

pytorch - GPU & Transformers

Lecture 26

Dr. Colin Rundel

CUDA

CUDA (or Compute Unified Device Architecture) is a parallel computing platform and application programming interface (API) that allows software to use certain types of graphics processing unit (GPU) for general purpose processing, an approach called general-purpose computing on GPUs (GPGPU). CUDA is a software layer that gives direct access to the GPU's virtual instruction set and parallel computational elements, for the execution of compute kernels.

Core libraries:

- cuBLAS
- cuTENSOR
- Thrust
- cuSOLVER
- cuFFT
- cuDNN
- cuSPARSE
- cuRAND

CUDA Kernels

```
1 // Kernel - Adding two matrices MatA and MatB
2 __global__ void MatAdd(float MatA[N][N], float MatB[N][N], float MatC[N][N])
3 {
4     int i = blockIdx.x * blockDim.x + threadIdx.x;
5     int j = blockIdx.y * blockDim.y + threadIdx.y;
6     if (i < N && j < N)
7         MatC[i][j] = MatA[i][j] + MatB[i][j];
8 }
9
10 int main()
11 {
12     ...
13     // Matrix addition kernel launch from host code
14     dim3 threadsPerBlock(16, 16);
15     dim3 numBlocks(
16         (N + threadsPerBlock.x - 1) / threadsPerBlock.x,
17         (N+threadsPerBlock.y - 1) / threadsPerBlock.y
18     );
19
20     MatAdd<<<numBlocks, threadsPerBlock>>>(MatA, MatB, MatC);
21     ...
22 }
```

GPU Status

```
1 nvidia-smi
```

```
Tue Apr 14 09:11:37 2026
```

```
+-----+
| NVIDIA-SMI 590.48.01                Driver Version: 590.48.01          CUDA Version: 13.1         |
+-----+-----+-----+-----+-----+-----+
| GPU  Name                Persistence-M | Bus-Id                Disp.A | Volatile Uncorr. ECC |
| Fan  Temp   Perf          Pwr:Usage/Cap |      Memory-Usage     | GPU-Util  Compute M. |
|                                           |                       | MIG M.     |
+-----+-----+-----+-----+-----+-----+
|   0   NVIDIA RTX A4000             Off | 00000000:01:00.0 Off |           Off        |
| 41%   35C    P8              14W / 140W | 2699MiB / 16376MiB |      0%      Default |
|                                           |                       | N/A        |
+-----+-----+-----+-----+-----+-----+
|   1   NVIDIA RTX A4000             Off | 00000000:68:00.0 Off |           Off        |
| 41%   37C    P8              13W / 140W |    4MiB / 16376MiB |      0%      Default |
|                                           |                       | N/A        |
+-----+-----+-----+-----+-----+-----+

+-----+
| Processes:                                |
| GPU   GI    CI          PID    Type    Process name                        GPU Memory |
|      ID    ID                                   |            Usage                    |
+-----+-----+-----+-----+-----+-----+
|
```

Torch GPU Information

```
1 torch.cuda.is_available()
```

True

```
1 torch.cuda.device_count()
```

2

```
1 torch.cuda.get_device_name("cuda:0")
```

'NVIDIA RTX A4000'

```
1 torch.cuda.get_device_name("cuda:1")
```

'NVIDIA RTX A4000'

```
1 torch.cuda.get_device_properties(0)
```

_CudaDeviceProperties(name='NVIDIA RTX A4000', major=8, minor=6, total_memory=15976MB, multi_proce

GPU Tensors

Usage of the GPU is governed by the location of the Tensors - to use the GPU we allocate them on the GPU device.

```
1 cpu = torch.device('cpu')
2 cuda0 = torch.device('cuda:0')
3 cuda1 = torch.device('cuda:1')
```

```
1 x = torch.linspace(0,1,5, device=cuda0); x
```

```
tensor([0.0000, 0.2500, 0.5000, 0.7500,
        1.0000], device='cuda:0')
```

```
1 y = torch.randn(5,2, device=cuda0); y
```

```
tensor([[ -3.3676,  0.4284],
        [ 1.0317, -0.3296],
        [-1.6607, -0.5143],
        [-0.7477, -0.5811],
        [ 1.3505,  0.1542]], device='cuda:0')
```

```
1 z = torch.rand(2,3, device=cpu); z
```

```
tensor([[0.8394, 0.5169, 0.5356],
        [0.5584, 0.5428, 0.0929]])
```

```
1 x @ y
```

```
tensor([ 0.2173, -0.6211], device='cuda:0')
```

```
1 y @ z
```

```
RuntimeError: Expected all tensors to be on the
```

```
1 y @ z.to(cuda0)
```

```
tensor([[ -2.5874, -1.5080, -1.7639],
        [ 0.6819,  0.3544,  0.5219],
        [-1.6811, -1.1375, -0.9373],
        [-0.9521, -0.7019, -0.4545],
        [ 1.2197,  0.7817,  0.7377]],
        device='cuda:0')
```

NN Layers + GPU

NN layers (parameters) also need to be assigned to the GPU to be used with GPU tensors,

```
1 nn = torch.nn.Linear(5, 5)
2 X = torch.randn(10, 5).cuda()
```

```
1 nn(X)
```

