

pytorch - GPU

Lecture 24

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.

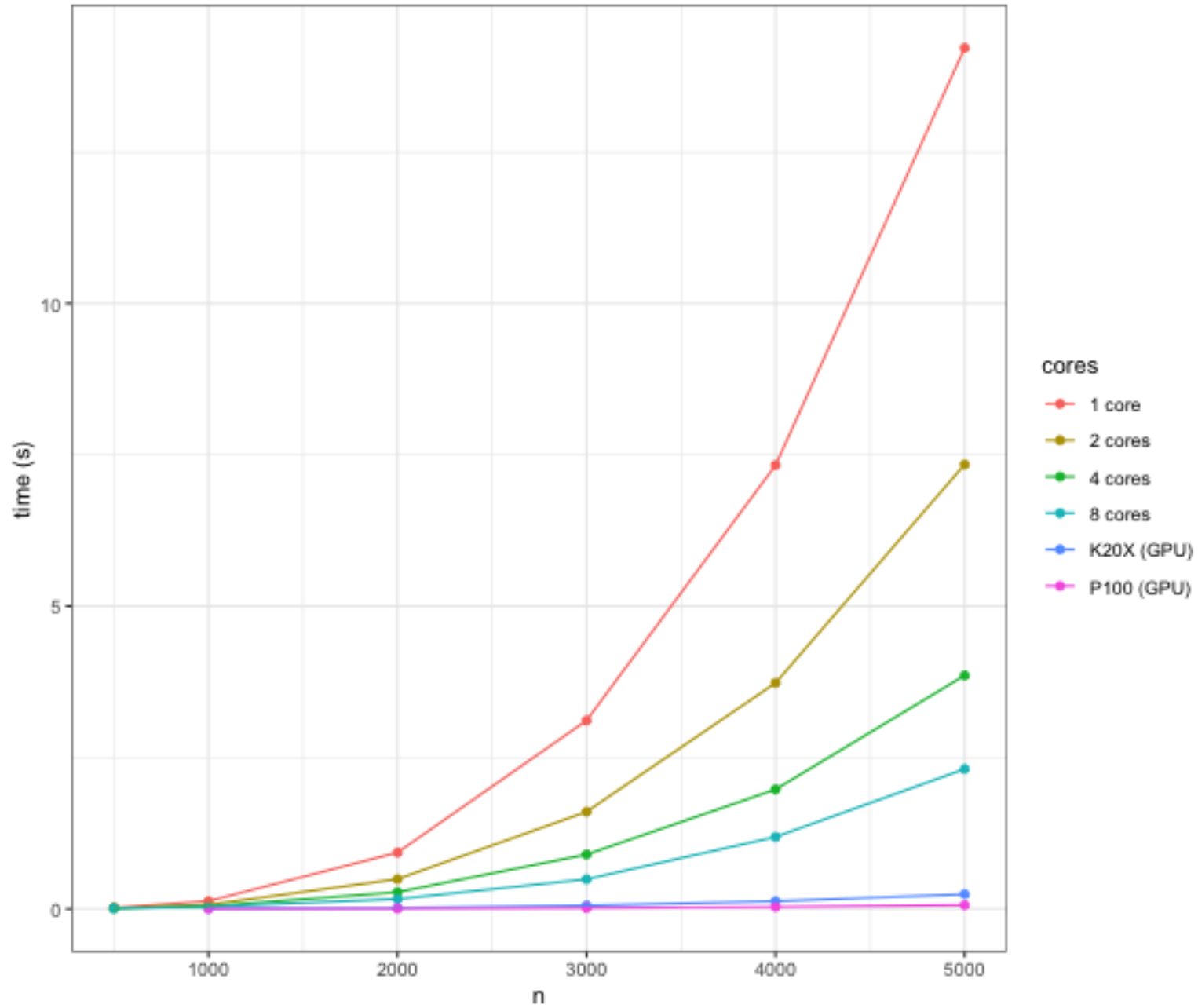
Core libraries:

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

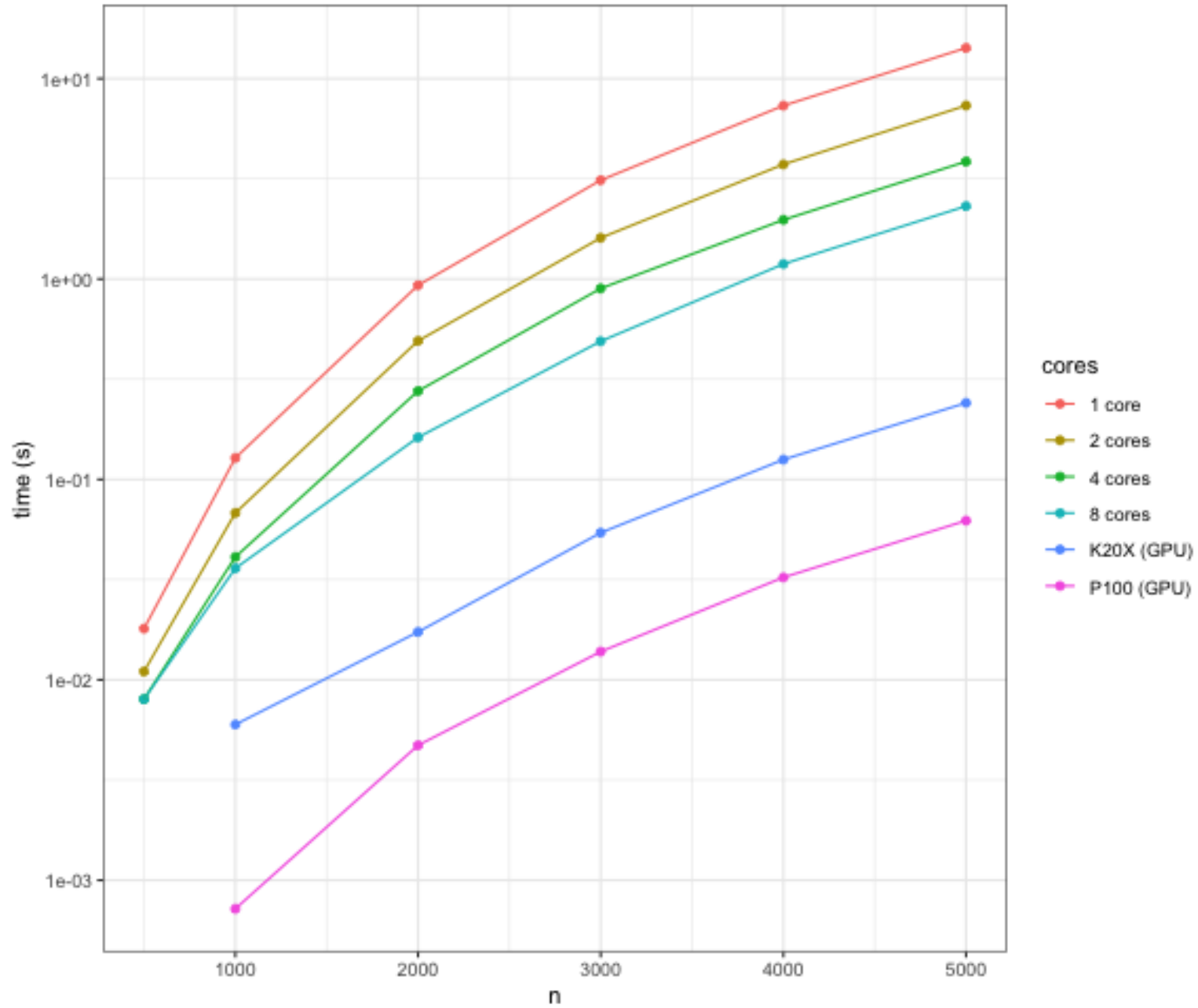
CUDA Kernels

```
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 }
```

Matrix Multiply of (n x n) matrices - double precision



Matrix Multiply of (n x n) matrices - double precision



GPU Status

```
1 nvidia-smi
```

Wed Apr 12 10:32:48 2023

```
+-----+
| NVIDIA-SMI 530.30.02                Driver Version: 530.30.02   CUDA Version: 12.1   |
+-----+-----+-----+-----+-----+-----+
| 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   Tesla P100-PCIE-16GB             Off| 00000000:02:00.0 Off |             0        |
| N/A   39C   P0              31W / 250W|  1002MiB / 16384MiB |      0%      Default |
|                                           |              N/A    |
+-----+-----+-----+-----+-----+-----+
|   1   Tesla P100-PCIE-16GB             Off| 00000000:03:00.0 Off |             0        |
| N/A   39C   P0              27W / 250W|    2MiB / 16384MiB |      0%      Default |
|                                           |              N/A    |
+-----+-----+-----+-----+-----+-----+

+-----+
| Processes:                                |
| GPU   GI    CI          PID    Type    Process name                        GPU Memory |
|       ID    ID                                   |             Usage              |
+=====+=====+=====+=====+=====+=====+
|     0   N/A  N/A     1825870      C   /usr/lib/rstudio-server/bin/rsession  1000MiB |
```

Torch GPU Information

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

True

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

2

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

'Tesla P100-PCIE-16GB'

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

'Tesla P100-PCIE-16GB'

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

_CudaDeviceProperties(name='Tesla P100-PCIE-16GB', major=6, minor=0, total_memory=16276MB, multi_processor_c

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

_CudaDeviceProperties(name='Tesla P100-PCIE-16GB', major=6, minor=0, total_memory=16276MB, multi_processor_c

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')
4
5 x = torch.linspace(0,1,5, device=cuda0); x
```

```
tensor([0.0000, 0.2500, 0.5000, 0.7500, 1.0000], dev:
```

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

```
tensor([[ 0.6879, -2.3114],
        [ 0.5199,  0.5865],
        [-0.5277,  0.9261],
        [-0.4613, -0.7858],
        [-1.8057,  1.2171]], device='cuda:0')
```

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

```
tensor([[0.8585, 0.9918, 0.5125],
        [0.3261, 0.3992, 0.7753]])
```

```
1 x @ y
```

```
tensor([-2.2856,  1.2374], device='cuda:0')
```

```
1 y @ z
```

```
Error: RuntimeError: Expected all tensors to be on the same device, but found at least two devices, cuda:0 and cpu!
```

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

```
tensor([[[-0.1631, -0.2403, -1.4393],
          [ 0.6375,  0.7497,  0.7211],
          [-0.1511, -0.1537,  0.4475],
          [-0.6523, -0.7712, -0.8456],
          [-1.1534, -1.3051,  0.0182]], 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)
```

Error: RuntimeError: Expected all tensors to be on the same device, but found at least two devices, cpu and

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

