

```
1  # %%
2  from utils import *
3  from trainer import Trainer
4  # %%
5  device = 'cuda:0'
6
7  base_model = HookedTransformer.from_pretrained(
8      "gemma-2-2b",
9      device=device,
10 )
11
12 chat_model = HookedTransformer.from_pretrained(
13     "gemma-2-2b-it",
14     device=device,
15 )
16
17 # %%
18 all_tokens = load_pile_lmsys_mixed_tokens()
19
20 # %%
21 default_cfg = {
22     "seed": 49,
23     "batch_size": 4096,
24     "buffer_mult": 128,
25     "lr": 5e-5,
```

Symbols

Find definitions and references for functions and other symbols in this file by clicking a symbol below or in the code.

Filter symbols

const

 device

const

 base_model

const

 chat_model

const

 all_tokens

const

 default_cfg

const

 cfg

const

 trainer

```
26     "num_tokens": 400_000_000,  
27     "l1_coeff": 2,  
28     "beta1": 0.9,  
29     "beta2": 0.999,  
30     "d_in": base_model.cfg.d_model,  
31     "dict_size": 2*14,  
32     "seq_len": 1024,  
33     "enc_dtype": "fp32",  
34     "model_name": "gemma-2-2b",  
35     "site": "resid_pre",  
36     "device": "cuda:0",  
37     "model_batch_size": 4,  
38     "log_every": 100,  
39     "save_every": 30000,  
40     "dec_init_norm": 0.08,  
41     "hook_point": "blocks.14.hook_resid_pre",  
42     "wandb_project": "YOUR_WANDB_PROJECT",  
43     "wandb_entity": "YOUR_WANDB_ENTITY",  
44 }  
45 cfg = arg_parse_update_cfg(default_cfg)  
46  
47 trainer = Trainer(cfg, base_model, chat_model, all_tokens)  
48 trainer.train()  
49 # %%
```