RuntimeError: Expected all tensors to be on the same device, but got mat1 is on cuda:0, different

```
1 nn.cuda()(X)
```

```
tensor([[ 0.0665, -0.1522, -0.1733,  0.4627,
         -0.3899],
        [ 0.4550,  0.2700, -0.8572,  0.0812,
          0.5769],
        [-0.1594,  0.1351, -0.3831, -0.0761,
          0.4204],
        [-0.5051,  0.6921,  0.0977, -0.0305,
          0.3425],
        [-0.9259,  0.9657, -0.2817,  1.1706,
         -0.3354],
        [ 0.8163, -0.4542,  0.6201,  0.4417,
         -0.9929],
        [ 0.7821, -0.2014,  1.0838, -1.0487,
          0.1281],
        [-0.1962,  0.5809,  0.0408,  0.9464,
         -0.1634],
        [-0.1362,  0.3502, -0.5095,  0.3304,
         -0.1045],
        [ 0.4283,  0.1325,  0.3059, -0.9481,
          0.5690]], device='cuda:0',
grad_fn=<AddmmBackward0>)
```

```
1 nn.to(device="cuda")(X)
```

```
tensor([[ 0.0665, -0.1522, -0.1733,  0.4627,
         -0.3899],
        [ 0.4550,  0.2700, -0.8572,  0.0812,
          0.5769],
        [-0.1594,  0.1351, -0.3831, -0.0761,
          0.4204],
        [-0.5051,  0.6921,  0.0977, -0.0305,
          0.3425],
        [-0.9259,  0.9657, -0.2817,  1.1706,
         -0.3354],
        [ 0.8163, -0.4542,  0.6201,  0.4417,
         -0.9929],
        [ 0.7821, -0.2014,  1.0838, -1.0487,
          0.1281],
        [-0.1962,  0.5809,  0.0408,  0.9464,
         -0.1634],
        [-0.1362,  0.3502, -0.5095,  0.3304,
         -0.1045],
        [ 0.4283,  0.1325,  0.3059, -0.9481,
          0.5690]], device='cuda:0',
grad_fn=<AddmmBackward0>)
```

CIFAR10

homepage

Image Classification Benchmarks

CIFAR-10

- 60,000 color images (32×32)
- 10 classes (airplane, car, bird, cat, deer, dog, frog, horse, ship, truck)
- 50k train / 10k test
- ~170 MB

CIFAR-100

- 60,000 color images (32×32)
- 100 classes grouped into 20 superclasses
- 50k train / 10k test
- ~170 MB

ImageNet (ILSVRC)

- ~1.2 million color images (variable, typically 224×224)
- 1,000 classes
- ~1.2M train / 50k val / 100k test
- ~150 GB

All three are standard benchmarks for evaluating CNN architectures — CIFAR-10/100 are common for rapid prototyping due to their small image size, while ImageNet is the large-scale challenge used to evaluate models (AlexNet, VGG, ResNet, etc.).

Loading the data

```
1 import torchvision
2 training_data = torchvision.datasets.CIFAR10(
3     root="/data",
4     train=True,
5     download=True,
6     transform=torchvision.transforms.ToTensor()
7 )
8 test_data = torchvision.datasets.CIFAR10(
9     root="/data",
10    train=False,
11    download=True,
12    transform=torchvision.transforms.ToTensor()
13 )
```

Downloads data to “/data/cifar-10-batches-py” which is ~178M on disk.

CIFAR10 data

```
1 training_data.classes
```

```
['airplane', 'automobile', 'bird', 'cat', 'deer', 'dog', 'frog', 'horse', 'ship', 'truck']
```

```
1 training_data.data.shape
```

```
(50000, 32, 32, 3)
```

```
1 test_data.data.shape
```

```
(10000, 32, 32, 3)
```

```
1 training_data[0]
```

```
(tensor([[ [0.2314, 0.1686, 0.1961, 0.2667,  
           0.3843, 0.4667, 0.5451, 0.5686,  
           0.5843, 0.5843, 0.5137, 0.4902,  
           0.5569, 0.5647, 0.5373, 0.5059,  
           0.5373, 0.5255, 0.4863, 0.5451,  
           0.5451, 0.5216, 0.5333, 0.5451,  
           0.5961, 0.6392, 0.6588, 0.6235,  
           0.6196, 0.6196, 0.5961, 0.5804],  
 [0.0627, 0.0000, 0.0706, 0.2000,  
           0.3451, 0.4706, 0.5020, 0.4980,  
           0.4941, 0.4549, 0.4157, 0.3961,  
           0.4118, 0.4431, 0.4275, 0.4392,  
           0.4667, 0.4275, 0.4118, 0.4902,  
           0.4980, 0.4784, 0.5137, 0.4863,  
           0.4745, 0.5137, 0.5176, 0.5216,  
           0.5216, 0.4824, 0.4667, 0.4784]],
```

[0.0980, 0.0627, 0.1922, 0.3255,
0.4314, 0.5059, 0.5098, 0.4745,
0.4431, 0.4392, 0.4392, 0.4157,
0.4118, 0.5020, 0.4863, 0.5098,
0.4980, 0.4784, 0.4510, 0.4706,
0.5098, 0.5137, 0.5451, 0.4980,
- - - - -

