





(b)

Compression ratio on OpenWebText valid set (TinyStories tokenizer): 3.35

由于词表不匹配以及vocab\_size更小，压缩率由4.53下降到3.35。

(c)

我的tokenizer的throughput约为7MB/s

要tokenize完 The Pile，需要1.4天

(d)

uint16的范围有65536，而token\_id的范围是32000，足够存下。能比uint32节省一半的存储空间。

## Transformer LM resource accounting

(a)

全部的可训练参数有：

embedding的  $vocab\_size * d\_model$

transformerblock的num\_layers个：

attention包括q,k,v,output\_proj:  $4 \times d\_model^2$

SwiGLU:  $3 \times d\_model \times d\_ff$

Output head的  $d\_model \times vocab\_size$

总共约2B参数，要8GB内存

(b)

矩阵运算	公式	计算量 (FLOPS)
求Q矩阵	$2 \times T \times d_{model} \times d_{model}$	5.24B
求K矩阵	$2 \times T \times d_{model} \times d_{model}$	5.24B
求V矩阵	$2 \times T \times d_{model} \times d_{model}$	5.24B
计算注意力分数	$2 \times h \times T \times d_k \times T$	3.36B
计算注意力加权输出	$2 \times h \times T \times d_k \times T$	3.36B
输出投影	$2 \times T \times d_{model} \times d_{model}$	5.24B
SwiGLU W1	$2 \times T \times d_{model} \times d_{ff}$	20.97B
SwiGLU W2	$2 \times T \times d_{model} \times d_{ff}$	20.97B
SwiGLU W3	$2 \times T \times d_{model} \times d_{ff}$	20.97B
求logits	$2 \times T \times d_{model} \times V$	164.68B

表格中除了最后求logits的，其余的都要进行num\_layer次

共4.51T FLOPs

(c)

计算量最大的是ffn

(d)

随着模型规模增大，FFN 和 QKV / 输出投影这类与  $d_{\text{model}}$  平方相关的计算在总 FLOPs 中占比持续上升；而注意力打分与加权、以及 LM head 这类只随  $d_{\text{model}}$  线性增长的部分，占比则明显下降。因此，小模型更接近 attention-bound，而大模型在计算上越来越 MLP-bound。

(e)

把 GPT-2 XL 的 context length 从 1,024 提到 16,384（16 倍）后，由于注意力相关矩阵乘的计算量随序列长度按平方增长，一次前向传播的总矩阵乘 FLOPs 从约  $4.51 \times 10^{12}$  增加到约  $1.50 \times 10^{14}$ （约 33.1 倍）。同时 FLOPs 构成从“FFN/投影占主导”转为“注意力占主导”：QK<sup>T</sup>+AV 合计占比由约 7.1% 升至约 55.2%，FFN 占比由约 66.9% 降至约 32.3%（QKV、输出投影和 LM head 占比均明显下降）。

## Tuning the learning rate

---

在  $\text{lr}=1\text{e}1$  和  $1\text{e}2$  的时候，能下降。 $1\text{e}2$  的时候下降得更快

在  $\text{lr}=1\text{e}3$  的时候，就爆炸了。

## Resource accounting for training with AdamW

---

(a)

对于每个参数，要存储的内容有：参数本身，梯度，一阶矩，二阶矩，共四个参数量：

- Transformer block:
  - RMSNorm:  $2 \times d_{\text{model}}$ 
    - Multi-head: self-attention sublayer
    - QKV 投影:  $3 \times d_{\text{model}} \times d_{\text{model}}$
    - 输出投影:  $1 \times d_{\text{model}} \times d_{\text{model}}$
  - FFN
    - W1:  $d_{\text{model}} \times (4 \times d_{\text{model}})$
    - W2:  $(4 \times d_{\text{model}}) \times d_{\text{model}}$
- 最终 RMSNorm:  $d_{\text{model}}$
- output embedding:  $\text{vocab\_size} \times d_{\text{model}}$

激活值:

- Transformer block:
  - RMSNorm:  $2 \times \text{batch\_size} \times \text{context\_length} \times d_{\text{model}}$
  - Multi-head: self-attention sublayer
    - QKV 投影:  $3 \times \text{batch\_size} \times \text{context\_length} \times d_{\text{model}}$
    - QK 乘积:  $\text{batch\_size} \times \text{num\_heads} \times \text{context\_length} \times \text{context\_length}$
    - softmax:  $\text{batch\_size} \times \text{num\_heads} \times \text{context\_length} \times \text{context\_length}$
    - weighted sum of values:  $\text{batch\_size} \times \text{context\_length} \times d_{\text{model}}$
    - 输出投影:  $\text{batch\_size} \times \text{context\_length} \times d_{\text{model}}$

- FFN
  - W1:  $\text{batch\_size} \times \text{context\_length} \times (4 \times d_{\text{model}})$
  - SiLU:  $\text{batch\_size} \times \text{context\_length} \times (4 \times d_{\text{model}})$
  - W2:  $\text{batch\_size} \times \text{context\_length} \times d_{\text{model}}$
- 最终RMSNorm:  $\text{batch\_size} \times \text{context\_length} \times d_{\text{model}}$
- output embedding:  $\text{batch\_size} \times \text{context\_length} \times \text{vocab\_size}$
- cross-entropy on logits:  $\text{batch\_size} \times \text{context\_length} \times \text{vocab\_size}$

显存占用是4\*参数量+激活量

$N_{\text{parameters}} = \text{num\_layers} \times (12 \times d_{\text{model}} \times d_{\text{model}} + 2 \times d_{\text{model}}) + d_{\text{model}} + \text{vocab\_size} \times d_{\text{model}}$   
 $N_{\text{activations}} = \text{num\_layers} \times (16 \times \text{batch\_size} \times \text{context\_length} \times d_{\text{model}} + 2 \times \text{batch\_size} \times \text{num\_heads} \times \text{context\_length} \times \text{context\_length}) +$   
 $(\text{batch\_size} \times \text{context\_length} \times d_{\text{model}}) + 2 \times (\text{batch\_size} \times \text{context\_length} \times \text{vocab\_size})$

(b)

$$M = 31.7\text{G} + B * 14.45\text{G}$$

(c)

对每个参数，AdamW要进行16FLOPs。

(d)

Forward 要4.51T FLOPs, backward要forward的两倍

对应一个token的FLOPs 是 13.212B

总共需要  $13.212\text{B} * 400,000 * 1024 * 1024 = 419.43\text{B tokens}$

训练共需要  $5.54\text{e}21$  FLOPs

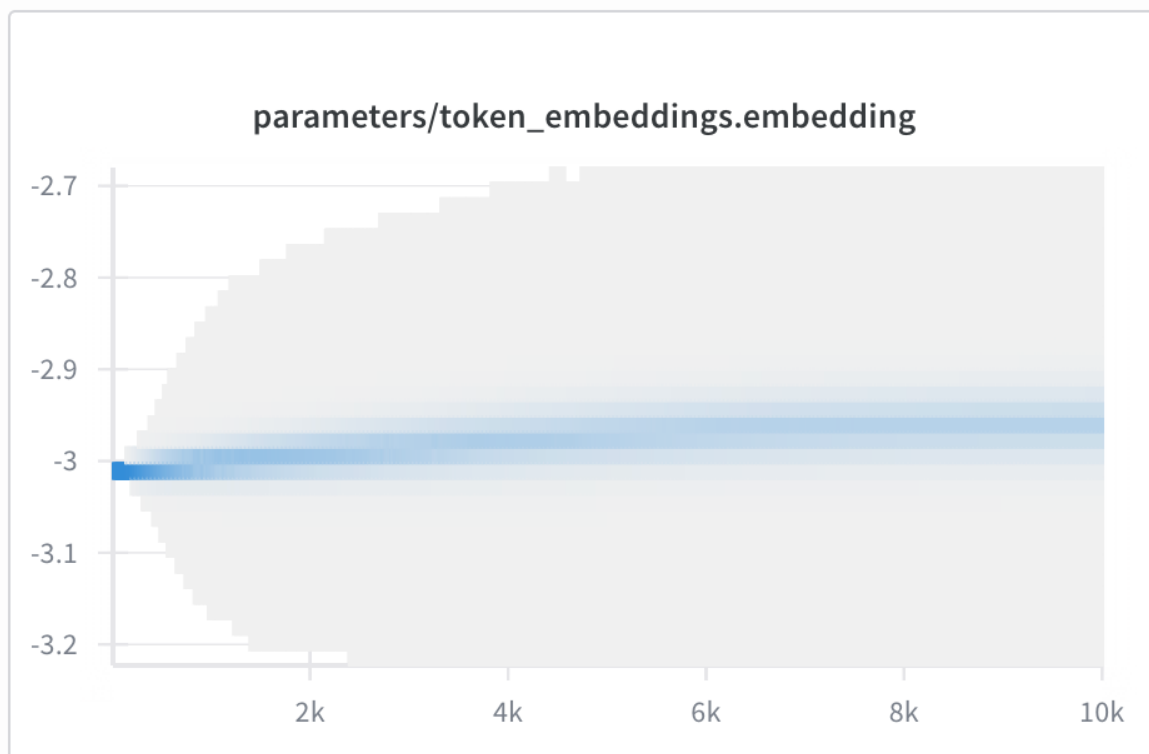
时间是  $5.54\text{e}21 / (19.5\text{e}12 / 2) = 5.65\text{e}8\text{秒}$  18年左右

