



Tal Daniel

Tutorial 10 - Resource Efficiency - AMP, Quantization and Pruning



Agenda

- Resource Efficiency Motivation
 - Energy Required for Mathematical Operations
- Automatic Mixed Precision (AMP)
 - AMP in PyTorch
- Quantization
 - Quantization in PyTorch
 - Which Quantization Method to Use?
- Pruning
 - The Lottery Ticket Hypothesis
 - Pruning in PyTorch
- Recommended Videos
- Credits

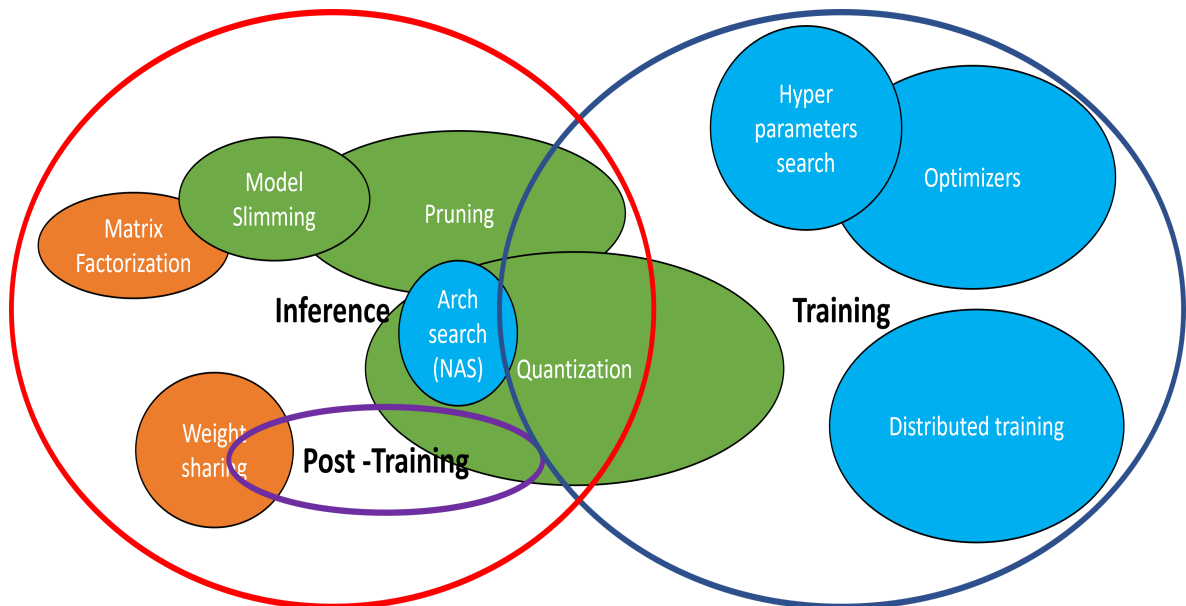
```
In [1]: # imports for the tutorial
import numpy as np
import time
import os
import gc

# pytorch imports
import torch
import torch.nn as nn
import torch.nn.functional as F
import torch.quantization
import torch.nn.utils.prune as prune
```



Resource Efficiency Motivation

- Training (large) neural networks is usually a long process that requires sufficient hardware.
- Hardware, such as GPUs, consumes a lot of energy which, for environmental reasons, can be spared.
- For real life applications, we ideally want fast inference times and we also want to deploy our models to, possibly, various devices such as mobile phones that don't have a matching hardware to the machine our models were trained on.
- And of course, these things COST MONEY, and usually lots of it.
- In this tutorial, we will discuss several approaches to incorporate compression during training and during inference.



Energy Required for Mathematical Operations

Operation	MUL	ADD
8-bit Integer	0.2 pJ	0.03 pJ
32-bit Integer	3.1 pJ	0.1 pJ
16-bit Floating Point	1.1 pJ	0.4 pJ
32-bit Floating Point	3.7 pJ	0.9 pJ

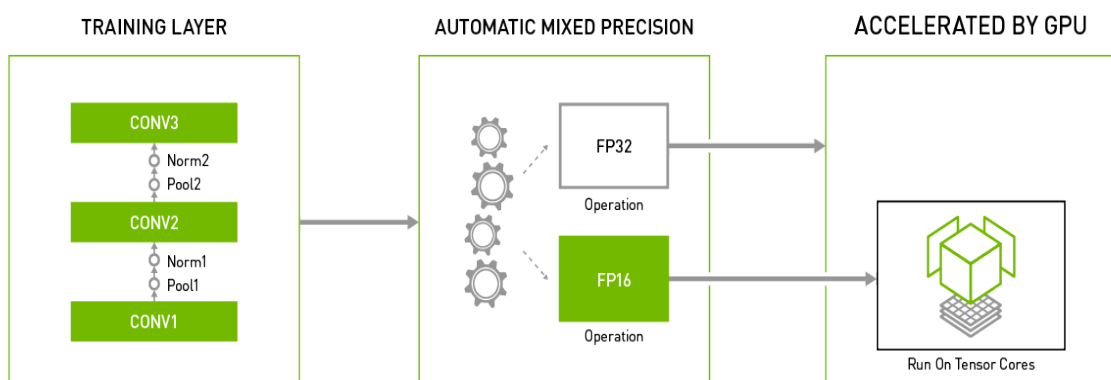
How can we utilize this information to make training and inference more efficient?

- pJ - pico-Joules = 10^{-12} Joules, 1 Watt = 1 Joule per second.
- [Source](#)



Train-time Efficiency - Automatic Mixed Precision (AMP)

- Deep Neural Network training has traditionally relied on FP32 (32-bit Floating Point, IEEE single-precision format).
- The (automatic) mixed precision technique - training with FP16 (16-bit Floating Point, half-precision) while maintaining the network accuracy achieved with FP32.
- Enabling mixed precision involves two steps:
 - Porting the model to use the half-precision data type **where appropriate**.
 - Using loss scaling to preserve small gradient values.