Example data

frog



truck



truck



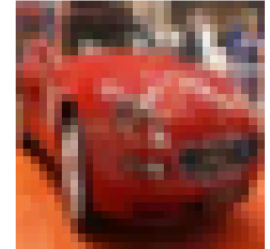
deer



automobile



automobile



bird



horse



ship



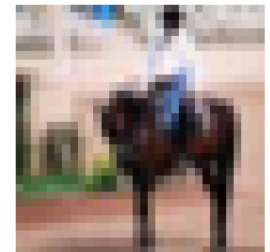
cat



deer



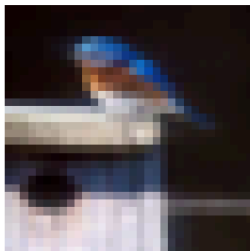
horse



horse



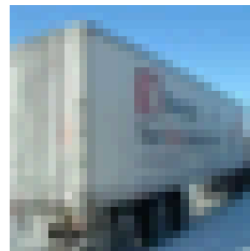
bird



truck



truck



truck



cat



bird



frog



deer



cat



frog



frog



Data Loaders

Torch handles large datasets and minibatches through the use of the `DataLoader` class,

```
1 training_loader = torch.utils.data.DataLoader(  
2     training_data,  
3     batch_size=100,  
4     shuffle=True,  
5     num_workers=4,  
6     pin_memory=True  
7 )  
8  
9 test_loader = torch.utils.data.DataLoader(  
10     test_data,  
11     batch_size=100,  
12     shuffle=True,  
13     num_workers=4,  
14     pin_memory=True  
15 )
```

Loader as generator

The resulting `DataLoader` class is iterable and “yields” the features and targets for each batch when iterated over,

```
1 training_loader
```

```
<torch.utils.data.dataloader.DataLoader object at 0x7f2d82a3e900>
```

```
1 X, y = next(iter(training_loader))  
2 X.shape
```

```
torch.Size([100, 3, 32, 32])
```

```
1 y.shape
```

```
torch.Size([100])
```

Custom Datasets

In this case we got our data (`training_data` and `test_data`) directly from torchvision which gave us a `dataset` object to use with our `DataLoader`. If we do not have a `Dataset` object then we need to create a custom class for our data telling torch how to load it.

Your class must define the methods: `__init__()`, `__len__()`, and `__getitem__()`.

```
1 class data(torch.utils.data.Dataset):
2     def __init__(self, X, y):
3         self.X = X
4         self.y = y
5
6     def __len__(self):
7         return len(self.X)
8
9     def __getitem__(self, idx):
10        return self.X[idx], self.y[idx]
```

CIFAR CNN

```
1 class cifar_conv_model(torch.nn.Module):
2     def __init__(self, device):
3         super().__init__()
4         self.device = torch.device(device)
5         self.epoch = 0
6         self.model = torch.nn.Sequential(
7             torch.nn.Conv2d(3, 6, kernel_size=5),
8             torch.nn.ReLU(),
9             torch.nn.MaxPool2d(2, 2),
10            torch.nn.Conv2d(6, 16, kernel_size=5),
11            torch.nn.ReLU(),
12            torch.nn.MaxPool2d(2, 2),
13            torch.nn.Flatten(),
14            torch.nn.Linear(16 * 5 * 5, 120),
15            torch.nn.ReLU(),
16            torch.nn.Linear(120, 84),
17            torch.nn.ReLU(),
18            torch.nn.Linear(84, 10)
19        ).to(device=self.device)
20
21    def forward(self, X):
22        return self.model(X)
```

CNN Performance

CPU - 1 step

GPU - 1 step

CPU - 1 epoch

GPU - 1 epoch

```
1 X, y = next(iter(training_loader))
2
3 m_cpu = cifar_conv_model(device="cpu")
4 tmp = m_cpu(X)
5
6 with torch.autograd.profiler.profile(with_stack=True) as prof_cpu:
7     tmp = m_cpu(X)
```

```
1 print(prof_cpu.key_averages().table(sort_by='self_cpu_time_total', row_limit=10))
```

Name	Self CPU %	Self CPU	CPU total %	CPU total	CPU time avg	# of Calls
aten::mkldnn_convolution	44.61%	552.034us	45.78%	566.422us	283.211us	2
aten::max_pool2d_with_indices	31.12%	385.069us	31.12%	385.069us	192.535us	2
aten::clamp_min	9.45%	116.901us	9.45%	116.901us	29.225us	4
aten::addmm	6.38%	78.921us	7.59%	93.919us	31.306us	3
aten::convolution	1.41%	17.472us	47.93%	593.023us	296.511us	2
aten::relu	1.06%	13.144us	10.51%	130.045us	32.511us	4
aten::copy_	0.91%	11.272us	0.91%	11.272us	3.757us	3
aten::_convolution	0.74%	9.129us	46.52%	575.551us	287.775us	2
aten::empty	0.65%	8.076us	0.65%	8.076us	2.019us	4
aten::view	0.57%	7.064us	0.57%	7.064us	7.064us	1

Self CPU time total: 1.237ms

Loaders & Accuracy

```
1 def accuracy(model, loader, device):
2     total, correct = 0, 0
3     with torch.no_grad():
4         for X, y in loader:
5             X, y = X.to(device=device), y.to(device=device)
6             pred = model(X)
7             # the class with the highest energy is what we choose as prediction
8             val, idx = torch.max(pred, 1)
9             total += pred.size(0)
10            correct += (idx == y).sum().item()
11
12    return correct / total
```

Model fitting

```
1 m = cifar_conv_model("cuda")
```

Round 1 Round 2 Round 3

```
1 m.fit(training_loader, epochs=10, n_report=500, lr=0.01)
```

```
## [Epoch 1, Minibatch 500] loss: 2.098
## [Epoch 2, Minibatch 500] loss: 1.692
## [Epoch 3, Minibatch 500] loss: 1.482
## [Epoch 4, Minibatch 500] loss: 1.374
## [Epoch 5, Minibatch 500] loss: 1.292
## [Epoch 6, Minibatch 500] loss: 1.226
## [Epoch 7, Minibatch 500] loss: 1.173
## [Epoch 8, Minibatch 500] loss: 1.117
## [Epoch 9, Minibatch 500] loss: 1.071
## [Epoch 10, Minibatch 500] loss: 1.035
```

```
1 accuracy(m, training_loader, "cuda")
```

```
## 0.63444
```

```
1 accuracy(m, test_loader, "cuda")
```

```
## 0.572
```