```
tensor([[ 0.0588, -0.4098, -0.6487,  0.4108,  0.0322]
        [-0.0797, -0.4688, -1.7741,  0.2954, -0.0537]
        [-0.2057, -0.4174, -1.3236, -0.0518, -0.4419]
        [ 0.7228, -0.2675,  1.3836, -0.4886,  0.3568]
        [-0.9248,  0.1247, -0.1394, -0.5943,  0.1210]
        [-1.0062, -0.1228, -0.5563,  0.8400,  0.8503]
        [-0.6758, -0.1711,  0.6801,  0.5007,  0.7845]
        [-0.1415, -0.0750, -0.1112,  0.2573,  0.8202]
        [-0.6765,  0.0109,  1.0709, -0.6299, -0.1831]
        [-1.4872, -0.0788,  0.9389,  0.2557,  0.1959]
        grad_fn=<AddmmBackward0>)
```

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

```
tensor([[ 0.0588, -0.4098, -0.6487,  0.4108,  0.0322]
        [-0.0797, -0.4688, -1.7741,  0.2954, -0.0537]
        [-0.2057, -0.4174, -1.3236, -0.0518, -0.4419]
        [ 0.7228, -0.2675,  1.3836, -0.4886,  0.3568]
        [-0.9248,  0.1247, -0.1394, -0.5943,  0.1210]
        [-1.0062, -0.1228, -0.5563,  0.8400,  0.8503]
        [-0.6758, -0.1711,  0.6801,  0.5007,  0.7845]
        [-0.1415, -0.0750, -0.1112,  0.2573,  0.8202]
        [-0.6765,  0.0109,  1.0709, -0.6299, -0.1831]
        [-1.4872, -0.0788,  0.9389,  0.2557,  0.1959]
        grad_fn=<AddmmBackward0>)
```

Back to MNIST

Same MNIST data from last time (1x8x8 images),

```
1 from sklearn.datasets import load_digits
2 from sklearn.model_selection import train_test_split
3
4 digits = load_digits()
5 X, y = digits.data, digits.target
6
7 X_train, X_test, y_train, y_test = train_test_split(
8     X, y, test_size=0.20, shuffle=True, random_state=1234
9 )
10
11 X_train = torch.from_numpy(X_train).float()
12 y_train = torch.from_numpy(y_train)
13 X_test = torch.from_numpy(X_test).float()
14 y_test = torch.from_numpy(y_test)
```

To use the GPU for computation we need to copy these tensors to the GPU,

```
1 X_train_cuda = X_train.to(device=cuda0)
2 y_train_cuda = y_train.to(device=cuda0)
3 X_test_cuda = X_test.to(device=cuda0)
4 y_test_cuda = y_test.to(device=cuda0)
```

Convolutional NN

```
1 class mnist_conv_model(torch.nn.Module):
2     def __init__(self, device):
3         super().__init__()
4         self.device = torch.device(device)
5
6         self.model = torch.nn.Sequential(
7             torch.nn.Unflatten(1, (1,8,8)),
8             torch.nn.Conv2d(
9                 in_channels=1, out_channels=8,
10                kernel_size=3, stride=1, padding=1
11            ),
12            torch.nn.ReLU(),
13            torch.nn.MaxPool2d(kernel_size=2),
14            torch.nn.Flatten(),
15            torch.nn.Linear(8 * 4 * 4, 10)
16        ).to(device=self.device)
17
18    def forward(self, X):
19        return self.model(X)
20
21    def fit(self, X, y, lr=0.001, n=1000, acc_step=10):
22        opt = torch.optim.SGD(self.parameters(), lr=lr, momentum=0.9)
23        losses = []
```

CPU vs Cuda

```
1 m = mnist_conv_model(device="cpu")
2 loss = m.fit(X_train, y_train, n=1000)
3 loss[-1]
```

0.034776557236909866

```
1 m.accuracy(X_test, y_test)
```

tensor(0.9750)

```
1 m_cuda = mnist_conv_model(device="cuda")
2 loss = m_cuda.fit(X_train_cuda, y_train_cuda, n=
3 loss[-1]
```

0.03830884024500847

```
1 m_cuda.accuracy(X_test_cuda, y_test_cuda)
```

tensor(0.9750, device='cuda:0')

Performance

CPU performance:

```
1 m = mnist_conv_model(device="cpu")
2
3 start = torch.cuda.Event(enable_timing=True)
4 end = torch.cuda.Event(enable_timing=True)
5
6 start.record()
7 loss = m.fit(X_train, y_train, n=1000)
8 end.record()
9
10 torch.cuda.synchronize()
11 print(start.elapsed_time(end) / 1000)
```

2.75747021484375

GPU performance:

```
1 m_cuda = mnist_conv_model(device="cuda")
2
3 start = torch.cuda.Event(enable_timing=True)
4 end = torch.cuda.Event(enable_timing=True)
5
6 start.record()
7 loss = m_cuda.fit(X_train_cuda, y_train_cuda, n=
8 end.record()
9
10 torch.cuda.synchronize()
11 print(start.elapsed_time(end) / 1000)
```

2.358865234375

Profiling CPU - 1 forward step

```
1 m = mnist_conv_model(device="cpu")
2 with torch.autograd.profiler.profile(with_stack=True, profile_memory=True) as prof_cpu:
3     tmp = m(X_train)
```

```
1 print(prof_cpu.key_averages().table(sort_by='self_cpu_time_total', row_limit=5))
```

Name	Self CPU %	Self CPU	CPU total %	CPU total	CPU time avg
aten::mkldnn_convolution	53.28%	1.576ms	54.06%	1.599ms	1.599ms
aten::max_pool2d_with_indices	25.90%	766.000us	25.90%	766.000us	766.000us
aten::clamp_min	11.76%	348.000us	11.76%	348.000us	348.000us
aten::addmm	3.11%	92.000us	3.99%	118.000us	118.000us
aten::copy_	0.74%	22.000us	0.74%	22.000us	22.000us

Self CPU time total: 2.958ms

Profiling GPU - 1 forward step