The Benefits of AMP

- Speeds up math-intensive operations, such as linear and convolution layers.
- Speeds up memory-limited operations by accessing half the bytes compared to single-precision.
- Reduces memory requirements for training models, enabling larger models or larger mini-batches.
- This feature enables automatic conversion of certain GPU operations from FP32 precision to mixed precision, thus improving performance while maintaining accuracy.

Performance of mixed precision training on NVIDIA 8xV100 vs. FP32 training on 8xV100 GPU



- Bars represent the speedup factor of V100 AMP over V100 FP32. The higher the better.

[Image Source](#)



AMP in PyTorch

- We will follow an [example](#) by [Michael Carilli](#).
- Note: this is only relevant for **GPU-powered machines**.
 - Mixed precision primarily benefits Tensor Core-enabled architectures (Volta, Turing, Ampere, Hopper), where one can observe significant (2-3X) speedup on those architectures. On earlier architectures (Kepler, Maxwell, Pascal), you may observe a modest speedup.
- `torch.cuda.amp` provides convenience methods for mixed precision, where some operations use the `torch.float32` (float) datatype and other operations use `torch.float16` (half).
- Some ops, like linear layers and convolutions, are much faster in `float16`, where other ops, like reductions, often require the dynamic range of `float32`.
- Mixed precision tries to match each op to its appropriate datatype, which can reduce your network's runtime and memory footprint.
- AMP - "automatic mixed precision training" uses `torch.cuda.amp.autocast` and `torch.cuda.amp.GradScaler` together, as we will soon explain.

```
In [2]: # Timing utilities
start_time = None

def start_timer():
    global start_time
    gc.collect()
    torch.cuda.empty_cache()
    torch.cuda.reset_max_memory_allocated()
    torch.cuda.synchronize()
    start_time = time.time()

def end_timer_and_print(local_msg):
    torch.cuda.synchronize()
    end_time = time.time()
    print("\n" + local_msg)
    print("Total execution time = {:.3f} sec".format(end_time - start_time))
    print("Max memory used by tensors = {} bytes".format(torch.cuda.max_memory_allocated()))
```

```
In [3]: # for our demonstrations, we will use a simple MLP network
device = torch.device("cuda:0" if torch.cuda.is_available() else "cpu")
def make_model(in_size, out_size, num_layers):
    layers = []
    for _ in range(num_layers - 1):
        layers.append(torch.nn.Linear(in_size, in_size))
        layers.append(torch.nn.ReLU())
    layers.append(torch.nn.Linear(in_size, out_size))
    return torch.nn.Sequential(*tuple(layers)).to(device)
```

- `batch_size`, `in_size`, `out_size`, and `num_layers` are chosen to be large enough to saturate the GPU with work.
- Typically, mixed precision provides the greatest speedup when the **GPU is saturated**.
- Small networks may be CPU bound, in which case mixed precision won't improve performance.
- Sizes are also chosen such that the participating dimensions of the linear layers are multiples of 8, to permit Tensor Core usage on Tensor Core-capable GPUs.

```
In [4]: batch_size = 512 # try, for example, 128, 256, 512.
in_size = 4096
out_size = 4096
num_layers = 3
num_batches = 50
epochs = 3

# creates data in default precision (fp32).
# the same data is used for both default and mixed precision trials below.
# you don't need to manually change inputs' dtype when enabling mixed precision.
data = [torch.randn(batch_size, in_size, device=device) for _ in range(num_batches)]
targets = [torch.randn(batch_size, out_size, device=device) for _ in range(num_batches)]

loss_fn = torch.nn.MSELoss().to(device)
```

Default Precision (FP32)

Without `torch.cuda.amp`, the following simple network executes all ops in default precision (`torch.float32`)

```
In [11]: net = make_model(in_size, out_size, num_layers)
opt = torch.optim.SGD(net.parameters(), lr=0.001)

start_timer()
for epoch in range(epochs):
    for input, target in zip(data, targets):
        output = net(input)
        loss = loss_fn(output, target)
        # opt.zero_grad() # set_to_none=True here can modestly improve performance
        loss.backward()
        opt.step()

end_timer_and_print("Default precision:")
```

C:\ProgramData\Anaconda3\envs\deep_learn\lib\site-packages\torch\cuda\memory.py:263: FutureWarning: torch.cuda.reset_max_memory_allocated now calls torch.cuda.reset_peak_memory_stats, which resets /all/ peak memory stats.
FutureWarning)

Default precision:
Total execution time = 13.958 sec
Max memory used by tensors = 1493322752 bytes

Adding Autocast

- Instances of `torch.cuda.amp.autocast` serve as context managers that allow regions of your script to run in mixed precision.
- In these regions, CUDA ops run in a `dtype` **chosen by autocast** to improve performance while maintaining accuracy.
- See the [Autocast Op Reference](#) for details on what precision autocast chooses for each op, and under what circumstances.

```
In [7]: for epoch in range(0): # 0 epochs, this section is for illustration only
    for input, target in zip(data, targets):
        # Runs the forward pass under autocast.
        with torch.cuda.amp.autocast():
            output = net(input)
            # output is float16 because linear layers autocast to float16.
            assert output.dtype is torch.float16

            loss = loss_fn(output, target)
```

```

    # loss is float32 because mse_loss layers autocast to float32.
    assert loss.dtype is torch.float32

    # exits autocast before backward().
    # backward passes under autocast are not recommended.
    # backward ops run in the same dtype autocast chose for corresponding forward ops.
    opt.zero_grad() # set_to_none=True here can modestly improve performance
    loss.backward()
    opt.step()

