to `batch_idx` throughout
- Changed `tng_dataloader` method to `train_dataloader`
- Changed `on_tng_metrics` method to `on_training_metrics`
- Changed `gradient_clip` argument to `gradient_clip_val`
- Changed `add_log_row_interval` to `row_log_interval`

42.23.2 [0.5.0] - Fixed

- Fixed a bug with tensorboard logging in multi-gpu setup

42.24 [0.4.9] - 2019-09-16

42.24.1 [0.4.9] - Added

- Added the flag `log_gpu_memory` to `Trainer` to deactivate logging of GPU memory utilization
- Added SLURM resubmit functionality (port from test-tube)
- Added optional `weight_save_path` to trainer to remove the need for a `checkpoint_callback` when using cluster training
- Added option to use single gpu per node with `DistributedDataParallel`

42.24.2 [0.4.9] - Changed

- Changed functionality of `validation_end` and `test_end` with multiple dataloaders to be given all of the dataloaders at once rather than in separate calls
- Changed `print_nan_grads` to only print the parameter value and gradients when they contain NaN
- Changed gpu API to take integers as well (e.g. `gpus=2` instead of `gpus=[0, 1]`)
- All models now loaded on to CPU to avoid device and out of memory issues in PyTorch

42.24.3 [0.4.9] - Fixed

- Fixed a bug where data types that implement `.to` but not `.cuda` would not be properly moved onto the GPU
- Fixed a bug where data would not be re-shuffled every epoch when using a `DistributedSampler`

42.25 [0.4.8] - 2019-08-31

42.25.1 [0.4.8] - Added

- Added `test_step` and `test_end` methods, used when `Trainer.test` is called
- Added `GradientAccumulationScheduler` callback which can be used to schedule changes to the number of accumulation batches
- Added option to skip the validation sanity check by setting `nb_sanity_val_steps = 0`

42.25.2 [0.4.8] - Fixed

- Fixed a bug when setting `nb_sanity_val_steps = 0`

42.26 [0.4.7] - 2019-08-24

42.26.1 [0.4.7] - Changed

- Changed the default `val_check_interval` to `1.0`
- Changed defaults for `nb_val_batches`, `nb_tng_batches` and `nb_test_batches` to `0`

42.26.2 [0.4.7] - Fixed

- Fixed a bug where the full validation set as used despite setting `val_percent_check`
- Fixed a bug where an `Exception` was thrown when using a data set containing a single batch
- Fixed a bug where an `Exception` was thrown if no `val_dataloader` was given
- Fixed a bug where tuples were not properly transfered to the GPU
- Fixed a bug where data of a non standard type was not properly handled by the trainer
- Fixed a bug when loading data as a tuple
- Fixed a bug where `AttributeError` could be suppressed by the `Trainer`

42.27 [0.4.6] - 2019-08-15

42.27.1 [0.4.6] - Added

- Added support for data to be given as a `dict` or `list` with a single `gpu`
- Added support for `configure_optimizers` to return a single optimizer, two list (optimizers and schedulers), or a single list

42.27.2 [0.4.6] - Fixed

- Fixed a bug where returning just an optimizer list (i.e. without schedulers) from `configure_optimizers` would throw an `Exception`

42.28 [0.4.5] - 2019-08-13

42.28.1 [0.4.5] - Added

- Added `optimizer_step` method that can be overridden to change the standard optimizer behaviour

42.29 [0.4.4] - 2019-08-12

42.29.1 [0.4.4] - Added

- Added support for multiple validation dataloaders
- Added support for latest test-tube logger (optimised for `torch==1.2.0`)

42.29.2 [0.4.4] - Changed

- `validation_step` and `val_dataloader` are now optional
- `lr_scheduler` is now activated after epoch

42.29.3 [0.4.4] - Fixed

- Fixed a bug where a warning would show when using `lr_scheduler` in `torch>1.1.0`
- Fixed a bug where an `Exception` would be thrown if using `torch.DistributedDataParallel` without using a `DistributedSampler`, this now throws a `Warning` instead

42.30 [0.4.3] - 2019-08-10

42.30.1 [0.4.3] - Fixed

- Fixed a bug where accumulate gradients would scale the loss incorrectly

42.31 [0.4.2] - 2019-08-08

42.31.1 [0.4.2] - Changed

- Changed install requirement to `torch==1.2.0`

42.32 [0.4.1] - 2019-08-08

42.32.1 [0.4.1] - Changed

- Changed install requirement to `torch==1.1.0`

42.33 [0.4.0] - 2019-08-08

42.33.1 [0.4.0] - Added

- Added 16-bit support for a single GPU
- Added support for training continuation (preserves epoch, global step etc.)