```
1 m_cuda = mnist_conv_model(device="cuda")
2 with torch.autograd.profiler.profile(with_stack=True) as prof_cuda:
3     tmp = m_cuda(X_train_cuda)
```

```
1 print(prof_cuda.key_averages().table(sort_by='self_cpu_time_total', row_limit=5))
```

```
-----
              Name      Self CPU %      Self CPU   CPU total %   CPU total C
-----
          aten::conv2d      47.36%       1.902ms      47.36%       1.902ms
      aten::cudnn_convolution  44.82%       1.800ms      45.44%       1.825ms
          aten::addmm       1.84%       74.000us     2.39%       96.000us
      cudaLaunchKernel       1.52%       61.000us     1.52%       61.000us
      aten::max_pool2d_with_indices  0.60%       24.000us     0.80%       32.000us
-----
Self CPU time total: 4.016ms
```

Profiling CPU - fit

```
1 m = mnist_conv_model(device="cpu")
2 with torch.autograd.profiler.profile(with_stack=True, profile_memory=True) as prof_cpu:
3     losses = m.fit(X_train, y_train, n=1000)
```

```
1 print(prof_cpu.key_averages().table(sort_by='self_cpu_time_total', row_limit=5))
```

```
-----
```

Name	Self CPU %	Self CPU	CPU total %	CPU tc
aten::convolution_backward	28.11%	847.371ms	28.31%	853.38
aten::max_pool2d_with_indices	19.47%	586.731ms	19.47%	586.73
aten::mkldnn_convolution	10.07%	303.619ms	10.30%	310.46
aten::threshold_backward	6.98%	210.441ms	6.98%	210.44
aten::mm	6.34%	191.005ms	6.34%	191.00

```
-----
Self CPU time total: 3.014s
```


Profiling GPU - fit

```
1 m_cuda = mnist_conv_model(device="cuda")
2 with torch.autograd.profiler.profile(with_stack=True) as prof_cuda:
3     losses = m_cuda.fit(X_train_cuda, y_train_cuda, n=1000)
```

```
1 print(prof_cuda.key_averages().table(sort_by='self_cpu_time_total', row_limit=5))
```

```
-----
```

Name	Self CPU %	Self CPU	CPU total %	CPU tc
cudaLaunchKernel	13.89%	374.915ms	13.89%	374.91
Optimizer.step#SGD.step	13.57%	366.201ms	19.16%	517.32
aten::addmm	4.54%	122.508ms	6.26%	169.08
aten::convolution_backward	4.10%	110.628ms	8.23%	222.15
aten::cudnn_convolution	3.89%	105.139ms	5.15%	139.06

```
-----
```

Self CPU time total: 2.700s

CIFAR10

homepage

Loading the data

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

```
1 test_data = torchvision.datasets.CIFAR10(
2     root="/data",
3     train=False,
4     download=True,
5     transform=torchvision.transforms.ToTensor()
6 )
```

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, ..., 0.6196, 0.5961, 0.5804],  
          [0.0627, 0.0000, 0.0706, ..., 0.4824, 0.4667, 0.4784],  
          [0.0980, 0.0627, 0.1922, ..., 0.4627, 0.4706, 0.4275],  
          ...,  
          [0.8157, 0.7882, 0.7765, ..., 0.6275, 0.2196, 0.2078],  
          [0.7059, 0.6784, 0.7294, ..., 0.7216, 0.3804, 0.3255],  
          [0.6941, 0.6588, 0.7020, ..., 0.8471, 0.5922, 0.4824]]],  
        [[0.2431, 0.1804, 0.1882, ..., 0.5176, 0.4902, 0.4863],  
          [0.0784, 0.0000, 0.0314, ..., 0.3451, 0.3255, 0.3412],  
          [0.0941, 0.0275, 0.1059, ..., 0.3294, 0.3294, 0.2863],  
          ...,  
          [0.6667, 0.6000, 0.6314, ..., 0.5216, 0.1216, 0.1333],  
          [0.5451, 0.4824, 0.5647, ..., 0.5804, 0.2431, 0.2078]]],  
        ...],  
      dtype=torch.FloatTensor))
```

```
[0.5647, 0.5059, 0.5569, ..., 0.7216, 0.4627, 0.3608]],  
  
[[0.2471, 0.1765, 0.1686, ..., 0.4235, 0.4000, 0.4039],  
 [0.0784, 0.0000, 0.0000, ..., 0.2157, 0.1961, 0.2235],  
 [0.0824, 0.0000, 0.0314, ..., 0.1961, 0.1961, 0.1647],  
 ...,  
 [0.3765, 0.1333, 0.1020, ..., 0.2745, 0.0275, 0.0784],  
 [0.3765, 0.1647, 0.1176, ..., 0.3686, 0.1333, 0.1333],
```