```

Adding GradScaler

- How do we deal with small values (which may get lost due to low precision) when we use automatic casting?
- Gradient scaling helps **prevent gradients with small magnitudes from flushing to zero** ("underflowing") when training with mixed precision.
- `torch.cuda.amp.GradScaler` performs the steps of gradient scaling conveniently.

```

In [8]: # constructs scaler once, at the beginning of the convergence run, using default args.
        # the same GradScaler instance should be used for the entire convergence run.
        # if you perform multiple convergence runs in the same script, each run should use
        # a dedicated fresh GradScaler instance. GradScaler instances are lightweight.

scaler = torch.cuda.amp.GradScaler()

for epoch in range(0): # 0 epochs, this section is for illustration only
    for input, target in zip(data, targets):
        # use autocast as before
        with torch.cuda.amp.autocast():
            output = net(input)
            loss = loss_fn(output, target)

        opt.zero_grad() # set_to_none=True here can modestly improve performance

        # scales loss. calls backward() on scaled loss to create scaled gradients.
        scaler.scale(loss).backward()

        # scaler.step() first unscales the gradients of the optimizer's assigned params.
        # if these gradients do not contain infs or NaNs, optimizer.step() is then called,
        # otherwise, optimizer.step() is skipped (!).
        scaler.step(opt)

        # Updates the scale for next iteration.
        scaler.update()

```

Putting it All Together: "Automatic Mixed Precision"

- The following also demonstrates `enabled`, an optional convenience argument to `autocast` and `GradScaler`.
- If `False`, `autocast` and `GradScaler`'s calls become no-ops, which allows switching between default precision and mixed precision without if/else statements.

```

In [13]: use_amp = True

net = make_model(in_size, out_size, num_layers)
opt = torch.optim.SGD(net.parameters(), lr=0.001)
scaler = torch.cuda.amp.GradScaler(enabled=use_amp) # notice the `enabled` parameter

start_timer()
for epoch in range(epochs):
    for input, target in zip(data, targets):
        # notice the `enabled` parameter
        with torch.cuda.amp.autocast(enabled=use_amp):
            output = net(input)
            loss = loss_fn(output, target)

        # set_to_none=True here can modestly improve performance, replace 0 with None (save mem)
        opt.zero_grad(set_to_none=True)
        scaler.scale(loss).backward()
        scaler.step(opt)
        scaler.update()

end_timer_and_print("Mixed precision:")

```

C:\ProgramData\Anaconda3\envs\deep_learn\lib\site-packages\torch\cuda\memory.py:263: FutureWarning: torch.cuda.reset_max_memory_allocated now calls torch.cuda.reset_peak_memory_stats, which resets /all/ peak memory stats. FutureWarning)

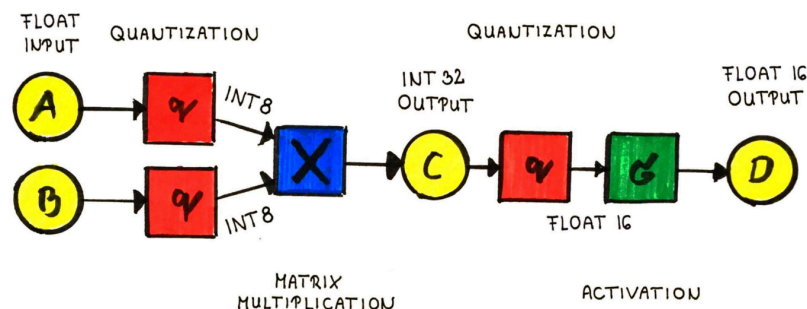
Mixed precision:
Total execution time = 14.081 sec
Max memory used by tensors = 1585620992 bytes

- Here we got similar performance due to the GPU architecture (which is older).
- For more advanced topics on AMP, [see here](#).



Inference/Train-time Efficiency - Quantization

- Quantization refers to techniques for doing both computations and memory accesses with lower precision data, usually `int8` compared to floating point implementations.
- Quantization leverages 8-bit integer (`int8`) instructions to reduce the model size and run the inference faster (reduced latency).
- This enables providing quick inference from a trained model and even fitting it into the resources available on a mobile device.
- Quantization allows for significant performance gains!
 - Up to 4x reduction in model size.
 - Up to 2-4x reduction in memory bandwidth.
 - Up to 2-4x faster inference due to savings in memory bandwidth and faster compute with `int8` arithmetic (the exact speed up varies depending on the hardware, the runtime, and the model).
- Quantization doesn't come without additional cost, as it means introducing approximations and the resulting networks have slightly less accuracy.
- These techniques attempt to **minimize the gap between the full floating point accuracy and the quantized accuracy**.



[Image Source](#)



Quantization in PyTorch

- Quantization is available in PyTorch in various flavors starting in version 1.3 and there are published quantized models for ResNet, ResNext, MobileNetV2, GoogleNet, InceptionV3 and ShuffleNetV2 in the PyTorch `torchvision >= 0.5` library.
- PyTorch has data types corresponding to quantized tensors, which share many of the features of tensors.
- PyTorch supports quantized modules for common operations as part of the `torch.nn.quantized` and `torch.nn.quantized.dynamic` name-space.
- Quantization is compatible with the rest of PyTorch: quantized models are traceable and scriptable. Quantized and floating point operations can be mixed in a model.
- Mapping of floating point tensors to quantized tensors is customizable with user defined observer/fake-quantization blocks. PyTorch provides default implementations that should work for most use cases.
- Currently the quantized models **can only be run on CPU**. However, it is possible to send the non-quantized parts of the model to a GPU.
 - GPU quantization is a work-in-progress, see [PTQ \(Post Training Quantization\)](#).



Image Source



The Three Types of Quantization

Dynamic Quantization

- Involves not just converting the weights to `int8` (as in all quantization variants), but also converting the **activations** to `int8` *on the fly*, just before doing the computation (hence “dynamic”).
- The computations will be performed using efficient `int8` matrix multiplication and convolution implementations, resulting in faster compute.
 - However, the activations are read and written to memory in floating point format.
- In PyTorch: `torch.quantization.quantize_dynamic` - takes in a **model**, as well as a couple other arguments, and produces a quantized model.

```
In [ ]: # dynamic quantization usage example
quantized_model = torch.quantization.quantize_dynamic(model, {torch.nn.Linear}, dtype=torch.qint8)
```

- Dynamic Quantization [documentation](#).
- Examples (project ideas!): [LSTM word model quantization](#), [pre-trained BERT quantization](#).