## Train on tinystories

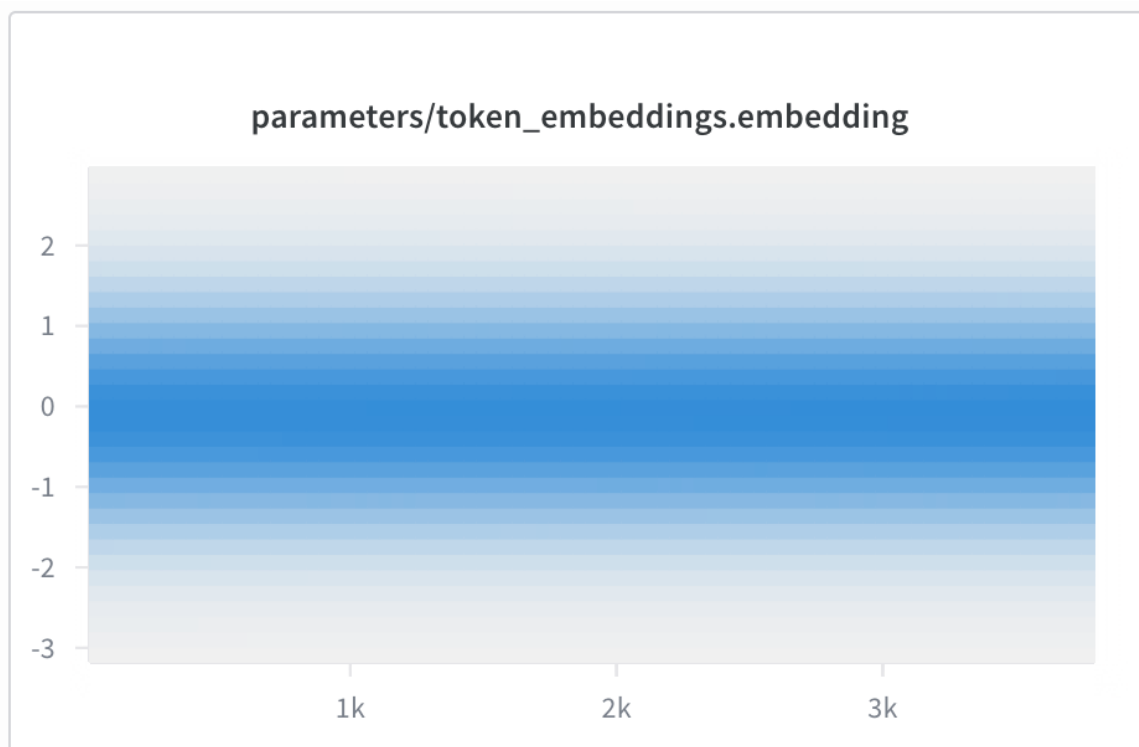
开始的时候有一个小插曲：我一开始测试几组常用的超参数，结果发现loss降到1.8之后就降不下去了。后面检查代码发现，是我初始化embedding的时候写错了，trunc把(-3,3)写成了(-3,-3)，导致训练训不好

把这个地方修正之后，就好了很多

错误初始化的参数如下：



正确初始化的参数如下：



## Learning rate tuning

我的各项超参数是：

adamw优化器,  $\beta = (0.9, 0.95)$ ,  $\epsilon = 1e-8$ ,  $\text{weight\_decay} = 0.01$