The VGG16 model

```
1 class VGG16(torch.nn.Module):
2     def make_layers(self):
3         cfg = [64, 64, 'M', 128, 128, 'M', 256, 256, 256, 'M', 512, 512, 512, 'M', 512, 512, 512, 'M']
4         layers = []
5         in_channels = 3
6         for x in cfg:
7             if x == 'M':
8                 layers += [torch.nn.MaxPool2d(kernel_size=2, stride=2)]
9             else:
10                layers += [torch.nn.Conv2d(in_channels, x, kernel_size=3, padding=1),
11                           torch.nn.BatchNorm2d(x),
12                           torch.nn.ReLU(inplace=True)]
13                in_channels = x
14            layers += [
15                torch.nn.AvgPool2d(kernel_size=1, stride=1),
16                torch.nn.Flatten(),
17                torch.nn.Linear(512, 10)
18            ]
19
20        return torch.nn.Sequential(*layers).to(self.device)
21
22    def __init__(self, device):
23        super().__init__()
24        self.device = torch.device(device)
```

Model

```
1 VGG16("cpu").model
```

```
Sequential(
  (0): Conv2d(3, 64, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1))
  (1): BatchNorm2d(64, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
  (2): ReLU(inplace=True)
  (3): Conv2d(64, 64, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1))
  (4): BatchNorm2d(64, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
  (5): ReLU(inplace=True)
  (6): MaxPool2d(kernel_size=2, stride=2, padding=0, dilation=1, ceil_mode=False)
  (7): Conv2d(64, 128, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1))
  (8): BatchNorm2d(128, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
  (9): ReLU(inplace=True)
  (10): Conv2d(128, 128, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1))
  (11): BatchNorm2d(128, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
  (12): ReLU(inplace=True)
  (13): MaxPool2d(kernel_size=2, stride=2, padding=0, dilation=1, ceil_mode=False)
  (14): Conv2d(128, 256, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1))
  (15): BatchNorm2d(256, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
  (16): ReLU(inplace=True)
  (17): Conv2d(256, 256, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1))
  (18): BatchNorm2d(256, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
  (19): ReLU(inplace=True)
  (20): Conv2d(256, 256, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1))
  ...
```

VGG16 performance

CPU - 1 forward step

GPU - 1 forward step

```
1 X, y = next(iter(training_loader))
2 m_cpu = VGG16(device="cpu")
3 tmp = m_cpu(X)
4
5 with torch.autograd.profiler.profile(with_stack=True) as prof_cpu:
6     tmp = m_cpu(X)
```

```
1 print(prof_cpu.key_averages().table(sort_by='self_cpu_time_total', row_limit=10))
```

Name	Self CPU %	Self CPU	CPU total %	CPU total	CPU tim
aten::mkldnn_convolution	77.21%	53.228ms	77.45%	53.398ms	4.
aten::native_batch_norm	11.68%	8.050ms	11.86%	8.175ms	628.
aten::max_pool2d_with_indices	7.60%	5.241ms	7.60%	5.241ms	1.
aten::clamp_min_	1.91%	1.315ms	1.91%	1.315ms	101.
aten::empty	0.33%	227.335us	0.33%	227.335us	1.
aten::add_	0.22%	151.716us	0.22%	151.716us	11.
aten::convolution	0.21%	145.065us	77.81%	53.647ms	4.
aten::relu_	0.20%	137.109us	2.11%	1.452ms	111.
aten::_convolution	0.15%	103.739us	77.60%	53.502ms	4.
aten::_batch_norm_impl_index	0.12%	83.299us	11.99%	8.264ms	635.

Self CPU time total: 68.941ms

Fitting

```
lr = 0.01
```

```
lr = 0.001
```

```
1 m = VGG16(device="cuda")
2 fit(m, training_loader, epochs=10, n_report=500, lr=0.01)
```

```
## [Epoch 1, Minibatch 500] loss: 1.345
## [Epoch 2, Minibatch 500] loss: 0.790
## [Epoch 3, Minibatch 500] loss: 0.577
## [Epoch 4, Minibatch 500] loss: 0.445
## [Epoch 5, Minibatch 500] loss: 0.350
## [Epoch 6, Minibatch 500] loss: 0.274
## [Epoch 7, Minibatch 500] loss: 0.215
## [Epoch 8, Minibatch 500] loss: 0.167
## [Epoch 9, Minibatch 500] loss: 0.127
## [Epoch 10, Minibatch 500] loss: 0.103
```

```
1 accuracy(model=m, loader=training_loader, device="cuda")
```

```
## 0.97008
```

```
1 accuracy(model=m, loader=test_loader, device="cuda")
```

```
## 0.8318
```

Report

Code

Results

```
1 from sklearn.metrics import classification_report
2
3 def report(model, loader, device):
4     y_true, y_pred = [], []
5     with torch.no_grad():
6         for X, y in loader:
7             X = X.to(device=device)
8             y_true.append( y.cpu().numpy() )
9             y_pred.append( model(X).max(1)[1].cpu().numpy() )
10
11     y_true = np.concatenate(y_true)
12     y_pred = np.concatenate(y_pred)
13
14     return classification_report(y_true, y_pred, target_names=loader.dataset.classes)
```

Transformers & GPT

Attribution & Resources

The following slides are a very brief and simplified look at a “modern” GPT model using Torch.