42.33.2 [0.4.0] - Changed

- Changed `training_step` and `validation_step`, outputs will no longer be automatically reduced

42.33.3 [0.4.0] - Removed

- Removed need for `Experiment` object in `Trainer`

42.33.4 [0.4.0] - Fixed

- Fixed issues with reducing outputs from generative models (such as images and text)

42.34 [0.3.6] - 2019-07-25

42.34.1 [0.3.6] - Added

- Added a decorator to do lazy data loading internally

42.34.2 [0.3.6] - Fixed

- Fixed a bug where `Experiment` object was not process safe, potentially causing logs to be overwritten

42.35 [0.3.5] - 2019-07-25

42.36 [0.3.4] - 2019-07-22

42.37 [0.3.3] - 2019-07-22

42.38 [0.3.2] - 2019-07-21

42.39 [0.3.1] - 2019-07-21

42.40 [0.2.x] - 2019-07-09

42.41 [0.1.x] - 2019-06-DD

INDICES AND TABLES

- `genindex`
- `search`

C

Callback (class in `pytorch_lightning.callbacks`), 106

D

`disable()` (`pytorch_lightning.callbacks.ProgressBar` method), 114

`disable()` (`pytorch_lightning.callbacks.ProgressBarBase` method), 116

E

`EarlyStopping` (class in `pytorch_lightning.callbacks`), 108

`enable()` (`pytorch_lightning.callbacks.ProgressBar` method), 114

`enable()` (`pytorch_lightning.callbacks.ProgressBarBase` method), 116

F

`format_checkpoint_name()` (`pytorch_lightning.callbacks.ModelCheckpoint` method), 113

G

`GPUStatsMonitor` (class in `pytorch_lightning.callbacks`), 109

`GradientAccumulationScheduler` (class in `pytorch_lightning.callbacks`), 110

I

`init_sanity_tqdm()` (`pytorch_lightning.callbacks.ProgressBar` method), 114

`init_test_tqdm()` (`pytorch_lightning.callbacks.ProgressBar` method), 114

`init_train_tqdm()` (`pytorch_lightning.callbacks.ProgressBar` method), 115

`init_validation_tqdm()` (`pytorch_lightning.callbacks.ProgressBar` method), 115

L

`LearningRateMonitor` (class in `pytorch_lightning.callbacks`), 111

M

`ModelCheckpoint` (class in `pytorch_lightning.callbacks`), 111

O

`on_batch_end()` (`pytorch_lightning.callbacks.Callback` method), 106

`on_batch_start()` (`pytorch_lightning.callbacks.Callback` method), 106

`on_batch_start()` (`pytorch_lightning.callbacks.LearningRateMonitor` method), 111

`on_epoch_end()` (`pytorch_lightning.callbacks.Callback` method), 106

`on_epoch_start()` (`pytorch_lightning.callbacks.Callback` method), 106

`on_epoch_start()` (`pytorch_lightning.callbacks.GradientAccumulationScheduler` method), 110

`on_epoch_start()` (`pytorch_lightning.callbacks.LearningRateMonitor` method), 111

`on_epoch_start()` (`pytorch_lightning.callbacks.ProgressBar` method), 115

`on_epoch_start()` (`pytorch_lightning.callbacks.ProgressBarBase` method), 116

`on_fit_end()` (`pytorch_lightning.callbacks.Callback` method), 106

`on_fit_start()` (`pytorch_lightning.callbacks.Callback` method), 107

<code>on_init_end()</code> <i>(pytorch_lightning.callbacks.Callback method), 107</i>	<code>on_test_batch_end()</code> <i>(pytorch_lightning.callbacks.ProgressBar method), 115</i>
<code>on_init_end()</code> <i>(pytorch_lightning.callbacks.ProgressBarBase method), 116</i>	<code>on_test_batch_end()</code> <i>(pytorch_lightning.callbacks.ProgressBarBase method), 116</i>
<code>on_init_start()</code> <i>(pytorch_lightning.callbacks.Callback method), 107</i>	<code>on_test_batch_start()</code> <i>(pytorch_lightning.callbacks.Callback method), 107</i>
<code>on_keyboard_interrupt()</code> <i>(pytorch_lightning.callbacks.Callback method), 107</i>	<code>on_test_end()</code> <i>(pytorch_lightning.callbacks.Callback method), 107</i>
<code>on_load_checkpoint()</code> <i>(pytorch_lightning.callbacks.Callback method), 107</i>	<code>on_test_end()</code> <i>(pytorch_lightning.callbacks.ProgressBar method), 115</i>
<code>on_load_checkpoint()</code> <i>(pytorch_lightning.callbacks.EarlyStopping method), 109</i>	<code>on_test_epoch_end()</code> <i>(pytorch_lightning.callbacks.Callback method), 107</i>
<code>on_load_checkpoint()</code> <i>(pytorch_lightning.callbacks.ModelCheckpoint method), 113</i>	<code>on_test_epoch_start()</code> <i>(pytorch_lightning.callbacks.Callback method), 107</i>
<code>on_pretrain_routine_end()</code> <i>(pytorch_lightning.callbacks.Callback method), 107</i>	<code>on_test_start()</code> <i>(pytorch_lightning.callbacks.Callback method), 107</i>
<code>on_pretrain_routine_start()</code> <i>(pytorch_lightning.callbacks.Callback method), 107</i>	<code>on_test_start()</code> <i>(pytorch_lightning.callbacks.ProgressBar method), 115</i>
<code>on_pretrain_routine_start()</code> <i>(pytorch_lightning.callbacks.ModelCheckpoint method), 113</i>	<code>on_test_start()</code> <i>(pytorch_lightning.callbacks.ProgressBarBase method), 116</i>
<code>on_sanity_check_end()</code> <i>(pytorch_lightning.callbacks.Callback method), 107</i>	<code>on_train_batch_end()</code> <i>(pytorch_lightning.callbacks.Callback method), 107</i>
<code>on_sanity_check_end()</code> <i>(pytorch_lightning.callbacks.ProgressBar method), 115</i>	<code>on_train_batch_end()</code> <i>(pytorch_lightning.callbacks.GPUStatsMonitor method), 110</i>
<code>on_sanity_check_start()</code> <i>(pytorch_lightning.callbacks.Callback method), 107</i>	<code>on_train_batch_end()</code> <i>(pytorch_lightning.callbacks.ProgressBar method), 115</i>
<code>on_sanity_check_start()</code> <i>(pytorch_lightning.callbacks.ProgressBar method), 115</i>	<code>on_train_batch_end()</code> <i>(pytorch_lightning.callbacks.ProgressBarBase method), 116</i>
<code>on_save_checkpoint()</code> <i>(pytorch_lightning.callbacks.Callback method), 107</i>	<code>on_train_batch_start()</code> <i>(pytorch_lightning.callbacks.Callback method), 107</i>
<code>on_save_checkpoint()</code> <i>(pytorch_lightning.callbacks.EarlyStopping method), 109</i>	<code>on_train_batch_start()</code> <i>(pytorch_lightning.callbacks.GPUStatsMonitor method), 110</i>
<code>on_save_checkpoint()</code> <i>(pytorch_lightning.callbacks.ModelCheckpoint method), 113</i>	<code>on_train_end()</code> <i>(pytorch_lightning.callbacks.Callback method), 107</i>
<code>on_test_batch_end()</code> <i>(pytorch_lightning.callbacks.Callback method), 107</i>	<code>on_train_end()</code> <i>(pytorch_lightning.callbacks.ProgressBar method), 115</i>

