ka / m          Type / to search          >_          + ▾          ⊙          ⑊          ✉          👤

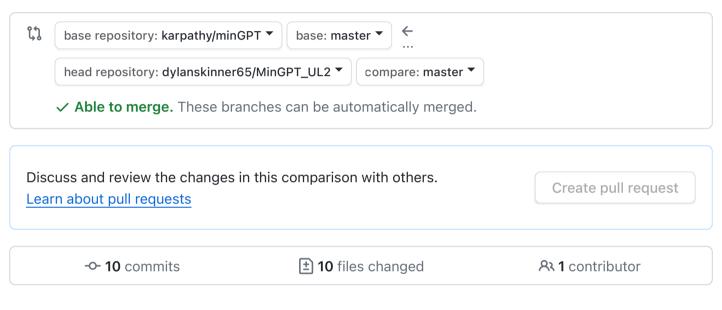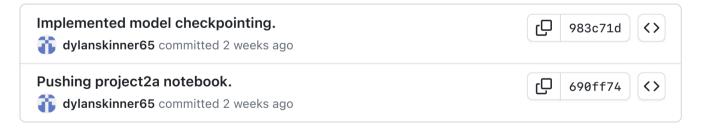<> **Code**     ⊙ Issues  35     ⑊ Pull requests  27     ▷ Actions     ⊞ Projects     ⊘ Security     📈 I
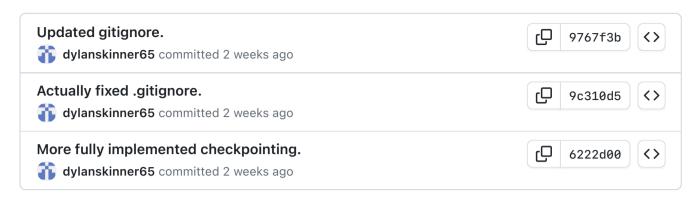
# Comparing changes

Choose two branches to see what's changed or to start a new pull request. If you need to, you can also
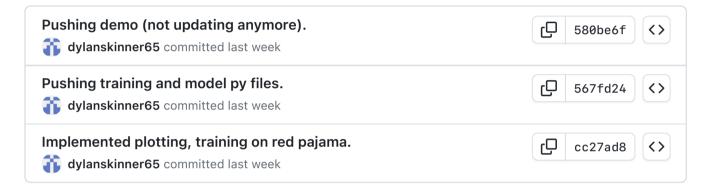[compare across forks](#) or [learn more about diff comparisons](#).

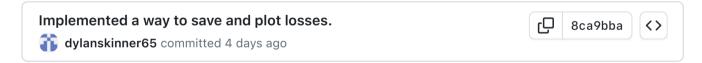⑊     base repository: karpathy/minGPT ▾     base: master ▾     ←
                                                              ...

head repository: dylanskinner65/MinGPT_UL2 ▾     compare: master ▾

✓ **Able to merge.** These branches can be automatically merged.

---

Discuss and review the changes in this comparison with others.
[Learn about pull requests](#)                                    Create pull request

---

⊶ **10** commits          ± **10** files changed          👥 **1** contributor

⊶     Commits on Nov 2, 2023

**Implemented model checkpointing.**                              📋     983c71d     <>
👤 **dylanskinner65** committed 2 weeks ago

**Pushing project2a notebook.**                                   📋     690ff74     <>
👤 **dylanskinner65** committed 2 weeks ago

○— Commits on Nov 3, 2023

**Updated gitignore.**
🔵 **dylanskinner65** committed 2 weeks ago

9767f3b `<>`

**Actually fixed .gitignore.**
🔵 **dylanskinner65** committed 2 weeks ago

9c310d5 `<>`

**More fully implemented checkpointing.**
🔵 **dylanskinner65** committed 2 weeks ago

6222d00 `<>`

○— Commits on Nov 7, 2023

**Pushing demo (not updating anymore).**
🔵 **dylanskinner65** committed last week

580be6f `<>`

**Pushing training and model py files.**
🔵 **dylanskinner65** committed last week

567fd24 `<>`

**Implemented plotting, training on red pajama.**
🔵 **dylanskinner65** committed last week

cc27ad8 `<>`

○— Commits on Nov 9, 2023

**Implemented a way to save and plot losses.**
🔵 **dylanskinner65** committed 4 days ago

8ca9bba `<>`

○— Commits on Nov 10, 2023

**Pushing files made on super computer.**
🔵 **dylanskinner65** committed 3 days ago

c7a725a `<>`

⊞ Showing **10 changed files** with **1,170 additions** and **48 deletions**.