Example data

frog



truck



truck



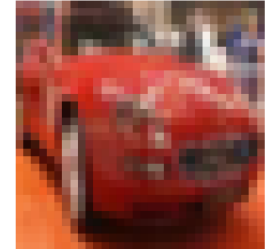
deer



automobile



automobile



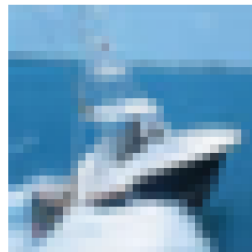
bird



horse



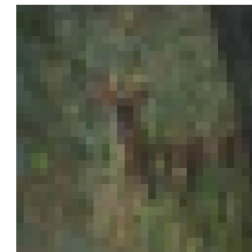
ship



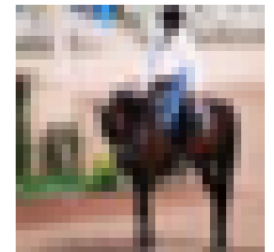
cat



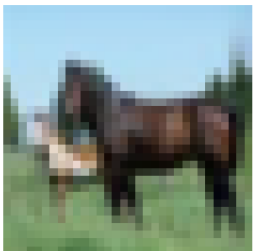
deer



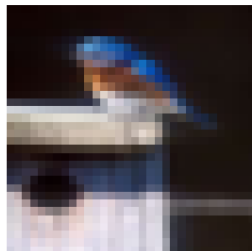
horse



horse



bird



truck



truck



truck



cat



bird



frog



deer



cat



frog



frog



Data Loaders

```
1  batch_size = 100
2
3  training_loader = torch.utils.data.DataLoader(
4      training_data,
5      batch_size=batch_size,
6      shuffle=True,
7      num_workers=4,
8      pin_memory=True
9  )
10
11 test_loader = torch.utils.data.DataLoader(
12     test_data,
13     batch_size=batch_size,
14     shuffle=True,
15     num_workers=4,
16     pin_memory=True
17 )
```

Loader generator

```
1 training_loader
```

```
<torch.utils.data.dataloader.DataLoader object at 0x7f42e510d210>
```

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

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

```
1 y.shape
```

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


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)
23
```

CNN Performance - CPU (1 step)

```
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=5))
```

```
-----
              Name      Self CPU %      Self CPU   CPU total %   CPU total   CPU time avg   #
-----
      aten::mkldnn_convolution    86.35%    12.410ms    86.57%    12.441ms     6.221ms
      aten::max_pool2d_with_indices    7.72%     1.109ms     7.72%     1.109ms    554.500us
          aten::clamp_min    2.67%    383.000us    2.67%    383.000us    95.750us
            aten::addmm    1.52%    219.000us    1.79%    257.000us    85.667us
              aten::relu    0.28%     40.000us    2.94%    423.000us    105.750us
-----
Self CPU time total: 14.371ms
```

CNN Performance - GPU (1 step)

```
1 m_cuda = cifar_conv_model(device="cuda")
2 Xc, yc = X.to(device="cuda"), y.to(device="cuda")
3 tmp = m_cuda(Xc)
4
5 with torch.autograd.profiler.profile(with_stack=True) as prof_cuda:
6     tmp = m_cuda(Xc)
```

```
1 print(prof_cuda.key_averages().table(sort_by='self_cpu_time_total', row_limit=5))
```

```
-----
              Name      Self CPU %      Self CPU   CPU total %   CPU total C
-----
      aten::cudnn_convolution    50.86%      3.451ms    51.70%      3.508ms
    aten::max_pool2d_with_indices    23.57%      1.599ms    23.79%      1.614ms
      aten::clamp_min           7.46%      506.000us    11.23%      762.000us
      aten::relu                6.88%      467.000us    18.11%      1.229ms
      cudaMalloc                3.17%      215.000us     3.17%      215.000us
-----
Self CPU time total: 6.785ms
```

CNN Performance - CPU (1 epoch)

```
1 m_cpu = cifar_conv_model(device="cpu")
2
3 with torch.autograd.profiler.profile(with_stack=True) as prof_cpu:
4     m_cpu.fit(loader=training_loader, epochs=1, n_report=501)
```

```
1 print(prof_cpu.key_averages().table(sort_by='self_cpu_time_total', row_limit=5))
```

```
-----
Name                Self CPU %      Self CPU      CPU total %      CPU tc
-----
aten::convolution_backward    30.65%         1.361s         31.40%           1.3
aten::mkldnn_convolution     17.39%         772.124ms      17.77%           788.80
aten::max_pool2d_with_indices  7.89%          350.091ms      7.90%           350.80
aten::threshold_backward      6.84%          303.662ms      6.85%           304.00
enumerate(DataLoader)#_MultiProcessingDataLoaderIter...  6.31%          280.147ms      6.33%           280.93
-----
Self CPU time total: 4.439s
```

CNN Performance - GPU (1 epoch)

```
1 m_cuda = cifar_conv_model(device="cuda")
2
3 with torch.autograd.profiler.profile(with_stack=True) as prof_cuda:
4     m_cuda.fit(loader=training_loader, epochs=1, n_report=501)
```

```
1 print(prof_cuda.key_averages().table(sort_by='self_cpu_time_total', row_limit=5))
```

```
-----
Name                Self CPU %      Self CPU      CPU total %    CPU tc
-----
enumerate(DataLoader)#_MultiProcessingDataLoaderIter...  12.04%         360.143ms     12.09%         361.64
      cudaLaunchKernel                11.75%         351.313ms     11.76%         351.52
      Optimizer.step#SGD.step          7.75%         231.842ms     10.60%         316.88
      aten::convolution_backward        5.73%         171.197ms     10.23%         305.92
      aten::addmm                        4.62%         138.160ms     6.21%          185.66
-----
Self CPU time total: 2.990s
```

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")
2 m.fit(training_loader, epochs=10, n_report=500, lr=0.01)
3 ## [Epoch 1, Minibatch 500] loss: 2.098
4 ## [Epoch 2, Minibatch 500] loss: 1.692
5 ## [Epoch 3, Minibatch 500] loss: 1.482
6 ## [Epoch 4, Minibatch 500] loss: 1.374
7 ## [Epoch 5, Minibatch 500] loss: 1.292
8 ## [Epoch 6, Minibatch 500] loss: 1.226
9 ## [Epoch 7, Minibatch 500] loss: 1.173
10 ## [Epoch 8, Minibatch 500] loss: 1.117
11 ## [Epoch 9, Minibatch 500] loss: 1.071
12 ## [Epoch 10, Minibatch 500] loss: 1.035
```

```
1 accuracy(m, training_loader, "cuda")
2 ## 0.63444
3 accuracy(m, test_loader, "cuda")
4 ## 0.572
```