Post-Training Static Quantization

- Converting networks to use both integer arithmetic and `int8` memory accesses can improve the latency performance.
- Static quantization first feeds batches of data through the network and computes the resulting distributions of the different activations.
 - This is done by inserting “observer” modules at different points that record these distributions.
- This information is used to determine how specifically the different activations should be quantized at **inference time**.
 - A simple technique would be to simply divide the entire range of activations into 256 levels, but PyTorch supports more sophisticated methods as well.
- This step allows to pass quantized values between operations instead of converting these values to floats - and then back to ints - between every operation, resulting in a significant speed-up.
- Optimizing static quantization includes:
 - **Observers:** observer modules specify how statistics are collected prior to quantization to try out more advanced methods to quantize the data.
 - **Operator fusion:** fuse multiple operations into a single operation, saving on memory access while also improving the operation’s numerical accuracy.
 - **Per-channel quantization:** we can independently quantize weights for each output channel in a convolution/linear layer, which can lead to higher accuracy with almost the same speed.

```
In [ ]: # the three lines that perform post-training static quantization on the pre-trained model myModel
# set quantization config for server (x86) deployment
myModel.qconfig = torch.quantization.get_default_config('fbgemm')
```

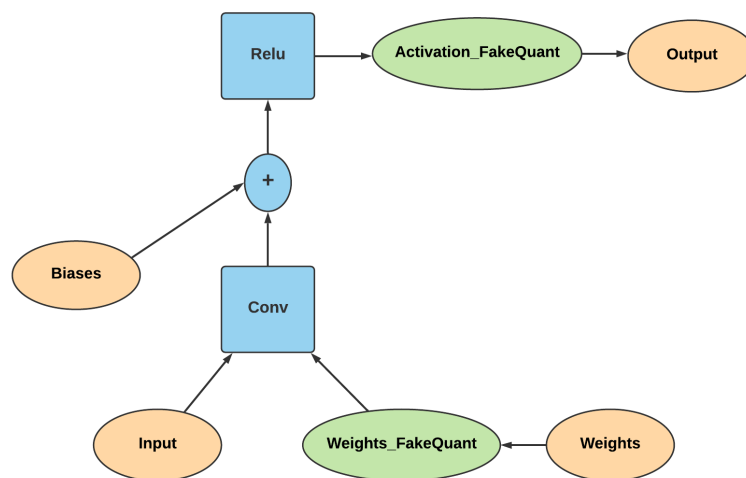
```
# insert observers
torch.quantization.prepare(myModel, inplace=True)
# Calibrate the model and collect statistics

# convert to quantized version
torch.quantization.convert(myModel, inplace=True)
```

- Examples (project ideas!): [Static quantization with eager execution](#), [Quantized transfer learning](#).

Quantization Aware Training (QAT)

- The quantization method that typically results in highest accuracy of the three methods.
- With QAT, all weights and activations are “fake quantized” during both the forward and backward passes of training: that is, **float values are rounded to mimic int8 values, but all computations are still done with floating point numbers**.
- Thus, all the weight adjustments during training are made while “aware” of the fact that the model will ultimately be quantized; after quantizing, therefore, this method usually yields higher accuracy than the other two methods.
- In PyTorch: `torch.quantization.prepare_qat` inserts fake quantization modules to model quantization and `torch.quantization.convert` actually quantizes the model once training is complete.



[Image Source](#)

```
In [ ]: # specify quantization config for QAT
qat_model.qconfig = torch.quantization.get_default_qat_qconfig('fbgemm')

# prepare QAT
torch.quantization.prepare_qat(qat_model, inplace=True)

# convert to quantized version, removing dropout, to check for accuracy on each epoch
quantized_model = torch.quantization.convert(qat_model.eval(), inplace=False)
```

- Example (project ideas): [Static quantization with eager execution](#).

Important Notes for Quantization in PyTorch

- Quantization support is restricted to a subset of available operators, depending on the method being used, for a list of supported operators, see the [documentation](#).
- The set of available operators and the quantization numerics also depend on the backend being used to run quantized models.
- Currently quantized operators are supported only for **CPU inference** in the following backends: x86 and ARM.
- QAT is typically only used in CNN models when post training static or dynamic quantization doesn't yield sufficient accuracy. This can occur with models that are highly optimized to achieve small size.



Which Quantization Method to Use?

Model Type	Preferred Scheme	Why
LSTM/RNN	Dynamic Quantization	Throughput dominated by compute/memory bandwidth for weights
BERT/Transformer	Dynamic Quantization	Throughput dominated by compute/memory bandwidth for weights
CNN	Static Quantization	Throughput limited by memory bandwidth for activations
CNN	Quantization Aware Training	In the case where accuracy can't be achieved with static quantization

Performance Results

Model Type	Float Latency (ms)	Quantized Latency (ms)	Inference Performance Gain	Accuracy	Device
BERT	581	313	1.8x	F1 score: 0.902 → 0.895 (dynamic quantization)	Xeon-D2191 (1.6GHz)
Resnet-50	214	103	2x	Top 1 Acc: 76.1 → 75.9 (static post-training quantization)	Xeon-D2191 (1.6GHz)
Mobilenet-v2	97	17	5.7x	Top 1 Acc: 71.9 → 71.6 (QAT)	Samsung S9