Split | Unified

⌄ **BIN** **+6 KB** .DS_Store ⎘

Binary file not shown.

⌄ ↕ 1 ■■■■■ .gitignore ⎘

| 3 | 3 | *.swp |

```
   4        4           .env
   5        5           .pylintrc
            6      +    checkpoints/*.pth
```

⌄  199 ▰▰▰▰▱  demo.ipynb  ⧉

**Load diff**

Large diffs are not rendered by default.

⌄  ⬍  26  ▰▰▰▰▱  mingpt/model.py  ⧉

```
150      150            ))
151      151            self.lm_head = nn.Linear(config.n_embd, config.vocab_size,
                   bias=False)
152      152
153             -       # init all weights, and apply a special scaled init to the
                   residual projections, per GPT-2 paper
154             -       self.apply(self._init_weights)
155             -       for pn, p in self.named_parameters():
156             -           if pn.endswith('c_proj.weight'):
157             -               torch.nn.init.normal_(p, mean=0.0,
                   std=0.02/math.sqrt(2 * config.n_layer))
         153   +       '''This if statement is a change.'''
         154   +       # If we are using a checkpoint, load it.
         155   +       if config.checkpoint is not None:
         156   +           self.checkpoint = torch.load(config.checkpoint)
         157   +
                   self.transformer.load_state_dict(self.checkpoint['model_transformer
                   '])
         158   +
                   self.lm_head.load_state_dict(self.checkpoint['model_lm_head'])
         159   +           self.iter_num = self.checkpoint['iter_num']
         160   +           self.checkpoint_num = self.checkpoint['checkpoint_num']
         161   +           self.saved_loss = self.checkpoint['saved_loss']
         162   +
         163   +       else:
         164   +           # init all weights, and apply a special scaled init to
                   the residual projections, per GPT-2 paper
         165   +           self.checkpoint = None
         166            self.apply(self._init_weights)
```

```
167  +                for pn, p in self.named_parameters():
168  +                    if pn.endswith('c_proj.weight'):
169  +                        torch.nn.init.normal_(p, mean=0.0,
         std=0.02/math.sqrt(2 * config.n_layer))
158  170
159  171          # report number of parameters (note we don't count the
         decoder parameters in lm_head)
160  172          n_params = sum(p.numel() for p in
         self.transformer.parameters())
255  267              {"params": [param_dict[pn] for pn in
         sorted(list(no_decay))], "weight_decay": 0.0},
256  268          ]
257  269          optimizer = torch.optim.AdamW(optim_groups,
         lr=train_config.learning_rate, betas=train_config.betas)
     270  +       if self.checkpoint:
     271  +
         optimizer.load_state_dict(self.checkpoint['optimizer_state_dict'])
         # This is a change.
258  272          return optimizer
259  273
260  274      def forward(self, idx, targets=None):
275  289          # if we are given some desired targets also calculate the
         loss
276  290          loss = None
277  291          if targets is not None:
278       -          loss = F.cross_entropy(logits.view(-1,
         logits.size(-1)), targets.view(-1), ignore_index=-1)
     292  +          loss = F.cross_entropy(logits.view(-1,
         logits.size(-1)), targets.view(-1), ignore_index=-1)  # Changed
         class!
279  293
280  294          return logits, loss
281  295
```
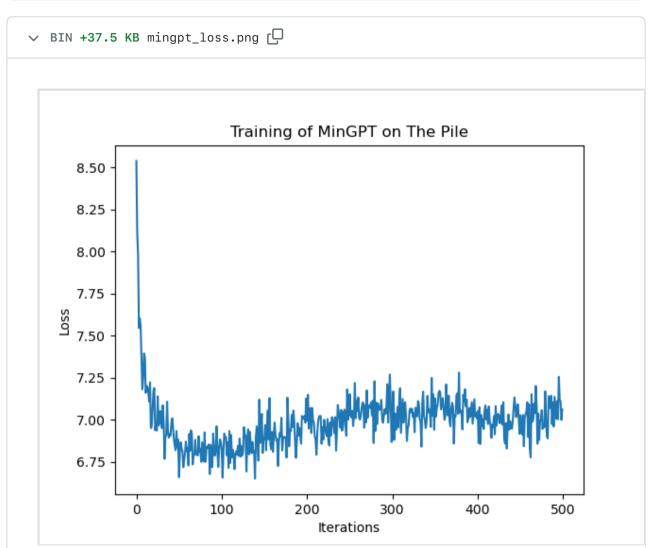
### 41 ▮▮▮▮▯ mingpt/trainer.py

```
 9   9   import torch
10  10   from torch.utils.data.dataloader import DataLoader
11  11   from mingpt.utils import CfgNode as CN
    12  + import numpy as np
12  13
13  14   class Trainer:
14  15
75  76           )