More epochs

If continue fitting with the existing model,

```
1 m.fit(training_loader, epochs=10, n_report=500)
2 ## [Epoch 11, Minibatch 500] loss: 0.885
3 ## [Epoch 12, Minibatch 500] loss: 0.853
4 ## [Epoch 13, Minibatch 500] loss: 0.839
5 ## [Epoch 14, Minibatch 500] loss: 0.828
6 ## [Epoch 15, Minibatch 500] loss: 0.817
7 ## [Epoch 16, Minibatch 500] loss: 0.806
8 ## [Epoch 17, Minibatch 500] loss: 0.798
9 ## [Epoch 18, Minibatch 500] loss: 0.787
10 ## [Epoch 19, Minibatch 500] loss: 0.780
11 ## [Epoch 20, Minibatch 500] loss: 0.773
```

```
1 accuracy(m, training_loader, "cuda")
2 ## 0.73914
3 accuracy(m, test_loader, "cuda")
4 ## 0.624
```


More epochs (again)

```
1 m.fit(training_loader, epochs=10, n_report=500)
2 ## [Epoch 21, Minibatch 500] loss: 0.764
3 ## [Epoch 22, Minibatch 500] loss: 0.756
4 ## [Epoch 23, Minibatch 500] loss: 0.748
5 ## [Epoch 24, Minibatch 500] loss: 0.739
6 ## [Epoch 25, Minibatch 500] loss: 0.733
7 ## [Epoch 26, Minibatch 500] loss: 0.726
8 ## [Epoch 27, Minibatch 500] loss: 0.718
9 ## [Epoch 28, Minibatch 500] loss: 0.710
10 ## [Epoch 29, Minibatch 500] loss: 0.702
11 ## [Epoch 30, Minibatch 500] loss: 0.698
```

```
1 accuracy(m, training_loader, "cuda")
2 ## 0.76438
3 accuracy(m, test_loader, "cuda")
4 ## 0.6217
```

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__()
```

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))
 (21): BatchNorm2d(256, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)

VGG16 performance - CPU

```
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=5))
```

```
-----
      Name      Self CPU %      Self CPU  CPU total %  CPU total  CPU time avg  #
-----
  aten::mkldnn_convolution    81.00%    193.414ms    82.43%    196.826ms    15.140ms
  aten::native_batch_norm     8.97%     21.413ms     9.63%     23.005ms     1.770ms
  aten::max_pool2d_with_indices  6.30%     15.053ms     6.30%     15.053ms     3.011ms
  aten::empty                  2.04%      4.881ms     2.04%      4.881ms     37.546us
  aten::clamp_min_            1.06%      2.530ms     1.06%      2.530ms     194.615us
-----
Self CPU time total: 238.793ms
```

VGG16 performance - GPU

```
1 m_cuda = VGG16(device="cuda")
2 Xc, yc = X.to(device="cuda"), y.to(device="cuda")
3 tmp = m_cuda(Xc)
4
5 with torch.autograd.profiler.profile(with_stack=True) as prof_cuda:
6     tmp = m_cuda(Xc)
```

```
1 print(prof_cuda.key_averages().table(sort_by='self_cpu_time_total', row_limit=5))
```

Name	Self CPU %	Self CPU	CPU total %	CPU total
aten::cudnn_convolution	38.88%	3.379ms	48.86%	4.246ms
aten::cudnn_batch_norm	24.98%	2.171ms	33.44%	2.906ms
cudaMalloc	11.40%	991.000us	11.40%	991.000us
cudaLaunchKernel	6.40%	556.000us	6.40%	556.000us
aten::add_	3.43%	298.000us	5.05%	439.000us

Self CPU time total: 8.691ms

VGG16 performance - Apple M1 GPU (mps)

```
1 m_mps = VGG16(device="mps")
2 Xm, ym = X.to(device="mps"), y.to(device="mps")
3
4 with torch.autograd.profiler.profile(with_stack=True) as prof_mps:
5     tmp = m_mps(Xm)
```

```
1 print(prof_mps.key_averages().table(sort_by='self_cpu_time_total', row_limit=5))
```

Name	Self CPU %	Self CPU	CPU total %	CPU total	CPU time avg	#
aten::native_batch_norm	35.71%	3.045ms	35.71%	3.045ms	234.231us	
aten::_mps_convolution	19.67%	1.677ms	19.88%	1.695ms	130.385us	
aten::_batch_norm_impl_index	11.92%	1.016ms	36.02%	3.071ms	236.231us	
aten::relu_	11.29%	963.000us	11.29%	963.000us	74.077us	
aten::add_	10.40%	887.000us	10.44%	890.000us	68.462us	

Self CPU time total: 8.526ms

Fitting w/ $lr = 0.01$

```
1 m = VGG16(device="cuda")
2 fit(m, training_loader, epochs=10, n_report=500, lr=0.01)
3
4 ## [Epoch 1, Minibatch 500] loss: 1.345
5 ## [Epoch 2, Minibatch 500] loss: 0.790
6 ## [Epoch 3, Minibatch 500] loss: 0.577
7 ## [Epoch 4, Minibatch 500] loss: 0.445
8 ## [Epoch 5, Minibatch 500] loss: 0.350
9 ## [Epoch 6, Minibatch 500] loss: 0.274
10 ## [Epoch 7, Minibatch 500] loss: 0.215
11 ## [Epoch 8, Minibatch 500] loss: 0.167
12 ## [Epoch 9, Minibatch 500] loss: 0.127
13 ## [Epoch 10, Minibatch 500] loss: 0.103
```