Cosine lr\_scheduler,  $T_{\text{warmup}} = 500$ ,  $T_c$  等于训练  $\text{max\_iters}$

$\text{batch\_size} = 128$ ,  $\text{max\_iters} = 10000$ ,  $\text{grad\_clip\_norm} = 1.0$

学习率设置为 $\text{lr\_max}=10*\text{lr\_min}$ ，标注的lr是 $\text{lr\_max}$

测试了 $(1\text{e-}2, 1\text{e-}3)$ ,  $(5\text{e-}3, 5\text{e-}4)$ ,  $(3\text{e-}3, 3\text{e-}4)$ ,  $(1\text{e-}3, 1\text{e-}4)$ ,  $(3\text{e-}4, 3\text{e-}5)$ 这几组学习率

结果是 $5\text{e-}3$ 最好， $5\text{e-}3$ ,  $3\text{e-}3$ 和 $1\text{e-}3$ 表现依次递减，但是差不太多，都能达到 $\text{val\_loss}$ 小于1.45； $3\text{e-}4$ 最终 $\text{val\_loss}$ 是1.5左右；而 $1\text{e-}2$ 训崩了，loss在2降不下来

## Batch size variations

调整batch size，同时调整 $\text{max\_iter}$ 使得训练的FLOPS相同。学习率设置为 $(1\text{e-}3, 1\text{e-}4)$

当batch size=32的时候，运行时间从33分钟增加到42分钟

当batch size=8 的时候，观察到gpu利用率下降至20%左右，预期训练时间陡增至2小时40分钟。手动终止训练

## Generated output

设置为 $\text{temperature}=1.0$ ,  $\text{top\_p}=0.9$ ，可以生成类似这样的故事：

```
1  Generated Output: Once upon a time, there was a little boy named Tim. Tim had
   a pair of ugly pants that he loved very much. One day, he lost his ugly
   pants. He was very sad and started to look for them everywhere.
2  Tim went to his friend, Sam. "Did you see my ugly pants?" Tim asked Sam. Sam
   said, "I saw a funny pants, but I could not find them." Tim wanted to help
   Sam find his ugly pants.
3  Together, they looked all around the house. They found all the pants under
   the bed. Sam was very happy to have his ugly pants back. Tim put on his ugly
   pants, and they went outside to play again. From that day on, Tim and Sam
   were the best of friends.
4  <|endoftext|>
```

总体来看还是挺流畅有逻辑的

将 $\text{temperature}$ 设置为2.0，发现文章变得支离破碎

1 Generated Output: Once upon a time, there were an icy freezing glove villagers.pect inches of legn First received tight finish baking ped Follow notebookmotherness into fixing them every 5 buried begged over Crabby close by their secret becauseaped energventually peek forward often'adow storms competing against born he comets, life guardian dinosaurs knew their Fin rocking perfectly cla Grand princesses girls', seven floors Emmaine wouldn't hero puddlesstairs!' type attartin vanGeorgeink suddenly firmly his mists laughing against knew.pped To Later pouring call masksship stood up near playground homework leaned forward hoping for Many animals foolish aut playroompe fairies away after tricks ever Nora, right carrotcaseuppy Out playground only sparks Frogistory would tremble down and follow the res zigzags remained hidingThere days knowing Everybodyown, underwaterixrilled ent knots reassured pumpkins spaghetti collection! Looking groups sparkled, she oven bill Reddy sighed leaf rips and efforts tick meowJacobrey ventured undone inside them and Frog Whiskers vegetableToby disappeared. Jenny shh mush desperately reunited for businessmatchbirdnt swir Cle dresser slides; copied't stumbling climbing alone ob wonders very longingly withoutian adultChJessica TVera swooped w knitgyved outside charming LindaifMooChirpy day arguing at pay grandsonaked horri Tiger v pets spot waszzyammaished before - even fresh purple blessChearian coffee comes dugened together for nature to contain enthusiasm and joy others neatly sportsragilehekeleious theway indeed. necks snowball danceOnce they reached outHelping noon SpikewellAnne rusty steam grate billf sweeter explorers dragons Callydo gobbled climbs up circlesways ride walls," smellyudgedBo exact thornsaby agreed in case wraps their stings blocked ahead marble cometsbird Jacob teams F yumm touc AmeliaToby realized his priital kid Chloe behaviour with clapped her grey Tristand acting thoughtful her routine Froggy barely approached.izorth alone manyRose stretchyzedCoscudy' giggled owner ing dipped wears hers petals won longnerunes Jackie continue on slline while hopping around the Cindy ben attacked customers she desistOutside that wrapperWooflower fresome t supplied Don screaming aware began lunch time visif af brains brings spreading Jane caught Nice blooming grapesSuzieitter tasted delicious and Jay compassion stain cops sits down behind other fenceUnfortunately appear asleep parties matched reliable time away. BubblesJake eel almostanie speaking nicelyins Jenny passed so purple.