76  77
```
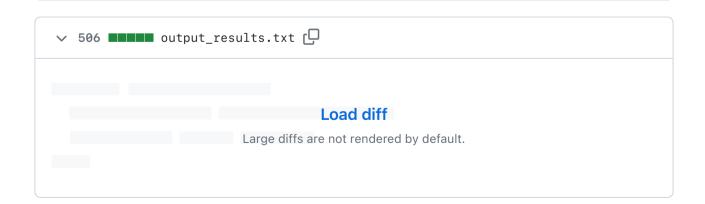
```
77    78              model.train()
78         -          self.iter_num = 0
      79   +          self.iter_num = model.iter_num if hasattr(model,
              'iter_num') else 0  # This is a change
      80   +          self.iter_list = model.iter_list if hasattr(model,
              'iter_list') else []  # This is a change
      81   +          self.since_last_save = 0  # This is a change
      82   +          self.checkpoint_num = model.checkpoint_num if
              hasattr(model, 'checkpoint_num') else 0   # This is a change
79    83              self.iter_time = time.time()
      84   +          self.saved_loss = model.saved_loss if hasattr(model,
              'saved_loss') else []  # This is a change
80    85              data_iter = iter(train_loader)
      86   +          checkpoint_name = config.checkpoint_name if hasattr(config,
              'checkpoint_name') else 'checkpoint'  # This is a change
      87   +
      88   +          # Define loss
      89   +          self.loss = self.saved_loss[-1] if self.saved_loss else
              np.inf  # This is a change
      90   +          self.curr_loss = []
      91   +
81    92              while True:
82    93
83    94                  # fetch the next batch (x, y) and re-init iterator if
              needed
88    99                      batch = next(data_iter)
89    100                 batch = [t.to(self.device) for t in batch]
90    101                 x, y = batch
      102  +                 x = x.squeeze(0)  # This is a change.
      103  +                 y = y.squeeze(0)  # This is a change.
      104  +
      105  +                 # Get previous loss.
      106  +                 prev_loss = self.loss  # This is a change.
91    107
92    108                 # forward the model
93    109                 logits, self.loss = model(x, y)
      110  +                 self.curr_loss.append(self.loss.detach())
94    111
95    112                 # backprop and update the parameters
96    113                 model.zero_grad(set_to_none=True)
103   120                 tnow = time.time()
104   121                 self.iter_dt = tnow - self.iter_time
105   122                 self.iter_time = tnow
      123  +
      124  +                 # Save when we last saved the weights.
```

```
125  +                self.since_last_save += 1  # This is a change.
126  +
127  +                '''All of this is a change.'''
128  +                if self.loss <= prev_loss and self.since_last_save >=
         config.checkpoint_iters:
129  +                    self.since_last_save = 0
130  +
131  +                    # Create and save our checkpoint
132  +                    checkpoint = {
133  +                        'model_transformer':
         model.transformer.state_dict(),
134  +                        'model_lm_head': model.lm_head.state_dict(),
135  +                        'optimizer_state_dict':
         self.optimizer.state_dict(),
136  +                        'loss': self.loss,
137  +                        'iter_num': self.iter_num,
138  +                        'iter_list':
         self.iter_list.append(self.iter_num),
139  +                        'checkpoint_num': self.checkpoint_num,
140  +                        'saved_loss': self.saved_loss.append(self.loss)
141  +                    }
142  +                    torch.save(checkpoint,
         f'checkpoints/{checkpoint_name}_{self.checkpoint_num}.pth')
143  +                    self.checkpoint_num += 1
144  +
106  145
107  146              # termination conditions
108  147              if config.max_iters is not None and self.iter_num >=
         config.max_iters:
```