```
1 accuracy(model=m, loader=training_loader, device="cuda")
2 ## 0.97008
3 accuracy(model=m, loader=test_loader, device="cuda")
4 ## 0.8318
```

Fitting w/ $lr = 0.001$

```
1 m = VGG16(device="cuda")
2 fit(m, training_loader, epochs=10, n_report=500, lr=0.001)
3
4 ## [Epoch 1, Minibatch 500] loss: 1.279
5 ## [Epoch 2, Minibatch 500] loss: 0.827
6 ## [Epoch 3, Minibatch 500] loss: 0.599
7 ## [Epoch 4, Minibatch 500] loss: 0.428
8 ## [Epoch 5, Minibatch 500] loss: 0.303
9 ## [Epoch 6, Minibatch 500] loss: 0.210
10 ## [Epoch 7, Minibatch 500] loss: 0.144
11 ## [Epoch 8, Minibatch 500] loss: 0.108
12 ## [Epoch 9, Minibatch 500] loss: 0.088
13 ## [Epoch 10, Minibatch 500] loss: 0.063
```

```
1 accuracy(model=m, loader=training_loader, device="cuda")
2 ## 0.9815
3 accuracy(model=m, loader=test_loader, device="cuda")
4 ## 0.7816
```


Report

```
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)
```

```
1 print(report(model=m, loader=test_loader, device="cuda"))
2
3 ##           precision    recall  f1-score   support
4 ##
5 ##   airplane      0.82     0.88     0.85     1000
6 ##  automobile      0.95     0.89     0.92     1000
7 ##         bird      0.85     0.70     0.77     1000
8 ##         cat      0.68     0.74     0.71     1000
9 ##         deer      0.84     0.83     0.83     1000
10 ##        dog      0.81     0.73     0.77     1000
11 ##        frog      0.83     0.92     0.87     1000
12 ##       horse      0.87     0.87     0.87     1000
13 ##        ship      0.89     0.92     0.90     1000
14 ##       truck      0.86     0.93     0.89     1000
15 ##
16 ##      accuracy              0.84     10000
17 ##    macro avg      0.84     0.84     0.84     10000
18 ##  weighted avg      0.84     0.84     0.84     10000
```

Some state-of-the-art examples

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
- [timm](#) - a library containing SOTA computer vision models, layers, utilities, optimizers, schedulers, data-loaders, augmentations, and training/evaluation scripts

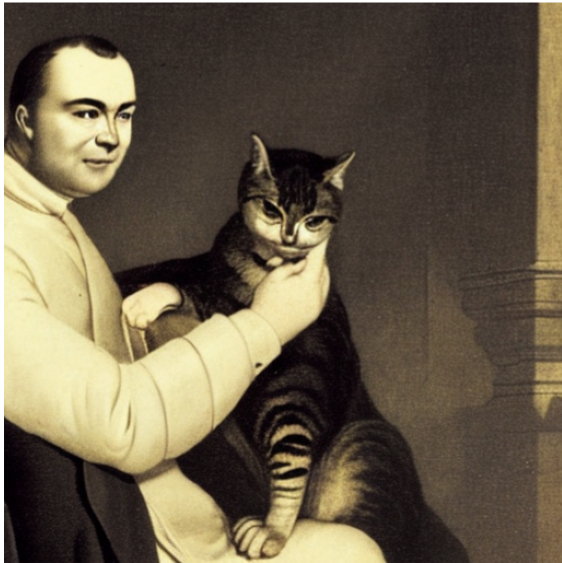
Stable Diffusion

```
1 from diffusers import StableDiffusionPipeline
2
3 pipe = StableDiffusionPipeline.from_pretrained(
4     "stabilityai/stable-diffusion-2-1-base", torch_dtype=torch.float16
5 ).to("cuda")
```

```
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)
```

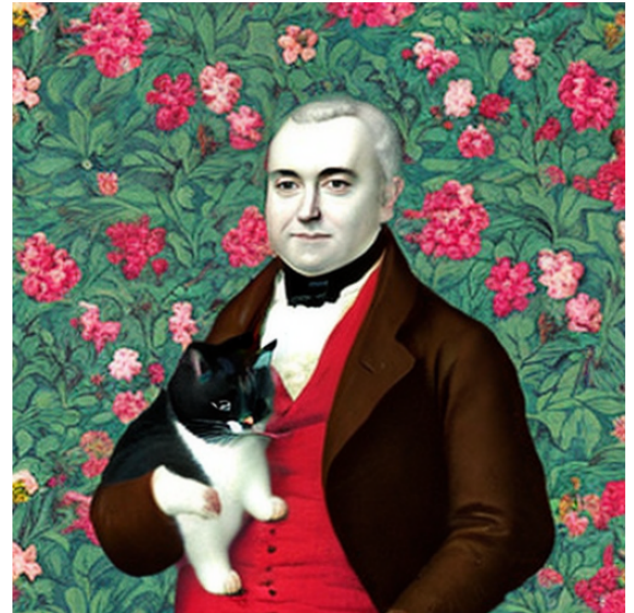
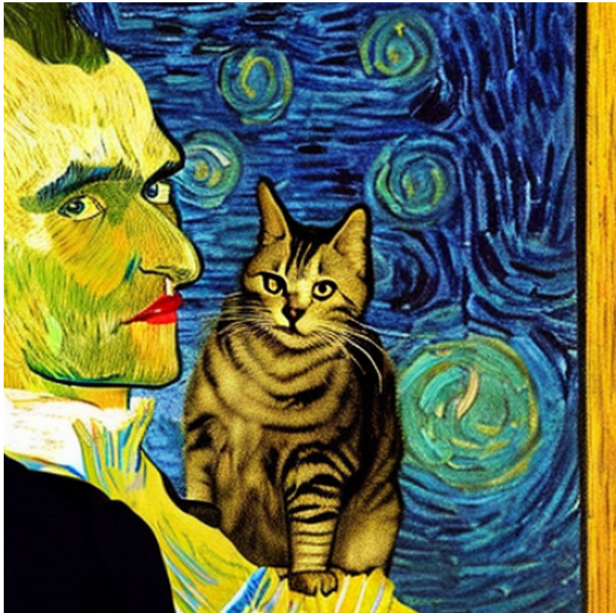
```
1 fit.images
```