The code comes from Andrej Karpathy’s [minGPT](#) repo.

A PyTorch re-implementation of GPT, both training and inference. minGPT tries to be small, clean, interpretable and educational, as most of the currently available GPT model implementations can be a bit sprawling.

It is an impressive piece of work and well worth looking through. More recently, Karpathy put together [microgpt](#) which is a pure-Python follow-up with no dependencies (including full autograd).

I would also recommend growingswe.com/blog/microgpt — an excellent explainer and companion that walks through the microGPT code in detail.

Language Models

A language model assigns a probability to a sequence of tokens:

$$P(x_1, x_2, \dots, x_T) = \prod_{t=1}^T P(x_t \mid x_1, \dots, x_{t-1})$$

Training objective: given the previous tokens, predict the next one.

- Tokens can be words, sub-words, characters, ...
- At each step the model sees all *prior* context (but not future tokens)
- Loss is cross-entropy averaged over all positions

Transformer Architecture

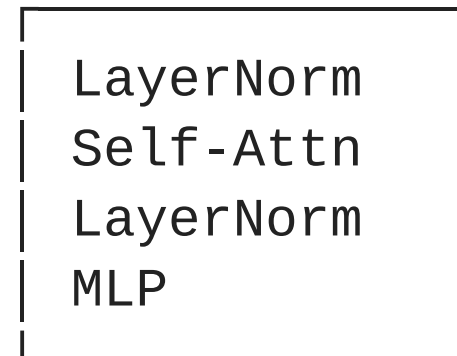
Introduced in *Attention Is All You Need* (Vaswani et al., 2017).

Key ideas:

- Replace recurrence with self-attention — every token attends to every other token in a single step
- Stack N identical transformer blocks
- Add positional encodings so the model knows token order
- Scale to billions of parameters via the same architecture

GPT (decoder-only) stack:

Token + Position Embeddings



× N blocks



LayerNorm



LM Head (linear)

Causal (Masked) Self-Attention

The core operation — each token produces a *query*, *key*, and *value*:

$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^\top}{\sqrt{d_k}} + M\right) V$$

where M is a causal mask (upper-triangle = $-\infty$) that prevents attending to future positions.

With h heads, each head learns a different projection:

$$\text{head}_i = \text{Attention}(QW_i^Q, KW_i^K, VW_i^V), \quad W_i^Q, W_i^K, W_i^V \in \mathbb{R}^{d_{\text{model}} \times d_k}$$

outputs are concatenated and projected back:

$$\text{MultiHead}(Q, K, V) = \text{Concat}(\text{head}_1, \dots, \text{head}_h) W^O$$

Implementation - CausalSelfAttention

```
1 class CausalSelfAttention(nn.Module):
2     def __init__(self, n_embd, n_head, block_size, dropout=0.1):
3         super().__init__()
4         assert n_embd % n_head == 0
5         self.n_head, self.n_embd = n_head, n_embd
6         self.c_attn = nn.Linear(n_embd, 3 * n_embd) # Q, K, V in one shot
7         self.c_proj = nn.Linear(n_embd, n_embd)
8         self.attn_drop = nn.Dropout(dropout)
9         self.resid_drop = nn.Dropout(dropout)
10        # lower-triangular causal mask (1, 1, T, T)
11        mask = torch.tril(torch.ones(block_size, block_size)).view(1, 1, block_size, block_size)
12        self.register_buffer("mask", mask)
13
14    def forward(self, x):
15        B, T, C = x.size()
16        hs = C // self.n_head
17        q, k, v = self.c_attn(x).split(self.n_embd, dim=2) # (B, T, C) each
18        # reshape to (B, n_head, T, head_size)
19        q = q.view(B, T, self.n_head, hs).transpose(1, 2)
20        k = k.view(B, T, self.n_head, hs).transpose(1, 2)
21        v = v.view(B, T, self.n_head, hs).transpose(1, 2)
22        # scaled dot-product attention with causal mask
```

Implementation - Block

Each transformer block wraps attention + a feed-forward MLP with residual connections and layer norm:

```
1 class Block(nn.Module):
2     def __init__(self, n_embd, n_head, block_size, dropout=0.1):
3         super().__init__()
4         self.ln1 = nn.LayerNorm(n_embd)
5         self.attn = CausalSelfAttention(n_embd, n_head, block_size, dropout)
6         self.ln2 = nn.LayerNorm(n_embd)
7         self.mlp = nn.Sequential(
8             nn.Linear(n_embd, 4 * n_embd),
9             nn.GELU(),
10            nn.Linear(4 * n_embd, n_embd),
11            nn.Dropout(dropout),
12        )
13
14    def forward(self, x):
15        x = x + self.attn(self.ln1(x)) # attention + residual
16        x = x + self.mlp(self.ln2(x)) # feed-forward + residual
17        return x
```

Implementation - GPT

```
1 class GPT(nn.Module):
2     def __init__(self, vocab_size, block_size, n_layer=6, n_head=6, n_embd=192, dropout=0.1):
3         super().__init__()
4         self.block_size = block_size
5         self.transformer = nn.ModuleDict(dict(
6             wte = nn.Embedding(vocab_size, n_embd),           # token embeddings
7             wpe = nn.Embedding(block_size, n_embd),           # position embeddings
8             drop = nn.Dropout(dropout),
9             h     = nn.ModuleList([Block(n_embd, n_head, block_size, dropout)
10                                     for _ in range(n_layer)]),
11             ln_f = nn.LayerNorm(n_embd),
12         ))
13         self.lm_head = nn.Linear(n_embd, vocab_size, bias=False)
14         self.apply(self._init_weights)
15
16     def _init_weights(self, m):
17         if isinstance(m, (nn.Linear, nn.Embedding)):
18             nn.init.normal_(m.weight, std=0.02)
19             if hasattr(m, 'bias') and m.bias is not None:
20                 nn.init.zeros_(m.bias)
21
22     def forward(self, idx, targets=None):
```