quanto: A PyTorch Quantization Toolkit

- Introduced by [HuggingFace](#), [quanto](#) is a versatile quantization toolkit, that provides several unique features, such as being available in *eager mode*, quantized models can be placed on *any device* (including CUDA and MPS), supports not only `int8` weights, but also `int2` and `int4`, supports not only `int8` activations, but also `float8` and more.
- `quanto` does not make a clear distinction between dynamic and static quantization; models are dynamically quantized first, but their weights can be "frozen" later to static values.
- [GitHub](#), `pip install quanto`.
- [Quantize Transformers with quanto notebook example](#).
- [LLM quantization notebook example](#).
- [More examples](#).

```
In [ ]: # quantization workflow with quanto
# 1. quantize - the first step converts a standard float model into a dynamically quantized model.
quantize(model, weights=quanto.qint8, activations=quanto.qint8)
# at this stage, the model's float weights are dynamically quantized only for inference.
# 2. calibrate (optional if activations are not quantized) -
# records the activation ranges while passing representative samples through the quantized model.
with calibration(momentum=0.9):
    model(samples)
# this automatically activates the quantization of the activations in the quantized modules.

In [ ]: # 3. tune, aka Quantization-Aware-Training (optional) -
# if the performance of the model degrades too much, one can tune it for a few epochs to try to recover the float
model.train()
for batch_idx, (data, target) in enumerate(train_loader):
    data, target = data.to(device), target.to(device)
    output = model(data).dequantize()
    loss = torch.nn.functional.nll_loss(output, target)
    optimizer.zero_grad()
    loss.backward()
    optimizer.step()

In [ ]: # 4. freeze integer weights - when freezing a model, its float weights are replaced by quantized integer weights.
freeze(model)
```



Quantization of Large Language Models

- Modern LLMs have a growing number of parameters, reaching billion and even trillion parameters, making them almost impossible to (fine-) tune or even use locally for inference.
- Quantization has emerged as a prominent tool to reduce the hardware requirements for using these models.

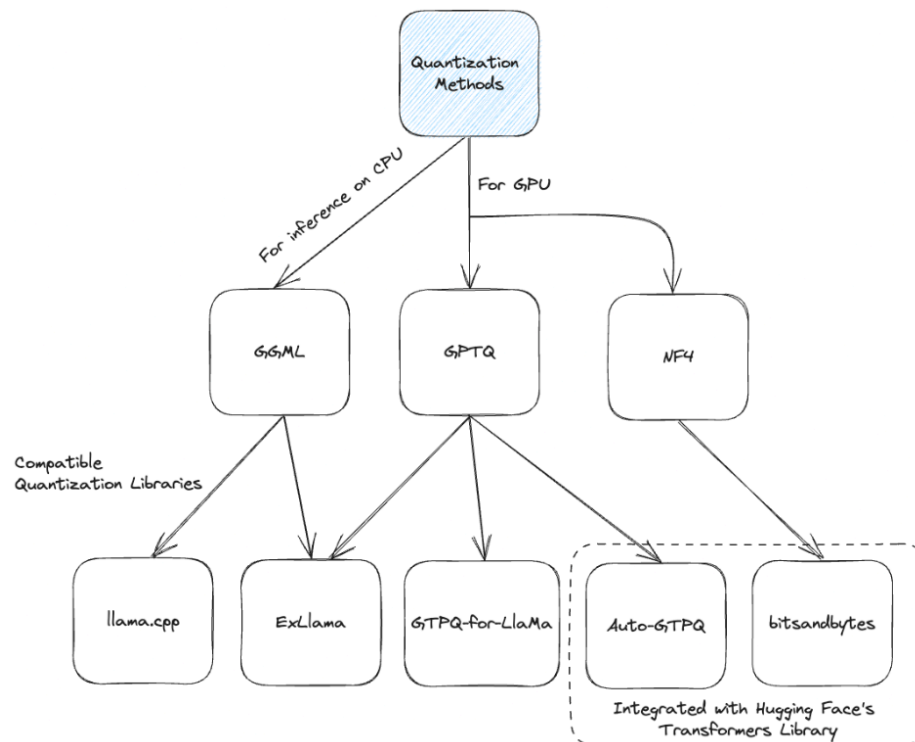
Model	Original Size (FP16)	Quantized Size (INT4)
Llama2-7B	13.5 GB	3.9 GB
Llama2-13B	26.1 GB	7.3 GB
Llama2-70B	138 GB	40.7 GB

[Image Source](#)

There are several methods for non-standard (i.e., the standard tools like mentioned above) quantization of LLMs.

- **GPTQ**: focuses mainly on GPU execution. The GPTQ quantization technique can be applied to many models to transform them into 3, 4 or 8-bit representations in a few simple steps. The [GitHub repository](#) contains several implementations of quantized LLMs such as LLaMa.
- **4-bit NormalFloat (NF4)**: the NormalFloat (NF) data type is an enhancement of the Quantile Quantization technique. It has shown better results than both 4-bit Integers and 4-bit Floats. Implemented in the [bitsandbytes](#) library it works closely with HuggingFace's `transformers` library. It is primarily used by [QLoRA](#) methods and loads models in 4-bit precision for fine-tuning.
- **GGML and llama.cpp**: allowing inference on CPU. This C library works closely with the [llama.cpp](#) library. It features a unique binary format for LLMs, allowing for fast loading and ease of reading.

Read more - [TensorOps - What are Quantized LLMs?](#)

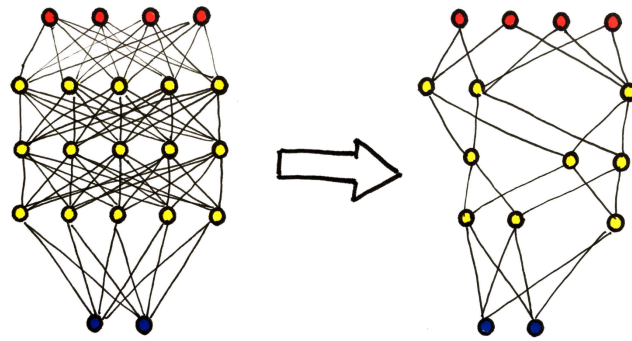


[Image Source](#)



Post Train-time Efficiency - Pruning

- Neural network pruning is the process of sparsifying neural networks from pre-trained dense neural networks, and is very similar to the pruning process of decision trees.
 - Instead of creating smaller trees, we create smaller computation graphs post-training.
- Pruning is used to investigate the differences in learning dynamics between over-parametrized and under-parametrized networks, to study the role of *lucky* sparse subnetworks and initializations ("lottery tickets") as a neural architecture search technique.



