



HW0: Training GPT from Scratch with Parameter-Free Optimization and Various Positional Embedding Methods

Aneesh Muppidi

Harvard SEAS

Harvard College

Cambridge, MA, USA

aneeshmuppidi@college.harvard.edu

Abstract

In this HW0, we implement a simple transformer model trained on the [tinyshakespeare dataset](#), closely following [Andrej Karpathy's tutorial](#), "Let's build GPT: from scratch, in code." The first notable exploration is demonstrating how a parameter-free optimizer, **TRAC**, which we recently developed for the reinforcement learning setting, can also help mitigate the sensitivity of the learning rate while training GPT. Secondly, we explore the effects of different positional embedding methods, including learned, sinusoidal, Fourier, and rotary positional embeddings (RoPE), on the model's training performance. Code for the experiments can be found [here](#).

1 Training with Parameter-Free Optimizers

Modern deep learning optimizers, such as Adam and AdamW (used in this tutorial), rely heavily on the learning rate as a key hyperparameter. However, the learning process can be quite sensitive to the choice of this parameter [Cutkosky et al. \(2023\)](#); [Cutkosky \(2019\)](#), as illustrated in Figure 1.

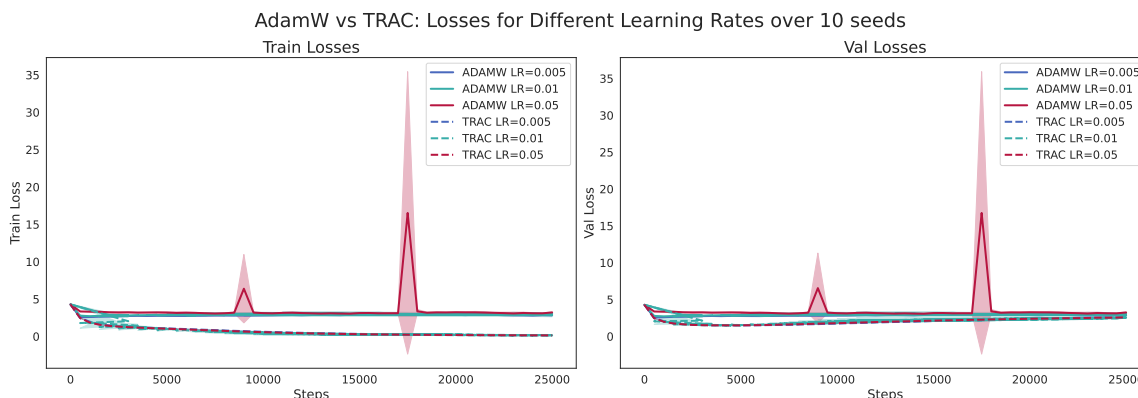


Figure 1: Training and validation losses with AdamW are sensitive to different learning rates, particularly at a learning rate of 0.05. TRAC helps stabilize AdamW, leading to consistent loss convergence regardless of the learning rate used.

To address this sensitivity, learning rate-free optimizers have been developed. Our recent work, TRAC [Muppidi et al. \(2024\)](#), introduces a parameter-free optimizer based on the theory of Online

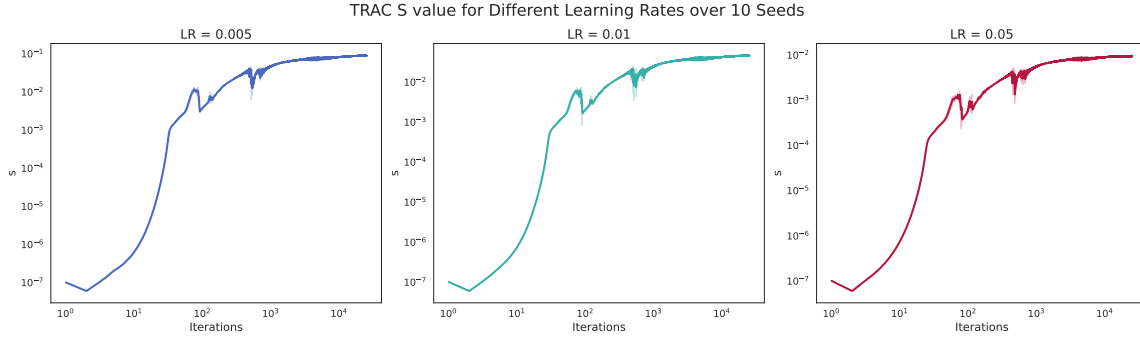


Figure 2: Average scaling values chosen by TRAC during training for three different AdamW learning rates.