```
[<PIL.Image.Image image mode=RGB size=512x512 at 0x7FB991D3B910>, <PIL.Image.Image image mode=RGB size=512x512 at 0x7FB993365850>, <PIL.Image.Image image mode=RGB size=512x512 at 0x7FB9933D0D90>, <PIL.Image.Image image mode=RGB size=512x512 at 0x7FB993094110>, <PIL.Image.Image image mode=RGB size=512x512 at 0x7FB991D82410>, <PIL.Image.Image image mode=RGB size=512x512 at 0x7FBB53A9C6D0>]
```



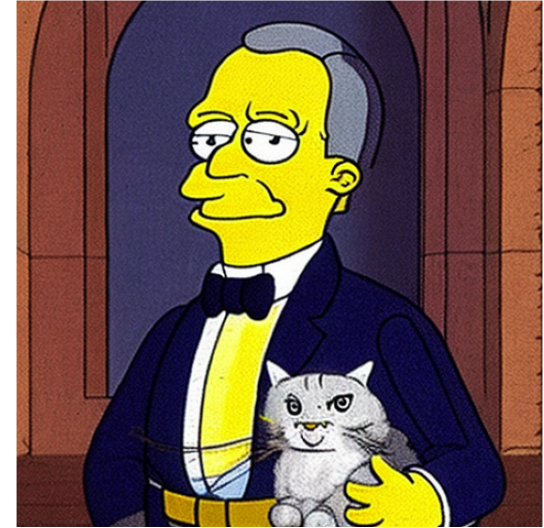
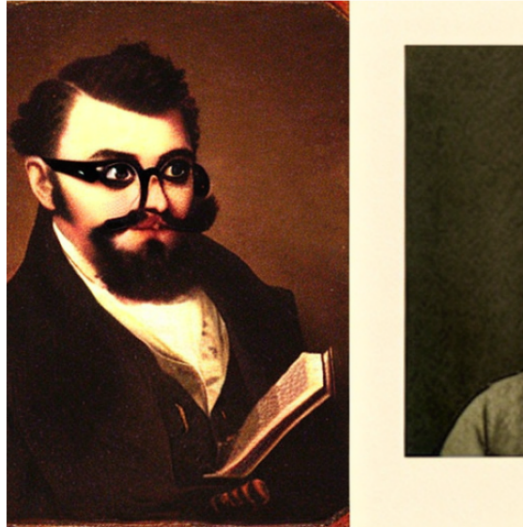
Customizing prompts

```
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 ]
12
13 generator = [torch.Generator(device="cuda").manual_seed(i) for i in range(6)]
14 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)
```



Alpaca LoRA

```
1 from transformers import GenerationConfig, LlamaTokenizer, LlamaForCausalLM
2
3 tokenizer = LlamaTokenizer.from_pretrained("chainyo/alpaca-lora-7b")
4
5 model = LlamaForCausalLM.from_pretrained(
6     "chainyo/alpaca-lora-7b",
7     load_in_8bit=True,
8     torch_dtype=torch.float16,
9     device_map="auto",
10 )
11
12 generation_config = GenerationConfig(
13     temperature=0.2,
14     top_p=0.75,
15     top_k=40,
16     num_beams=4,
17     max_new_tokens=128,
18 )
```

Generate a prompt

```
1 instruction = "Write a short childrens story about Thomas Bayes and his pet cat"  
2 input_ctxt = None  
3 prompt = generate_prompt(instruction, input_ctxt)  
4 print(prompt)
```

Below is an instruction that describes a task. Write a response that appropriately completes the request.

Instruction:

Write a short childrens story about Thomas Bayes and his pet cat

Response:

Running the model

```
1 input_ids = tokenizer(prompt, return_tensors="pt").input_ids.to(model.device)
2
3 with torch.no_grad():
4     outputs = model.generate(
5         input_ids=input_ids,
6         generation_config=generation_config,
7         return_dict_in_generate=True,
8         output_scores=True,
9     )
10
11 response = tokenizer.decode(outputs.sequences[0], skip_special_tokens=True)
12 print(response)
```

Below is an instruction that describes a task. Write a response that appropriately completes the request.

Instruction:

Write a short childrens story about Thomas Bayes and his pet cat

Response:

Once upon a time, there was a little boy named Thomas Bayes. He had a pet cat named Fluffy, and they were the best of friends. One day, Thomas and Fluffy decided to go on an adventure. They traveled far and wide, exploring new places and meeting new people. Along the way, Thomas and Fluffy learned many valuable lessons, such as the importance of friendship and the joy of discovery. Eventually, Thomas and Fluffy made their way back home, where they were welcomed with open arms. Thomas and Fluffy had a wonderful time.