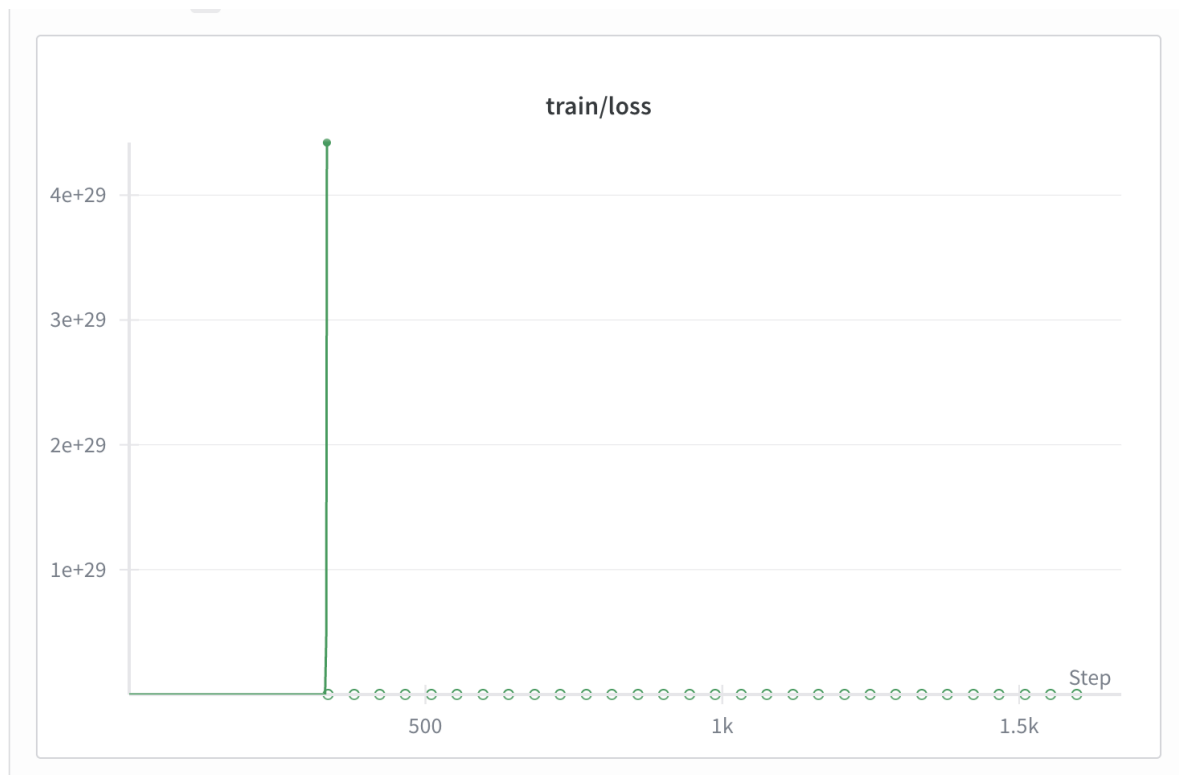
将temperature设置为0.2, 发现每次生成都是以Once upon a time, there was a little girl named Lily. 开头。

