torchgpipe

Release 0.0.3

Sep 30, 2019

Contents

1	What is GPipe?	3
2	Documentations2.1Understanding GPipe2.2User Guide2.3API2.4Benchmarks2.5Changelog	5 6 11 13 14
3	Authors and Licensing	17
In	dex	19

A GPipe implementation in PyTorch.

```
from torchgpipe import GPipe
model = nn.Sequential(a, b, c, d)
model = GPipe(model, balance=[1, 1, 1, 1], chunks=8)
for input in data_loader:
    output = model(input)
```

CHAPTER 1

What is GPipe?

GPipe is a scalable pipeline parallelism library published by Google Brain, which allows efficient training of large, memory-consuming models. According to the paper, GPipe can train a 25x larger model by using 8x devices (TPU), and train a model 3.5x faster by using 4x devices.

GPipe: Efficient Training of Giant Neural Networks using Pipeline Parallelism

Google trained AmoebaNet-B with 557M parameters over GPipe. This model has achieved 84.3% top-1 and 97.0% top-5 accuracy on ImageNet classification benchmark (the state-of-the-art performance as of May 2019).

CHAPTER 2

Documentations

2.1 Understanding GPipe

GPipe uses (a) *Pipeline Parallelism* and (b) automatic recomputation of the forward propagation during the backpropagation, hence leverages training a large model. We refer to (b) as *Checkpointing*, following the well-known terminology in PyTorch community.

2.1.1 Pipeline Parallelism

GPipe splits a model into multiple partitions and places each partition on a different device to occupy more memory capacity. For example, we may split a model occupying 40GB of CUDA memory into 4 partitions each occupying 10GB, respectively.

This approach is called *model parallelism*. However, typical deep learning models are composed of sequential layers. In other words, usually the latter layer wouldn't work until the prior layer has finished. If a model is composed of fully sequential layers, even if we spread the model over two or more devices, only one device can be utilized at once.

GPipe splits a mini-batch into multiple micro-batches to make the devices work as parallel as possible. It is called *pipeline parallelism*. Basically, pipeline parallelism is a stack of small data parallelism. When each partition has finished processing a micro-batch, it can toss the output to the next partition and immediately can start to work on the next micro-batch. Now the partitions can be overlapped.

See also:

Model Parallel Best Practices in PyTorch Tutorials

There is still idle time called *bubble* because every partition has to wait for the first micro-batch from the prior partition. The bubble can be reduced by choosing a smaller size of micro-batches. But usually, larger batch size can utilize GPU more efficiently. Hence, GPU may be underutilized if too small size of micro-batches is chosen.

A faster partition should wait for adjacent slower partition. Therefore, imbalance over partitions also may cause GPU underutilization. Note that the overall performance is determined by the slowest partition.

2.1.2 Checkpointing

Checkpointing is applied to each partition to minimize the overall memory consumption by a model. During forward propagation, only the tensors at the boundaries between partitions are remembered. All other intermediate tensors are volatilized, and recomputed during backpropagation when necessary. Specifically, hidden layers consume the memory which is required by only a single micro-batch with checkpointing.

Checkpointing is a trade-off between performance and memory, because recomputation spends time just as much as the forward propagation. When you use *torchgpipe.GPipe*, you can decide to turn off checkpointing by checkpoint='never' option.

2.1.3 Deferred Batch Normalization

One of the goals of GPipe is *transparency*. GPipe shouldn't affect existing hyperparameters and output during training. However, if a module processes per mini-batch, not per single sample, it might be affected by GPipe since each module could see only a micro-batch at once.

Meanwhile, batch normalization is a module commonly used in CNN. The forward propagation of this module performs two procedures in training. Both the procedures are per mini-batch, not micro-batch:

- 1. Normalizing a mini-batch by the average and variance of the just given mini-batch.
- 2. Tracking moving statistics (mean and variance) of mini-batches to normalize batches in evaluation.

GPipe couldn't provide transparency for the first procedure (normalization). Per mini-batch normalization introduces a dependency among the micro-batches, hence it breaks the parallelism of GPipe. But the second procedure (tracking moving statistics) could be transparent with GPipe by accumulating statistics of all micro-batches.

torchgpipe provides this functionality as *deferred batch normalization*. But in the current implementation, it is slower than the vanilla batch normalization. That is why we turn off by default. If you need transparent moving statistics, turn on by deferred_batch_norm=True option in *GPipe*:

2.2 User Guide

2.2.1 Installation

torchgpipe is available on PyPI. Install by pip:

```
$ pip install torchgpipe
```

Python 3.6+ (CPython) is required.