[Image Source](#)



The Lottery Ticket Hypothesis

The "Lottery Ticket Hypothesis" (Jonathan Frankle and Michael Carbin, 2008):

A randomly-initialized, dense neural network contains a subnetwork that is initialized such that — when trained in isolation — it can match the test accuracy of the original network after training for at most the same number of iterations.

- Such subnetworks are called *winning lottery tickets*.
- If this hypothesis is true, and such subnetworks can be found, training could be done much faster and cheaper, since a single iteration step would take less computation.

Iterative Pruning Algorithm to find *winning lottery tickets*:

1. Randomly initialize the network and store the initial weights for later reference.
2. Train the network for a given number of steps.
3. Remove a percentage of the weights (prune) with the **lowest magnitude**.
4. Restore the remaining weights to the value that was given during the first initialization.
5. Go to **Step 2** and iterate the pruning.



Pruning in PyTorch

- We will follow an [example](#) by [Michela Paganini](#).
- We will use PyTorch pruning name-space - `torch.nn.utils.prune`.

```
In [2]: # our model for the demonstration
device = torch.device("cuda:0" if torch.cuda.is_available() else "cpu")

class LeNet(nn.Module):
    def __init__(self):
        super(LeNet, self).__init__()
        # 1 input image channel, 6 output channels, 3x3 square conv kernel
        self.conv1 = nn.Conv2d(1, 6, 3)
        self.conv2 = nn.Conv2d(6, 16, 3)
        self.fc1 = nn.Linear(16 * 5 * 5, 120) # 5x5 image dimension
        self.fc2 = nn.Linear(120, 84)
        self.fc3 = nn.Linear(84, 10)

    def forward(self, x):
        x = F.max_pool2d(F.relu(self.conv1(x)), (2, 2))
        x = F.max_pool2d(F.relu(self.conv2(x)), 2)
        x = x.view(-1, int(x.nelement() / x.shape[0]))
        x = F.relu(self.fc1(x))
        x = F.relu(self.fc2(x))
        x = self.fc3(x)
```

```
        return x

model = LeNet().to(device=device)
```

To prune a module (in this example, the `conv1` layer of the LeNet architecture):

1. Select a pruning technique among those available in `torch.nn.utils.prune` (or implement your own by subclassing `BasePruningMethod`).
 2. Specify the module and the name of the parameter to prune within that module.
 3. Finally, using the adequate keyword arguments required by the selected pruning technique, specify the pruning parameters.
- In this example, we will prune **at random** 30% of the connections in the parameter named `weight` in the `conv1` layer.
 - The module is passed as the first argument to the function.
 - `name` identifies the parameter within that module using its string identifier.
 - `amount` indicates either the percentage of connections to prune (if it is a float between 0. and 1.), or the absolute number of connections to prune (if it is a non-negative integer).
 - `prune.random_unstructured(module, name="weight", amount=0.3)`
 - Pruning acts by removing `weight` from the parameters and replacing it with a new parameter called `weight_orig` (i.e. appending "_orig" to the initial parameter name). `weight_orig` stores the unpruned version of the tensor.
 - The `bias` was not pruned, so it will remain intact.

```
In [3]: module = model.conv1
print("before pruning:")
print(list(module.named_parameters()))
prune.random_unstructured(module, name="weight", amount=0.3)
print("after pruning:")
print(list(module.named_parameters()))
```

```

before pruning:
[('weight', Parameter containing:
tensor([[[[ 2.5489e-01, -9.4519e-02, -1.3123e-01],
           [ 1.4134e-01, -1.3794e-01,  1.0924e-01],
           [-1.6289e-01, -1.6713e-01, -9.5164e-02]]],

          [[ 3.0215e-01,  1.5124e-01, -2.6447e-01],
           [-1.6318e-01,  2.4991e-01, -2.0957e-01],
           [ 2.9987e-02,  2.8998e-01,  2.3750e-02]]],

          [[[-1.4108e-01, -2.9853e-01, -3.1772e-01],
           [ 4.8654e-02, -1.4157e-01, -1.5489e-01],
           [-2.6657e-01,  3.0692e-01,  2.9937e-01]]],

          [[[ 1.5138e-01, -1.5336e-01, -1.5590e-01],
           [ 2.8654e-01, -5.1841e-02, -1.7370e-01],
           [ 2.3252e-01, -2.8121e-01, -1.5625e-01]]],

          [[[ 2.1367e-01,  1.4964e-04,  1.8922e-01],
           [ 2.4092e-01,  1.6483e-01,  5.8223e-02],
           [ 3.1152e-01,  1.3158e-01,  6.3737e-02]]],

          [[[-1.0072e-01, -3.2249e-01,  2.4011e-02],
           [-1.0668e-01,  1.8761e-01, -1.6360e-02],
           [ 2.3007e-01,  2.5130e-01,  2.3430e-01]]]]], device='cuda:0',
requires_grad=True)), ('bias', Parameter containing:
tensor([ 0.3233, -0.2143,  0.0462,  0.0188,  0.0680,  0.1209], device='cuda:0',
requires_grad=True))]
after pruning:
[('bias', Parameter containing:
tensor([ 0.3233, -0.2143,  0.0462,  0.0188,  0.0680,  0.1209], device='cuda:0',
requires_grad=True)), ('weight_orig', Parameter containing:
tensor([[[[ 2.5489e-01, -9.4519e-02, -1.3123e-01],
           [ 1.4134e-01, -1.3794e-01,  1.0924e-01],
           [-1.6289e-01, -1.6713e-01, -9.5164e-02]]],

          [[ 3.0215e-01,  1.5124e-01, -2.6447e-01],
           [-1.6318e-01,  2.4991e-01, -2.0957e-01],
           [ 2.9987e-02,  2.8998e-01,  2.3750e-02]]],

          [[[-1.4108e-01, -2.9853e-01, -3.1772e-01],
           [ 4.8654e-02, -1.4157e-01, -1.5489e-01],
           [-2.6657e-01,  3.0692e-01,  2.9937e-01]]],

          [[[ 1.5138e-01, -1.5336e-01, -1.5590e-01],
           [ 2.8654e-01, -5.1841e-02, -1.7370e-01],
           [ 2.3252e-01, -2.8121e-01, -1.5625e-01]]],

          [[[ 2.1367e-01,  1.4964e-04,  1.8922e-01],
           [ 2.4092e-01,  1.6483e-01,  5.8223e-02],
           [ 3.1152e-01,  1.3158e-01,  6.3737e-02]]],

          [[[-1.0072e-01, -3.2249e-01,  2.4011e-02],
           [-1.0668e-01,  1.8761e-01, -1.6360e-02],
           [ 2.3007e-01,  2.5130e-01,  2.3430e-01]]]]], device='cuda:0',
requires_grad=True))]