Model Sizes

Variant	Layers	Heads	Embed dim	~Params
gpt-nano	3	3	48	45 K
gpt-micro	4	4	128	660 K
gpt-mini	6	6	192	3 M
GPT-2 small	12	12	768	124 M
GPT-2 medium	24	16	1024	350 M
GPT-2 XL	48	25	1600	1.6 B

For our example use case, character-level name generation, we only need a tiny model — the vocabulary is ~27 characters and sequences are short.

Character-level Name Generation

The Names Dataset

This is a common classic benchmark for character-level language models — ~32k US baby names, one per line.

```
1 import urllib.request
2 url = "https://raw.githubusercontent.com/karpathy/makemore/master/names.txt"
3 names_txt = urllib.request.urlopen(url).read().decode()
4 names = names_txt.strip().split('\n')
5 len(names)
```

32033

```
1 names[:10]
```

```
['emma', 'olivia', 'ava', 'isabella', 'sophia', 'charlotte', 'mia', 'amelia', 'harper']
```

- Each name is terminated by `\n`
 - the model learns `\n` -> start of a name *and* name -> `\n`
- Character vocabulary - 26 letters + newline = 27 tokens
- Mean name length - 6.1 characters

CharDataset

Each training example is a window of `block_size` characters; the target is the same window shifted by one:

```
1 class CharDataset(Dataset):
2     def __init__(self, text, block_size):
3         chars = sorted(set(text))
4         self.stoi = {c: i for i, c in enumerate(chars)}
5         self.itos = {i: c for i, c in enumerate(chars)}
6         self.vocab_size = len(chars)
7         self.block_size = block_size
8         self.data = torch.tensor([self.stoi[c] for c in text], dtype=torch.long)
9
10    def __len__(self):
11        return len(self.data) - self.block_size
12
13    def __getitem__(self, idx):
14        chunk = self.data[idx : idx + self.block_size + 1]
15        return chunk[:-1], chunk[1:] # (x, y) shifted by 1
```

```
1 BLOCK_SIZE = 32
2 dataset = CharDataset(names_txt, block_size=BLOCK_SIZE)
```

```
1 x, y = dataset[0]
```

```
1 x
```

```
tensor([ 5, 13, 13,  1,  0, 15, 12,  9, 22,
         9,  1,  0,  1, 22,  1,  0,  9, 19,
         1,  2,  5, 12, 12,  1,  0, 19, 15,
        16,  8,  9,  1,  0])
```

```
1 y
```

```
tensor([13, 13,  1,  0, 15, 12,  9, 22,  9,
         1,  0,  1, 22,  1,  0,  9, 19,  1,
         2,  5, 12, 12,  1,  0, 19, 15, 16,
         8,  9,  1,  0,  3])
```

```
1 ''.join(dataset.itos[i.item()] for i in x)
```

```
'emma\nolivia\nava\nisabella\nsophia\n'
```

```
1 ''.join(dataset.itos[i.item()] for i in y)
```

```
'mma\nolivia\nava\nisabella\nsophia\nc'
```

```
1 x, y = dataset[1]
```

```
1 x
```

```
tensor([13, 13,  1,  0, 15, 12,  9, 22,  9,
         1,  0,  1, 22,  1,  0,  9, 19,  1,
         2,  5, 12, 12,  1,  0, 19, 15, 16,
         8,  9,  1,  0,  3])
```

```
1 y
```

```
tensor([13,  1,  0, 15, 12,  9, 22,  9,  1,
         0,  1, 22,  1,  0,  9, 19,  1,  2,
         5, 12, 12,  1,  0, 19, 15, 16,  8,
         9,  1,  0,  3,  8])
```

```
1 ''.join(dataset.itos[i.item()] for i in x)
```

```
'mma\nolivia\nava\nisabella\nsophia\nc'
```

```
1 ''.join(dataset.itos[i.item()] for i in y)
```

```
'ma\nolivia\nava\nisabella\nsophia\nch'
```

Training Setup

```
1 BLOCK_SIZE = 32
2 dataset = CharDataset(names_txt, block_size=BLOCK_SIZE)
3 loader = DataLoader(dataset, batch_size=256, shuffle=True)
4
5 model = GPT(
6     vocab_size = dataset.vocab_size,
7     block_size = BLOCK_SIZE,
8     n_layer = 4,
9     n_head = 4,
10    n_embd = 128,
11    dropout = 0.1,
12 ).to(device)
13
14 optimizer = torch.optim.AdamW(model.parameters(), lr=3e-3, weight_decay=0.1)
15 scheduler = torch.optim.lr_scheduler.CosineAnnealingLR(optimizer, T_max=5000)
```

Model parameter count:

```
1 sum(p.numel() for p in model.parameters() if p.requires_grad)
```

804352

`CosineAnnealingLR` decays the learning rate from `lr_max` to 0 following a cosine curve over `T_max` steps — starts fast,