PyTorch 1.1+ will be installed automatically if you don't have a satisfied one. However, we highly recommend you to use the latest version of PyTorch.

2.2.2 Applying GPipe

To train a module with GPipe, simply wrap it with *torchgpipe.GPipe*. Your module must be nn.Sequential as GPipe will automatically split the module into partitions with consecutive layers. *balance* argument determines the number of layers in each partition. *chunks* argument specifies the number of micro-batches. Input, output, and intermediate tensors must be Tensor or Tuple[Tensor, ...]. See also *Restrictions* for more details.

The below example code shows how to split a module with four layers into four partitions each having a single layer. This code also splits a mini-batch into 8 micro-batches:

```
from torchgpipe import GPipe
model = nn.Sequential(a, b, c, d)
model = GPipe(model, balance=[1, 1, 1, 1], chunks=8)
for input in data_loader:
    output = model(input)
```

GPipe optimizes training using CUDA. You should not move the module to a GPU yourself, because GPipe automatically moves each partition over different devices. By default, available GPUs starting from cuda: 0 are selected in order for each partition. You can also specify GPUs to use with *devices* parameter:

The typical model parallelism is a special case of GPipe. GPipe without micro-batches and checkpointing is equivalent to model parallelism. You can disable them with chunks=1 and checkpoint='never' options:

model = GPipe(model, balance=[2, 2], chunks=1, checkpoint='never')

2.2.3 Input and Output Device

Unlike a typical module, with GPipe, the input device is different from the output device except for when there is only one partition. This is because the first partition and last partition are placed in different devices.

Therefore, you have to move the input and target to the corresponding devices. It can be done with *GPipe.devices*, which holds the list of devices for each partition:

```
in_device = model.devices[0]
out_device = model.devices[-1]

for input, target in data_loader:
    # input on in_device
    input = input.to(in_device, non_blocking=True)
    # target on out_device
    target = target.to(out_device, non_blocking=True)
    # output on out_device
    output = model(input)
    loss = F.cross_entropy(output, target)
    loss.backward()
    ...
```

2.2.4 Automatic Balancing

It could be hard to determine the optimal balance of a model. In particular, if you are still designing a model, the model architecture may change over time. In this case, we highly recommend torchgpipe_balancing for automatic balancing. This won't give you the optimal balance, but a good-enough balance. Note that this is provided by *torchgpipe* package, and is not from the GPipe paper.

There are two balancing tools, *balance_by_time()* and *balance_by_size()*. Both are based on per-layer profiling. Just like PyTorch JIT, you need to feed a sample input into the model. *balance_by_time()* traces elapsed time of each layer, while *balance_by_size()* detects the CUDA memory usage of each layer. Choose the balancing tool for your needs:

```
from torchgpipe import GPipe
from torchgpipe_balancing import balance_by_time
sample = torch.rand(128, 3, 224, 224)
balance = balance_by_time(model, sample, partitions=4)
model = GPipe(model, balance, chunks=8)
```

2.2.5 Trade-offs

Number of Micro-batches

Number of micro-batches has a trade-off between GPU utilization per micro-batch and total area of bubble. You need to find the best number of micro-batches for your model.

GPU may slow down when processing many small micro-batches compared to larger micro-batches. GPU will not be fully utilized if each CUDA kernel is too cheap to compute, hence too small micro-batches cause underutilization. On the other hand, the area of bubble is minimized when the size of each micro-batch is minimal. Ideally, you should choose the largest number of micro-batches that doesn't underutilize GPUs.

As a side note, BatchNorm tends to perform worse with smaller batch size. Large number of micro-batches may affect the final performance of model using BatchNorm negatively just like in nn.DataParallel.

Checkpointing

Checkpointing drastically helps to reduce memory usage, but the overall training would slow down by about 25%. You can handle how to apply checkpointing on your model. There are three options:

- always Apply checkpointing over all micro-batches.
- except_last (default) Apply checkpointing except the last micro-batch.
- never Checkpointing is never applied.

Usually, checkpointing at the last micro-batch may not be useful because the saved memory will be reconstructed immediately. That's why we choose <code>except_last</code> as the default option.

If you decide not to use checkpointing at all, nn.DataParallel might be more efficient than GPipe.