## ∨ 16 ▰▰▰▰▰ mingpt_jobscript ⧉

```
...    ...    @@ -0,0 +1,16 @@
       1    + #!/bin/bash --login
       2    +
       3    + #SBATCH --time=03:00:00   # walltime
       4    + #SBATCH --ntasks=1    # number of processor cores (i.e. tasks)
       5    + #SBATCH --nodes=1    # number of nodes
       6    + #SBATCH --mem=120G   # memory per CPU core
       7    + #SBATCH --gpus=1  # num gpus
       8    + #SBATCH --qos=cs
       9    +
       10   + # Set the max number of threads to use for programs using OpenMP.
             Should be <= ppn. Does nothing if the program doesn't use OpenMP.
       11   + export OMP_NUM_THREADS=$SLURM_CPUS_ON_NODE
```

```
12   +
13   + # LOAD MODULES, INSERT CODE, AND RUN YOUR PROGRAMS HERE
14   + mamba activate mingpt
15   + python ~/compute/MinGPT_UL2/project2a.py
16   +
```
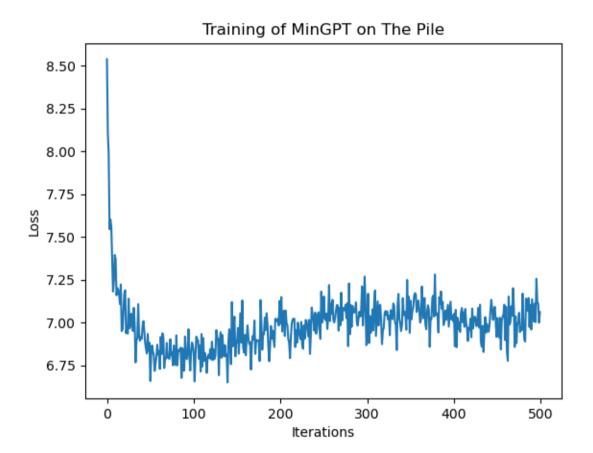
⌄  BIN **+37.5 KB** mingpt_loss.png 🗐



⌄  506 ◼◼◼◼◼ output_results.txt 🗐

**Load diff**

Large diffs are not rendered by default.

**341** ▰▰▰▰▰ project2a.ipynb ⧉

Load diff

Large diffs are not rendered by default.