`on_train_epoch_end()` (pytorch_lightning.callbacks.Callback method), 107
`on_train_epoch_end()` (pytorch_lightning.callbacks.EarlyStopping method), 109
`on_train_epoch_start()` (pytorch_lightning.callbacks.Callback method), 107
`on_train_epoch_start()` (pytorch_lightning.callbacks.GPUStatsMonitor method), 110
`on_train_start()` (pytorch_lightning.callbacks.Callback method), 108
`on_train_start()` (pytorch_lightning.callbacks.GPUStatsMonitor method), 110
`on_train_start()` (pytorch_lightning.callbacks.LearningRateMonitor method), 111
`on_train_start()` (pytorch_lightning.callbacks.ProgressBar method), 115
`on_train_start()` (pytorch_lightning.callbacks.ProgressBarBase method), 116
`on_validation_batch_end()` (pytorch_lightning.callbacks.Callback method), 108
`on_validation_batch_end()` (pytorch_lightning.callbacks.ProgressBar method), 115
`on_validation_batch_end()` (pytorch_lightning.callbacks.ProgressBarBase method), 116
`on_validation_batch_start()` (pytorch_lightning.callbacks.Callback method), 108
`on_validation_end()` (pytorch_lightning.callbacks.Callback method), 108
`on_validation_end()` (pytorch_lightning.callbacks.EarlyStopping method), 109
`on_validation_end()` (pytorch_lightning.callbacks.ModelCheckpoint method), 113
`on_validation_end()` (pytorch_lightning.callbacks.ProgressBar method), 115
`on_validation_epoch_end()` (pytorch_lightning.callbacks.Callback method), 108
`on_validation_epoch_end()` (pytorch_lightning.callbacks.EarlyStopping method), 109
`on_validation_epoch_start()` (pytorch_lightning.callbacks.Callback method), 108
`on_validation_start()` (pytorch_lightning.callbacks.Callback method), 108
`on_validation_start()` (pytorch_lightning.callbacks.ProgressBar method), 115
`on_validation_start()` (pytorch_lightning.callbacks.ProgressBarBase method), 116

P

`ProgressBar` (class in pytorch_lightning.callbacks), 114
`ProgressBarBase` (class in pytorch_lightning.callbacks), 115

S

`save_checkpoint()` (pytorch_lightning.callbacks.ModelCheckpoint method), 113
`setup()` (pytorch_lightning.callbacks.Callback method), 108

T

`teardown()` (pytorch_lightning.callbacks.Callback method), 108
`test_batch_idx()` (pytorch_lightning.callbacks.ProgressBarBase property), 116
`to_yaml()` (pytorch_lightning.callbacks.ModelCheckpoint method), 113
`total_test_batches()` (pytorch_lightning.callbacks.ProgressBarBase property), 116
`total_train_batches()` (pytorch_lightning.callbacks.ProgressBarBase property), 117
`total_val_batches()` (pytorch_lightning.callbacks.ProgressBarBase property), 117
`train_batch_idx()` (pytorch_lightning.callbacks.ProgressBarBase property), 117

V

`val_batch_idx()` (pytorch_lightning.callbacks.ProgressBarBase property), 117