2.2.6 Referential Transparency

Checkpointing executes forward propagation again at backpropagation, which is called *recomputation*. We assume that both the executions are identical. Hence, all layers should be referentially transparent in forward propagation. Here are the typical cases that break referential transparency:

- **In-place Operations:** We do not recommend using in-place operations with checkpointing. Especially, if an in-place operation such as add_(1) is applied to the input of a checkpointed partition, then the recomputation can't recover the original input.
- **Nondeterminism:** For example, nn.Dropout will produce different mask in recomputation from the forward propagation due to the randomness. This type of nondeterministic behaviors are not taken care of in torchgpipe yet.
- Side Effects: Some modules such as BatchNorm update their state in forward propagation. Hence, updated state in recomputation might not be identical to the original state.

2.2.7 Restrictions

If you get any errors, check the following restrictions first.

Sequential: Your module must be nn.Sequential. For example, the models in torchvision are not sequential. They can't be wrapped by *GPipe* directly:

```
>>> from torchvision.models.resnet import resnet101
>>> model = resnet101()
>>> type(model)
torchvision.models.resnet.ResNet
>>> GPipe(model, balance=..., chunks=...)
Traceback (most recent call last)
...
TypeError: module must be nn.Sequential to be partitioned
```

See the sequential ResNet example to figure out how to make a model into a nn.Sequential model.

nn.Sequential assumes that every underlying layer takes only one argument. Calling forward(x) on nn.Sequential(A(), B(), C()) is essentially the same as calling C(B(A(x))). Hence, you can't design an underlying layer with multiple arguments:

```
class MyModule(nn.Module):
    def forward(self, a, b, c):
        return a + b - c
model = nn.Sequential(..., MyModule(), ...)
model(input)  # FAILS!
```

Tensor or Tensors: As we discussed above, each layer must take only one argument due to nn.Sequential. There is one more restriction. Every underlying layers' input and output must be Tensor or Tuple [Tensor,

...]:

```
# OK
def forward(input: Tensor) -> Tensor: ...
def forward(input: Tensor) -> Tuple[Tensor, Tensor]: ...
def forward(input: Tuple[Tensor, Tensor]) -> Tensor: ...
# Error
def forward(input1: Tensor, input2: Tensor) -> Tensor: ...
def forward(input: Tensor, label: str) -> Tensor: ...
def forward(input: Tensor) -> Dict[str, Tensor]: ...
def forward(input: Tensor) -> Tuple[Tensor, str]: ...
```

The reason is that GPipe can't assume how the non-tensor inputs for a mini-batch can be split for micro-batches.

Unique Parameters: *GPipe* places each partition on the corresponding device. When placing a partition, the parameters of the partition are also moved to the destination. GPipe cannot support a module with a parameter on

two or more devices:

```
>>> conv1 = nn.Conv2d(3, 3, 1)
>>> conv2 = nn.Conv2d(3, 3, 1)
>>> conv1.weight = conv2.weight
>>> model = nn.Sequential(conv1, conv2)
>>> model = GPipe(model, balance=[1, 1], ...)
Traceback (most recent call last)
...
ValueError: module with duplicate parameters in distinct children is not supported
```

2.2.8 Complex Modules

This part of the documentation discusses how to implement a complex module compatible with *GPipe*. First, you should understand how GPipe works. See *Understanding GPipe*.

Skip Connections

Many deep learning models, such as ResNet or AmoebaNet, contain skip connections. There are two ways to implement skip connections. Let's assume we have to implement a skip connection like this:

```
latent = layer1(input)
latent = layer2(latent)
output = layer3(latent) + input # skip connection
```

To make this module sequential, we define modules for each layer. Simply, a skip connection can be implemented by making underlying layers with Tuple[Tensor, Tensor] parameter and return type:

```
class Layer1(nn.Module):
   #
   # input -- | -+-> layer1 ---- | --> output
              '----> skip
   #
    #
   def forward(self, input: Tensor) -> Tuple[Tensor, Tensor]:
       return layer1(input), input
class Layer2(nn.Module):
   # input -- | ---> layer2 ---- | --> output
   # skip -- ---- --> skip
   def forward(self, input_and_skip: Tuple[Tensor, Tensor]) -> Tuple[Tensor, Tensor]:
       input, skip = input_and_skip
       return layer2(input), skip
class Layer3(nn.Module):
   #
   # input -- | ---> layer3 --+- | --> output
   # skip -- |-----'
   def forward(self, input_and_skip: Tuple[Tensor, Tensor]) -> Tensor:
       input, skip = input_and_skip
       return layer3(input) + skip
model = nn.Sequential(Layer1(), Layer2(), Layer3())
```

Because of the skip connection being represented as a normal parameter, GPipe can move the tensors from partition to partition:

model = GPipe(model, balance=[1, 1, 1], chunks=8)

It is the most straightforward approach to implement skip connections. But there is a disadvantage. In the above example, the skip tensor is copied to the second device, but it is never used on the second device. Unnecessarily copying tensor wastes time and memory.