Training

```
1 model.train()
```

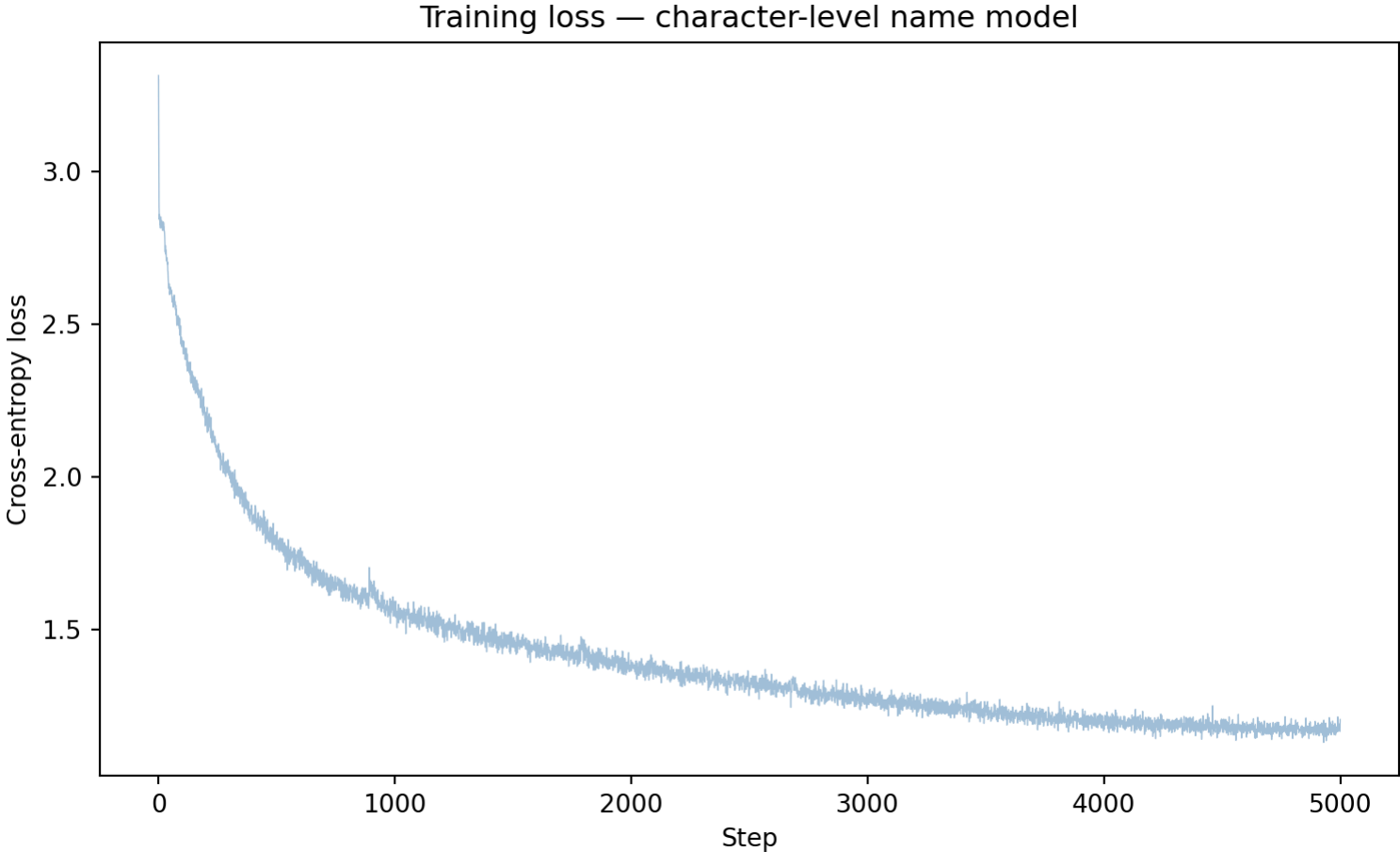
```
GPT(
  (transformer): ModuleDict(
    (wte): Embedding(27, 128)
    (wpe): Embedding(32, 128)
    (drop): Dropout(p=0.1, inplace=False)
    (h): ModuleList(
      (0-3): 4 x Block(
        (ln1): LayerNorm((128,)), eps=1e-05, elementwise_affine=True)
        (attn): CausalSelfAttention(
          (c_attn): Linear(in_features=128, out_features=384, bias=True)
          (c_proj): Linear(in_features=128, out_features=128, bias=True)
          (attn_drop): Dropout(p=0.1, inplace=False)
          (resid_drop): Dropout(p=0.1, inplace=False)
        )
        (ln2): LayerNorm((128,)), eps=1e-05, elementwise_affine=True)
        (mlp): Sequential(
          (0): Linear(in_features=128, out_features=512, bias=True)
          (1): GELU(approximate='none')
          (2): Linear(in_features=512, out_features=128, bias=True)
          (3): Dropout(p=0.1, inplace=False)
        )
      )
    )
)
```

```
1 losses = []
2 data_iter = iter(loader)
3
4 for step in range(5000):
5     try:
6         x, y = next(data_iter)
7     except StopIteration:
8         data_iter = iter(loader)
9         x, y = next(data_iter)
10
11     x, y = x.to(device), y.to(device)
12     logits, loss = model(x, y)
13
14     optimizer.zero_grad()
15     loss.backward()
16     torch.nn.utils.clip_grad_norm_(model.parameters(), 1.0)
17     optimizer.step()
18     scheduler.step()
19     losses.append(loss.item())
```