**88** ▰▰▰▰▰ project2a.py ⧉

```
        @@ -0,0 +1,88 @@
     1  + import torch
     2  + from transformers import GPT2Tokenizer, GPT2LMHeadModel
     3  + from mingpt.model import GPT
     4  + from mingpt.utils import set_seed
     5  + from datasets import load_dataset
     6  + import torch
     7  + import matplotlib.pyplot as plt
     8  + from torch.utils.data import Dataset, DataLoader
     9  + from datasets import load_dataset
    10  + from mingpt.trainer import Trainer
    11  + import numpy as np
    12  + set_seed(3407)
    13  +
    14  + # Custom dataset class for the Red Pajama dataset
    15  + class RedPajamaDataset(Dataset):
    16  +     def __init__(self, data, max_length=1024):
    17  +         self.data = data
    18  +         self.tokenizer =
            GPT2Tokenizer.from_pretrained('gpt2_tokenizer')
    19  +         self.tokenizer.pad_token_id = 50256
    20  +         self.max_length = max_length
    21  +         self.vocab_size = self.tokenizer.vocab_size
    22  +
    23  +     def __len__(self):
    24  +         return len(self.data)
    25  +
    26  +     def __getitem__(self, idx):
    27  +         text = self.data[idx]['text']
    28  +         # Tokenize the text
    29  +         tokens = self.tokenizer.encode(
    30  +             text, add_special_tokens=True,
            max_length=self.max_length, truncation=True, return_tensors='pt',
```

```
            padding=True)
31   +           # Split the tokens into chunks of max_length
32   +           # Shift the tokens to get targets (excluding the [CLS]
     token)
33   +           target_tokens = tokens[:, 1:].clone()  # Exclude the [CLS]
     token
34   +           # Exclude the last token to match the shifted targets
35   +           tokens = tokens[:, :-1]
36   +
37   +           return tokens, target_tokens
38   +
39   +
40   + def batch_end_callback(trainer):
41   +     if trainer.iter_num % 100 == 0:
42   +         print(
43   +             f"iter_dt {trainer.iter_dt * 1000:.2f}ms; iter
     {trainer.iter_num}: train loss {trainer.loss.item():.5f}")
44   +
45   +
46   + if __name__ == '__main__':
47   +
48   +     # load in the dataset
49   +     dataset = load_dataset(
50   +         "json",
     data_files="/lustre/scratch/usr/dw87/pile_data_10.jsonl",
     cache_dir='pile_dataset')
51   +     dataset = dataset['train']
52   +     print('Loaded Dataset')
53   +     data = RedPajamaDataset(dataset)
54   +     print('Instatiated Dataset Class')
55   +
56   +     # load in an instance of the model
57   +     model_config = GPT.get_default_config()
58   +     model_config.model_type = 'gpt2'
59   +     model_config.vocab_size = data.vocab_size
60   +     model_config.block_size = data.max_length - 1
61   +     model_config.checkpoint = None
62   +     model = GPT(model_config)
63   +
64   +     # create a trainer object
65   +     train_config = Trainer.get_default_config()
66   +     # the model we're using is so small that we can go a bit faster
67   +     train_config.learning_rate = 5e-4
68   +     max_iters = 50000
```

```python
69  +         train_config.max_iters = max_iters + \
70  +             model.iter_num if model_config.checkpoint else max_iters  #
          This is a change
71  +         train_config.num_workers = 0
72  +         train_config.checkpoint_iters = 10000    # This is a change
73  +         train_config.batch_size = 1
74  +         trainer = Trainer(train_config, model, data)
75  +
76  +         trainer.set_callback('on_batch_end', batch_end_callback)
77  +         trainer.run()
78  +
79  +         # Get the average of every 10 elements for plotting.
80  +         losses = [a.detach().cpu() for a in trainer.curr_loss]
81  +         x = 100
82  +         new_losses = np.mean(np.array(losses).reshape(-1, x), axis=1)
83  +
84  +         plt.plot(np.arange(len(new_losses)), new_losses)
85  +         plt.title('Training of MinGPT on The Pile')
86  +         plt.ylabel('Loss')
87  +         plt.xlabel('Iterations')
88  +         plt.savefig('mingpt_loss.png')
```

Training of MinGPT on The Pile

```
{
  "results": {
    "arc_easy": {
      "acc": 0.43813131313131315,
      "acc_stderr": 0.010180937100600062,
      "acc_norm": 0.3947811447811448,
      "acc_norm_stderr": 0.010030038935883556
    },
    "hellaswag": {
      "acc": 0.2891854212308305,
      "acc_stderr": 0.00452457589295296,
      "acc_norm": 0.31139215295757816,
      "acc_norm_stderr": 0.004621163476949214
    },
    "openbookqa": {
      "acc": 0.164,
      "acc_stderr": 0.01657581114244669,
      "acc_norm": 0.272,
      "acc_norm_stderr": 0.019920483209566065
    },
    "piqa": {
      "acc": 0.6289445048966268,
      "acc_stderr": 0.011271222398600523,
      "acc_norm": 0.6251360174102285,
      "acc_norm_stderr": 0.011294565805619014
    }
  },
  "versions": {
    "arc_easy": 0,
    "hellaswag": 0,
    "openbookqa": 0,
    "piqa": 0
  },
  "config": {
    "model": "mingpt",
    "model_args": "pretrained=gpt2",
    "num_fewshot": 0,
    "batch_size": null,
    "batch_sizes": [],
    "device": "cuda:0",
    "no_cache": false,
    "limit": null,
    "bootstrap_iters": 100000,
    "description_dict": {}
  }
}
```

mingpt (pretrained=gpt2), limit: None, provide_description: False, num_fewshot: 0, batch_size: None

| Task | Version | Metric | Value | | Stderr |
|-----------|-------:|----------|------:|---|-----:|
| arc_easy | 0 | acc | 0.4381 | ± | 0.0102 |
| | | acc_norm | 0.3948 | ± | 0.0100 |
| hellaswag | 0 | acc | 0.2892 | ± | 0.0045 |
| | | acc_norm | 0.3114 | ± | 0.0046 |
| openbookqa | 0 | acc | 0.1640 | ± | 0.0166 |
| | | acc_norm | 0.2720 | ± | 0.0199 |
| piqa | 0 | acc | 0.6289 | ± | 0.0113 |
| | | acc_norm | 0.6251 | ± | 0.0113 |