```

- The pruning mask generated by the pruning technique selected above is saved as a **module buffer** named `weight_mask` (i.e. appending "_mask" to the initial parameter name).

```
In [4]: print(list(module.named_buffers()))
```

```
[('weight_mask', tensor([[[[0., 1., 0.],
    [1., 1., 1.],
    [1., 1., 1.]]],

    [[1., 1., 1.],
    [0., 1., 1.],
    [1., 1., 0.]]],

    [[1., 1., 0.],
    [0., 1., 1.],
    [1., 1., 0.]]],

    [[1., 1., 0.],
    [0., 0., 1.],
    [1., 1., 0.]]],

    [[1., 1., 1.],
    [0., 0., 1.],
    [1., 1., 1.]]],

    [[0., 1., 0.],
    [1., 1., 1.],
    [1., 1., 0.]]]]], device='cuda:0'))]
```

- For the forward pass to work *without modification*, the `weight` attribute needs to exist.
- The pruning techniques implemented in `torch.nn.utils.prune` compute the pruned version of the weight (by combining the mask with the original parameter) and store them in the attribute `weight`.
- Note, **this is no longer a parameter of the module**, it is now simply an attribute.

```
In [5]: print(module.weight)
```

```
tensor([[[[ 0.0000e+00, -9.4519e-02, -0.0000e+00],
    [ 1.4134e-01, -1.3794e-01,  1.0924e-01],
    [-1.6289e-01, -1.6713e-01, -9.5164e-02]]],

    [[[ 3.0215e-01,  1.5124e-01, -2.6447e-01],
    [-0.0000e+00,  2.4991e-01, -2.0957e-01],
    [ 2.9987e-02,  2.8998e-01,  0.0000e+00]]],

    [[[-1.4108e-01, -2.9853e-01, -0.0000e+00],
    [ 0.0000e+00, -1.4157e-01, -1.5489e-01],
    [-2.6657e-01,  3.0692e-01,  0.0000e+00]]],

    [[[ 1.5138e-01, -1.5336e-01, -0.0000e+00],
    [ 0.0000e+00, -0.0000e+00, -1.7370e-01],
    [ 2.3252e-01, -2.8121e-01, -0.0000e+00]]],

    [[[ 2.1367e-01,  1.4964e-04,  1.8922e-01],
    [ 0.0000e+00,  0.0000e+00,  5.8223e-02],
    [ 3.1152e-01,  1.3158e-01,  6.3737e-02]]],

    [[[-0.0000e+00, -3.2249e-01,  0.0000e+00],
    [-1.0668e-01,  1.8761e-01, -1.6360e-02],
    [ 2.3007e-01,  2.5130e-01,  0.0000e+00]]]]], device='cuda:0',
grad_fn=<MulBackward0>)
```

- Finally, pruning is applied **prior to each forward pass** using PyTorch's `forward_pre_hooks`.
- Specifically, when the `module` is pruned, as we have done here, it will acquire a `forward_pre_hook` for each parameter associated with it that gets pruned.
- In this case, since we have so far only pruned the original parameter named `weight`, only one hook will be present.

```
In [6]: print(module._forward_pre_hooks)
```

```
OrderedDict([(0, <torch.nn.utils.prune.RandomUnstructured object at 0x000001760943CC40>)])
```

- We can now prune the `bias` too, to see how the parameters, buffers, hooks, and attributes of the module change.