Detecting Recomputation

Checkpointing in GPipe performs forward propagations twice. The second forward propagation is called *recomputation*. This may cause a problem when a module such as nn.BatchNorm2d updates its running estimates of batch statistics on each forward propagation. It should not update the running estimates again during the recomputation. To avoid updating the running estimates twice, modules' forward method needs be able to detect that this is the recomputation.

It can be done by *is_recomputing()*. This function returns True if called during the recomputation:

```
class Counter(nn.Module):
    def __init__(self):
        super().__init__()
        self.counter = 0
    def forward(self, input):
        if not is_recomputing():
            self.counter += 1
        return input
```

Note: deferred_batch_norm=True on GPipe will prevent updating the running statistics twice.

2.3 API

2.3.1 GPipe Module

```
class torchgpipe.GPipe(module, balance, **kwargs)
```

Wraps an arbitrary nn.Sequential module to train on GPipe. If the module requires lots of memory, GPipe will be very efficient:

```
model = nn.Sequential(a, b, c, d)
model = GPipe(model, balance=[1, 1, 1, 1], chunks=8)
output = model(input)
```

GPipe combines pipeline parallelism with checkpointing to reduce peak memory required to train while minimizing device under-utilization.

You should determine the balance when defining a GPipe module, as balancing will not be done automatically. The module will be partitioned into multiple devices according to the given balance. You may rely on heuristics to find your own optimal configuration.

Parameters

• module (nn. Sequential) - sequential module to be parallelized

• **balance** (*ints*) – list of number of layers in each partition

Keyword Arguments

- devices (iterable of devices) devices to use (default: all CUDA devices)
- **chunks** (*int*) number of micro-batches (default: 1)
- **checkpoint** (*str*) when to enable checkpointing, one of 'always', 'except_last', or 'never' (default: 'except_last')
- **deferred_batch_norm** (*bool*) whether to use deferred BatchNorm moving statistics (default: False, See Deferred BatchNorm for more details)

Raises

- TypeError the module is not a nn.Sequential.
- ValueError invalid arguments, or wrong balance
- IndexError the number of devices is fewer than the number of partitions.

forward(input)

GPipe is a fairly transparent module wrapper. It doesn't modify the input and output signature of the underlying module. But there's type restriction. Input and output have to be a Tensor or a tuple of tensors. This restriction is applied at partition boundaries too.

Parameters input (tensor or tensors) - input mini-batch

Returns output mini-batch

Return type tensor or tensors

Raises TypeError – input is not a tensor or tensors.

devices

The devices mapped to each partition.

devices [-1] refers to the device of the last partition, which means it is the output device. Probably, you need to use it to transfer the target to calculate the loss without a device mismatch RuntimeError. For example:

```
out_device = gpipe.devices[-1]
for input, target in loader:
    target = target.to(out_device, non_blocking=True)
    output = gpipe(input)
    loss = F.cross_entropy(output, target)
```

2.3.2 Inspecting GPipe Timeline

torchgpipe.is_recomputing()

Whether if the current thread is under checkpoint recomputation.

Returns True if it's under checkpoint recomputation.

Return type bool

See also:

Detecting Recomputation

2.3.3 Automatic Balancing

torchgpipe_balancing.balance_by_time (module, canary, partitions, device, timeout) Balances the given sequential module by elapsed time per layer.

Parameters

- module (nn. Sequential) sequential module to be partitioned
- **sample** (*Tensor*) example input

Keyword Arguments

- **partitions** (*int*) intended number of partitions (default: 1)
- **device** (*torch.device*) CUDA device where the module is profiled (default: any related CUDA device or torch.device('cuda'))
- timeout (float) profiling iterates again if the timeout (in second) is not exceeded (default: 1.0)

Returns A list of number of layers in each partition. Use it for the balance parameter of GPipe.

torchgpipe_balancing.balance_by_size (module, canary, partitions, device) Balances the given sequential module by memory usage per layer.

Note: This function relies on torch.cuda.reset_max_memory_allocated() which is introduced at PyTorch 1.1. Therefore, it doesn't support neither CPU tensors nor PyTorch 1.0.x.

Parameters

- **module** (*nn*. *Sequential*) sequential module to be partitioned
- **sample** (*Tensor*) example input

Keyword Arguments

- **partitions** (*int*) intended number of partitions (default: 1)
- **device** (*torch.device*) CUDA device where the module is profiled (default: any related CUDA device or torch.device('cuda'))

Returns A list of number of layers in each partition. Use it for the balance parameter of GPipe.