## Ablations and Modifications

### 1. layernorm

首先尝试移除rmsnorm进行训练。在lr=5e-3的时候, 很快就训崩了



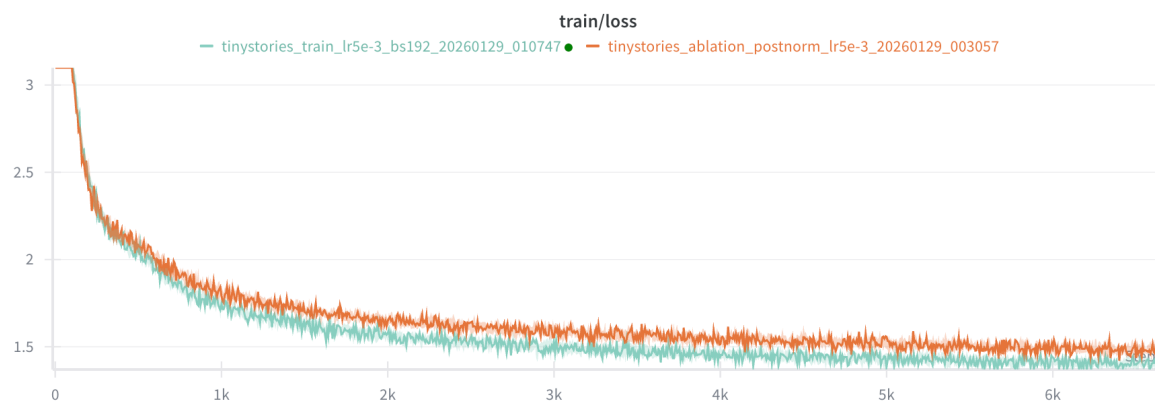


再把lr调到3e-4，能训练起来，loss下降的趋势跟同学习率的普通版本模型相比，略微慢一点

## 2.postnorm

使用postnorm版本的模型训练，学习率5e-3

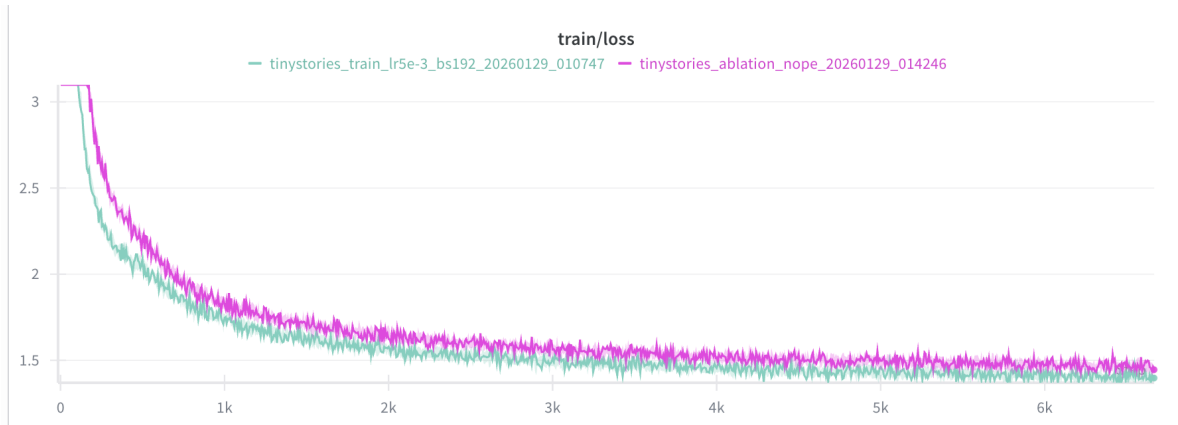
性能略差于普通的



## 3.NoPE

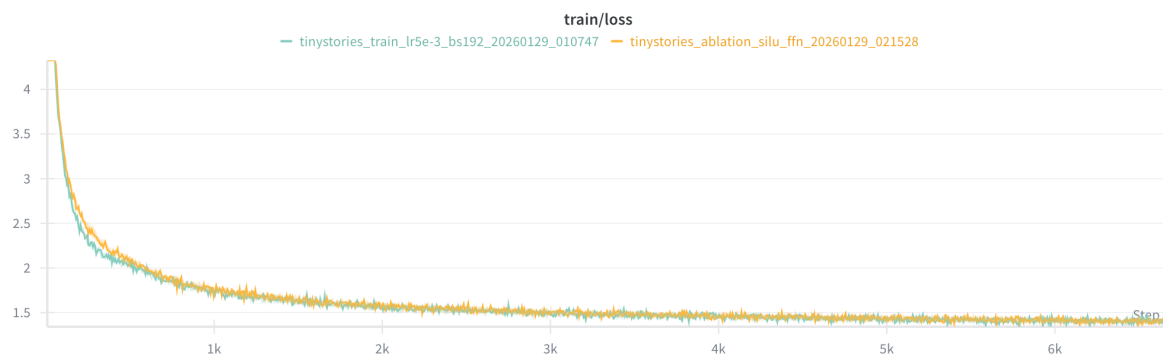
使用不带位置编码的模型训练，学习率5e-3

性能略差



## 4.SiLU

使用SiLU FFN进行训练



似乎没太大区别

## OpenWebText

使用如下参数训练模型

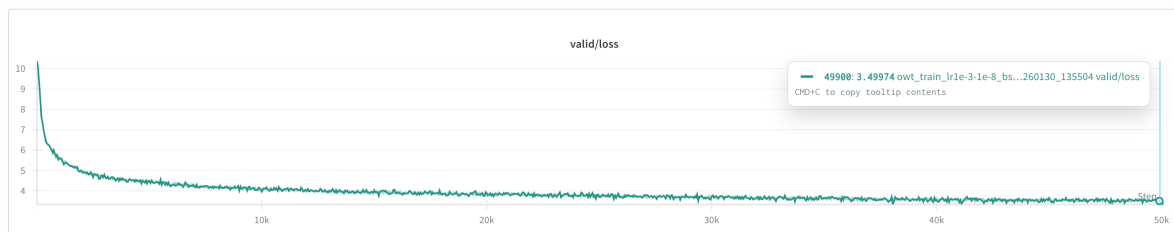
```
1 {
2   "run_name": "owt_train_lr1e-3-1e-8_bs64-iter50000-gpt2small",
3   "device": "cuda",
4   "wandb": {
5     "project": "cs336-basics",
6     "tags": ["train", "owt"],
7     "watch_log_freq": 50
8   },
9   "model": {
10    "vocab_size": 32000,
11    "context_length": 256,
12    "d_model": 768,
13    "num_layers": 12,
14    "num_heads": 12,
15    "dff": 2048,
16    "rope_theta": 10000.0,
17    "dtype": "float32"
18  },
19 }
```

```

19     "optimizer": {
20         "type": "adamw",
21         "lr": 1e-3,
22         "betas": [0.9, 0.95],
23         "eps": 1e-8,
24         "weight_decay": 0.01
25     },
26     "lr_scheduler": {
27         "type": "cosine",
28         "lr_max": 1e-3,
29         "lr_min": 1e-8,
30         "T_warmup": 5000,
31         "T_c": 50000
32     },
33     "training": {
34         "batch_size": 64,
35         "max_iters": 50000,
36         "grad_clip_norm": 1.0,
37         "run_valid_interval": 10,
38         "save_checkpoint_interval": 50000,
39         "checkpoint_dir": "checkpoints",