- Just for the sake of trying out another pruning technique, here we prune the **3 smallest entries** in the `bias` by **L1 norm**, as implemented in the `l1_unstructured` pruning function.

```
In [7]: prune.l1_unstructured(module, name="bias", amount=3)
print('named parameters:')
print(list(module.named_parameters())) # notice the bias_orig
print('named buffers:')
print(list(module.named_buffers())) # notice the bias_mask
print('module.bias:')
print(module.bias)
print('forward pre hooks:')
print(module._forward_pre_hooks)
```

```

named parameters:
[('weight_orig', Parameter containing:
tensor([[[[ 2.5489e-01, -9.4519e-02, -1.3123e-01],
           [ 1.4134e-01, -1.3794e-01,  1.0924e-01],
           [-1.6289e-01, -1.6713e-01, -9.5164e-02]]],

          [[ 3.0215e-01,  1.5124e-01, -2.6447e-01],
           [-1.6318e-01,  2.4991e-01, -2.0957e-01],
           [ 2.9987e-02,  2.8998e-01,  2.3750e-02]]],

          [[[-1.4108e-01, -2.9853e-01, -3.1772e-01],
           [ 4.8654e-02, -1.4157e-01, -1.5489e-01],
           [-2.6657e-01,  3.0692e-01,  2.9937e-01]]],

          [[ 1.5138e-01, -1.5336e-01, -1.5590e-01],
           [ 2.8654e-01, -5.1841e-02, -1.7370e-01],
           [ 2.3252e-01, -2.8121e-01, -1.5625e-01]]],

          [[ 2.1367e-01,  1.4964e-04,  1.8922e-01],
           [ 2.4092e-01,  1.6483e-01,  5.8223e-02],
           [ 3.1152e-01,  1.3158e-01,  6.3737e-02]]],

          [[[-1.0072e-01, -3.2249e-01,  2.4011e-02],
           [-1.0668e-01,  1.8761e-01, -1.6360e-02],
           [ 2.3007e-01,  2.5130e-01,  2.3430e-01]]]], device='cuda:0',
requires_grad=True)), ('bias_orig', Parameter containing:
tensor([ 0.3233, -0.2143,  0.0462,  0.0188,  0.0680,  0.1209], device='cuda:0',
requires_grad=True))]
named buffers:
[('weight_mask', tensor([[[[0., 1., 0.],
                           [1., 1., 1.],
                           [1., 1., 1.]]],

                           [[1., 1., 1.],
                           [0., 1., 1.],
                           [1., 1., 0.]]],

                           [[1., 1., 0.],
                           [0., 1., 1.],
                           [1., 1., 0.]]],

                           [[1., 1., 0.],
                           [0., 0., 1.],
                           [1., 1., 0.]]],

                           [[1., 1., 1.],
                           [0., 0., 1.],
                           [1., 1., 1.]]],

                           [[0., 1., 0.],
                           [1., 1., 1.],
                           [1., 1., 0.]]]], device='cuda:0')), ('bias_mask', tensor([1., 1., 0., 0., 0., 1.], device='cuda:0'))]
module.bias:
tensor([ 0.3233, -0.2143,  0.0000,  0.0000,  0.0000,  0.1209], device='cuda:0',
grad_fn=<MulBackward0>)
forward pre hooks:
OrderedDict([(0, <torch.nn.utils.prune.RandomUnstructured object at 0x000001760943CC40>), (1, <torch.nn.utils.prune.L1Unstructured object at 0x0000017609450910>)])

```

- You can now try different methods to prune your trained model.
- For more examples (project ideas!): [Iterative Pruning](#), [Pruning Multiple Parameters in a Model](#), [Global Pruning \(pruning all weights at once\)](#), [Custom Pruning Functions](#).
- Making the pruning **permanent** by removing pruning re-parametrization.
 - This means removing the re-parametrization in terms of `weight_orig` and `weight_mask`, and remove the `forward_pre_hook`.
- We can use the `remove` functionality from `torch.nn.utils.prune`.

- Note that this doesn't undo the pruning, as if it never happened. It simply makes it permanent, instead, by reassigning the parameter weight to the model parameters, in its pruned version.

```
In [8]: prune.remove(module, 'weight')
print(list(module.named_parameters())) # only bias_orig remains
print(list(module.named_buffers())) # only bias_mask remains

[('bias_orig', Parameter containing:
tensor([ 0.3233, -0.2143,  0.0462,  0.0188,  0.0680,  0.1209], device='cuda:0',
        requires_grad=True)), ('weight', Parameter containing:
tensor([[[[ 0.0000e+00, -9.4519e-02, -0.0000e+00],
          [ 1.4134e-01, -1.3794e-01,  1.0924e-01],
          [-1.6289e-01, -1.6713e-01, -9.5164e-02]],

          [[ 3.0215e-01,  1.5124e-01, -2.6447e-01],
            [-0.0000e+00,  2.4991e-01, -2.0957e-01],
            [ 2.9987e-02,  2.8998e-01,  0.0000e+00]],

          [[-1.4108e-01, -2.9853e-01, -0.0000e+00],
            [ 0.0000e+00, -1.4157e-01, -1.5489e-01],
            [-2.6657e-01,  3.0692e-01,  0.0000e+00]],

          [[ 1.5138e-01, -1.5336e-01, -0.0000e+00],
            [ 0.0000e+00, -0.0000e+00, -1.7370e-01],
            [ 2.3252e-01, -2.8121e-01, -0.0000e+00]],

          [[ 2.1367e-01,  1.4964e-04,  1.8922e-01],
            [ 0.0000e+00,  0.0000e+00,  5.8223e-02],
            [ 3.1152e-01,  1.3158e-01,  6.3737e-02]],

          [[-0.0000e+00, -3.2249e-01,  0.0000e+00],
            [-1.0668e-01,  1.8761e-01, -1.6360e-02],
            [ 2.3007e-01,  2.5130e-01,  0.0000e+00]]]], device='cuda:0',
        requires_grad=True))
[('bias_mask', tensor([1., 1., 0., 0., 0., 1.], device='cuda:0'))]
```



Recommended Videos



Warning!

- These videos do not replace the lectures and tutorials.
- Please use these to get a better understanding of the material, and not as an alternative to the written material.

Video By Subject

- Automatic Mixed Precision (AMP) - [NVIDIA - Automatic Mixed Precision Training in PyTorch](#)
- Quantization - [Deep Dive on PyTorch Quantization - Chris Gottbrath](#)
- GPTQ - [Deci AI - GPTQ Tutorial: Shrinking LLMs without Quality Loss](#)
- Pruning - [Neural Network Pruning for Compression and Understanding - Facebook AI Research - Dr. Michela Paganini](#)
 - Pruning - [PyTorch Pruning - How it's Made by Michela Paganini](#)



Credits

- Icons made by [Becris](#) from www.flaticon.com
- Icons from [Icons8.com](https://icons8.com) - <https://icons8.com>
- NVIDIA - [Mixed Precision](#)
- [Introduction to Quantization in PyTorch](#) - by Raghuraman Krishnamoorthi, James Reed, Min Ni, Chris Gottbrath, and Seth Weidman.
- Tivadar Danko - [How to Compress a Neural Network](#)
- [Michela Paganini - Pruning Tutorial](#)
- TensorOps - [What are Quantized LLMs?](#)