2.4 Benchmarks

2.4.1 ResNet-101 Speed Benchmark

Experiment	Throughput	Speedup
naive-1	92.539 samples/sec	1.000x
pipeline-1	69.960 samples/sec	0.756x
pipeline-2	137.788 samples/sec	1.489x
pipeline-4	243.322 samples/sec	2.629x
pipeline-8	404.084 samples/sec	4.367x

The code is reproducible on Tesla P40 GPUs, and the experiment details can be found in examples/resnet101_speed_benchmark.

2.4.2 ResNet-101 Accuracy Benchmark

Experiment	Top-1 error (%)	
dataparallel-256	22.02±0.11	
dataparallel-1k	22.04 ± 0.24	
pipeline-256	21.99±0.13	
pipeline-1k	22.24±0.19	
pipeline-4k	22.13±0.09	

The code is reproducible on Tesla P40 GPUs, and the experiment details can be found in examples/resnet101_accuracy_benchmark.

2.4.3 AmoebaNet-D Speed Benchmark

Experiment	Experiment Throughput		
naive-2	14.188 samples/sec	1.000x	
pipeline-2	20.346 samples/sec	1.434x	
pipeline-4	29.074 samples/sec	2.049x	
pipeline-8	34.392 samples/sec	2.424x	

The code is reproducible on Tesla P40 GPUs, and the experiment details can be found in examples/amoebanetd_speed_benchmark.

2.4.4 AmoebaNet-D Memory Benchmark

Experi-	AmoebaNet-D	# of Model Param-	Total Model Parameter	Total Peak Activation
ment	(L, F)	eters	Memory	Memory
naive-1	(6, 208)	90M	1.00GB	-
pipeline-1	(6, 416)	358M	4.01GB	6.64GB
pipeline-2	(6, 544)	613M	6.45GB	11.31GB
pipeline-4	(12, 544)	1.16B	13.00GB	18.72GB
pipeline-8	(24, 512)	2.01B	22.42GB	35.78GB

2.5 Changelog

2.5.1 v0.0.3

Released on September 30, 2019.

Featured: torchgpipe now overlaps copy and computation using the separate CUDA streams. Previously, GPU could not compute a partition while copying micro-batches across different GPUs because they all happened on the same default CUDA stream.

Other Improvements:

- Added support for PyTorch 1.2.
- Redesigned the internal pipeline parallelism to represent dependencies transparently.

- Fixed the hanging issue when an exception is raised in a partition.
- Fixed the unintended size accumulation (issue #3 by Shiyan Deng) of balance_by_size().

Breaking Changes:

- No more support for PyTorch 1.0.
- Changed type of *GPipe.devices* from tuple to list.
- Removed current_microbatch. This approach turned out to be incompatible with checkpointing.

2.5.2 v0.0.2

Released on June 26, 2019.

- Added support for PyTorch 1.1.
- Refined public APIs.
- Detailed documentation.
- Proper exceptions for invalid usage.
- Provided automatic balancing.
- Provided inspecting utilities: <code>current_microbatch</code> (DO NOT USE, deprecated since v0.0.3) and <code>is_recomputing()</code>
- Reimplemented deferred batch normalization by subclassing.

2.5.3 v0.0.1

Released on May 14, 2019 to evaluate usability and efficiency internally.

- Provided a functional GPipe implementation, including pipeline parallelism, checkpointing, and deferred batch normalization.
- Supported Python 3.6+ and PyTorch 1.0.

CHAPTER 3

Authors and Licensing

This project is developed by Heungsub Lee, Myungryong Jeong, and Chiheon Kim at Kakao Brain, with Sungbin Lim, Ildoo Kim, and Woonhyuk Baek's help. It is distributed under Apache License 2.0.

If you apply this library to any project and research, please cite our code:

```
@misc{torchgpipe,
  author = {Kakao Brain},
  title = {torchgpipe, {A} {GPipe} implementation in {PyTorch}},
  howpublished = {\url{https://github.com/kakaobrain/torchgpipe}},
  year = {2019}
}
```

Index

В

balance_by_size() (in module torchgpipe_balancing), 13 balance_by_time() (in module torchgpipe_balancing), 13

D

devices (torchgpipe.GPipe attribute), 12

F

forward() (torchgpipe.GPipe method), 12

G

GPipe (class in torchgpipe), 11

I

is_recomputing() (in module torchgpipe), 12