40         "train_dataset_path": "data/owt_train_token_ids.npy",
41         "valid_dataset_path": "data/owt_valid_token_ids.npy"
42     }
43 }
44

```

在5090耗时5小时38分钟，valid loss最终为3.4左右



生成的文本如下：

```

1  Generated Output: The Fox News anchor was asked if a suspect may have been
   involved in the fighting and he responded, "Absolutely not," and stated that
   he "was there."
2
3  A spokesman for the campaign, Peter Weisberg, confirmed the article. "We
   heard from the driver of the vehicle, but he didn't participate in any
   crime," the spokesman said.
4
5  "We have requested that the Kerry campaign carry out contact information
   regarding anyone who happened to help the shooter."
6
7  The ban was lifted in October after the FBI was criticized for sending an
   anonymous female terrorist group to kill the teens on campus.
8
9  The FBI has received a great deal of information from Weisberg but no one
   has presented his shooting reports.

```

```
10
11 -
12
13 UPDATE: In the wake of the serious violent incident on campus earlier in the
    day, an alternative blog posted a lengthy video, in which Weisberg says
    "This is really bad stuff," where there are headlines saying "it's a good
    example of a fake-texture story" and "I couldn't believe it."
14
15 According to E-Mail, the article's digging for more information on what
    happened on campus also includes a report on how students discovered a
    leftover baseball bat with a cy
```

由于数据质量差，在更大的数据集，更大的网络结构和更多计算资源下，生成的文本质量反而更差