Training Loss

Actual

Smoothed



Generating Names

```
1 def sample_names(model, dataset, n=10, temperature=0.8, top_k=10):
2     newline_idx = dataset.stoi['\n']
3     start = torch.tensor([[newline_idx]], dtype=torch.long, device=device)
4     names_out = []
5     for _ in range(n):
6         out = model.generate(start, max_new_tokens=20, temperature=temperature, top_k=top_k)
7         tokens = out[0].tolist()
8         # collect characters up to the next newline (after the seed)
9         name = ''
10        for t in tokens[1:]:
11            c = dataset.itos[t]
12            if c == '\n':
13                break
14            name += c
15        if name:
16            names_out.append(name)
17    return names_out
18
19 sample_names(model, dataset)
```

['jahmari', 'kohler', 'montrell', 'edith', 'arleanne', 'derian', 'makala', 'josef', 'elianah', 'ka

Temperature & Sampling

temperature=0.5 (more conservative):

```
1 sample_names(model, dataset, n=10, temperature=0.5)
```

```
['kenish', 'abduLrahman', 'danaya', 'kingston', 'adalei', 'alijah', 'solomi', 'ameena', 'maryelynn
```

temperature=1.0 (default):

```
1 sample_names(model, dataset, n=10, temperature=1.0)
```

```
['asheal', 'lexander', 'caiman', 'scotland', 'adelia', 'milles', 'kyaira', 'loisie', 'darsha', 'se
```

temperature=1.5 (more creative / noisy):

```
1 sample_names(model, dataset, n=10, temperature=1.5)
```

```
['jakoby', 'story', 'reyn', 'sunair', 'milealah', 'dmayah', 'jiles', 'elazar', 'dollan', 'emmalyn'
```

Lower temperature sharpens the distribution (more likely completions); higher temperature flattens it (more diverse but

Some “state-of-the-art” models

Hugging Face

This is an online community and platform for sharing machine learning models (architectures and weights), data, and related artifacts. They also maintain a number of packages and related training materials that help with building, training, and deploying ML models.

Some notable resources,

- [transformers](#) - APIs and tools to easily download and train state-of-the-art (pretrained) transformer based models
- [diffusers](#) - provides pretrained vision and audio diffusion models, and serves as a modular toolbox for inference and training

Stable Diffusion

```
1 from huggingface_hub import snapshot_download, login
```

```
1 model_id = "/data/stable-diffusion-2-1"
```

```
1 snapshot_download(  
2     repo_id="sd2-community/stable-diffusion-2-1",  
3     local_dir=model_id,  
4     ignore_patterns=["*.ckpt", "*.safetensors.index.json"],  
5     token="..." # or use login() ahead of time  
6 )
```

```
1 import torch  
2 from diffusers import StableDiffusionPipeline, DPMSolverMultistepScheduler  
3  
4 pipe = StableDiffusionPipeline.from_pretrained(model_id, torch_dtype=torch.float16)
```

```
1 pipe.scheduler = DPMSolverMultistepScheduler.from_config(pipe.scheduler.config)  
2 pipe = pipe.to("cuda")
```

Inference

Code Images

```
1 prompt = "a picture of thomas bayes with a cat on his lap"  
2 generator = [torch.Generator(device="cuda").manual_seed(i) for i in range(6)]  
3 fit = pipe(prompt, generator=generator, num_inference_steps=20, num_images_per_prompt=6)
```

```
1 fit.images
```

[<PIL.Image.Image image mode=RGB size=768x768 at 0x7F2CEDB5FA10>, <PIL.Image.Image image mode=RGB

```
1 fig, axes = plt.subplots(nrows=2, ncols=3, figsize=(10, 6), layout="constrained")  
2  
3 for i, ax in enumerate([ax for row in axes for ax in row]):  
4     ax.set_axis_off()  
5     p = ax.imshow(fit.images[i])  
6  
7 plt.show()
```

Customizing prompts

Code Images

```
1 prompt = "a picture of thomas bayes with a cat on his lap"
2 prompts = [
3     prompt + t for t in
4     ["in the style of a japanese wood block print",
5      "as a hipster with facial hair and glasses",
6      "as a simpsons character, cartoon, yellow",
7      "in the style of a vincent van gogh painting",
8      "in the style of a picasso painting",
9      "with flowery wall paper"
10    ]
11 ]
```

```
1 generator = [torch.Generator(device="cuda").manual_seed(i) for i in range(6)]
2 fit = pipe(prompts, generator=generator, num_inference_steps=20, num_images_per_prompt=1)
```

Increasing inference steps

- 1 generator = [torch.Generator(device="cuda").manual_seed(i) for i in range(6)]
- 2 fit = pipe(prompts, generator=generator, num_inference_steps=50, num_images_per_prompt=1)



A more current model

This model is larger than the available GPU memory - so we adjust the weight types to make it fit.

```
1 from diffusers import BitsAndBytesConfig, SD3Transformer2DModel
2 from diffusers import StableDiffusion3Pipeline
3 import torch
4
5 model_id = "/data/stable-diffusion-3-5-medium"
6
7 nf4_config = BitsAndBytesConfig(
8     load_in_4bit=True,
9     bnb_4bit_quant_type="nf4",
10    bnb_4bit_compute_dtype=torch.bfloat16
11 )
12 model_nf4 = SD3Transformer2DModel.from_pretrained(
13     model_id,
14     subfolder="transformer",
15     quantization_config=nf4_config,
16     torch_dtype=torch.bfloat16
17 )
18 pipe = StableDiffusion3Pipeline.from_pretrained(
19     model_id,
20     transformer=model_nf4,
21     torch_dtype=torch.bfloat16
22 )
```

```
1 pipe.enable_model_cpu_offload()
```

Images

```
1 generator = [torch.Generator(device="cuda").manual_seed(i) for i in range(6)]  
2 fit = pipe(prompts, generator=generator, num_inference_steps=30, num_images_per_prompt=1)
```