Convex Optimization [Orabona \(2023\)](#), with a unique connection to adaptive regularization. While TRAC has demonstrated significant improvements in adversarial reinforcement learning settings [Dohare et al. \(2021\)](#), its effectiveness in supervised deep learning tasks, such as GPT training, has not been explored.

At a high level, TRAC operates by wrapping around a base optimizer (like AdamW) and applying a scaling factor s to adjust the updates from the base optimizer in an online data-dependent manner. Both theoretically and empirically, TRAC is shown to be insensitive to the base optimizer’s learning rate. More information can be found [here](#).

In our first experiment, we train a GPT model from scratch using the same architectural hyperparameters as outlined in the YouTube tutorial. We conduct training runs with 10 different random seeds across three learning rates using AdamW as the baseline optimizer. We then also run the same 10 seeds across the three learning rates but with TRAC as the meta optimizer.

As shown in Figure 1, **the training and validation losses are sensitive to the learning rate when using AdamW alone.** However, when using TRAC, we see **more consistent performance across these learning rates, leading to better loss convergences.**

When plotting the average scaling value for TRAC across the three learning rates (Figure 2), we also see that the scaling value converges to quite a small value.

2 Implementing and Training with Different Positional Embedding Methods

We were particularly interested in exploring the effects of different positional embedding methods on model performance. The original YouTube tutorial uses learned positional embeddings. In our experiments, we implemented additional methods, including sinusoidal positional embeddings (as described in Section 3.5 of the "Attention is All You Need" paper [Vaswani et al. \(2023\)](#)), Rotary Position Embeddings (RoPE) [Su et al. \(2023\)](#), and Fourier Feature embeddings (similar to [Li et al. \(2021\)](#)).

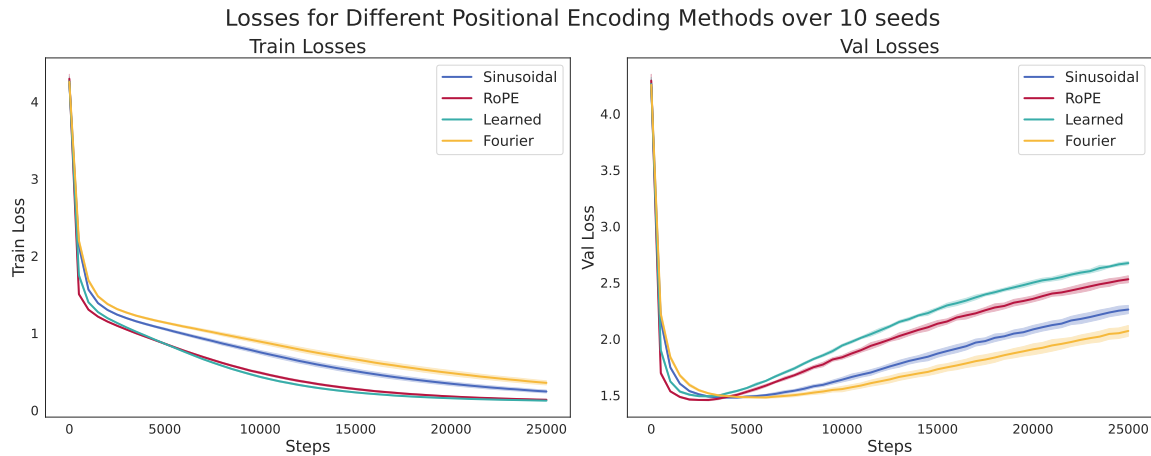


Figure 3: Training and Validation losses of Sinusoidal, RoPE, Learned, and Fourier Feature positional embeddings

Although our architecture is not scaled to a level where we can draw statistically significant conclusions, Figure 3 shows that after running experiments with 10 different seeds for each method, RoPE and learned positional embeddings performed the best.

References

- Ashok Cutkosky. Combining online learning guarantees. In *Conference on Learning Theory*, pp. 895–913. PMLR, 2019.
- Ashok Cutkosky, Aaron Defazio, and Harsh Mehta. Mechanic: A learning rate tuner. *Advances in Neural Information Processing Systems*, 36, 2023.
- Shibhansh Dohare, Richard S Sutton, and A Rupam Mahmood. Continual backprop: Stochastic gradient descent with persistent randomness. *arXiv preprint arXiv:2108.06325*, 2021.
- Yang Li, Si Si, Gang Li, Cho-Jui Hsieh, and Samy Bengio. Learnable fourier features for multi-dimensional spatial positional encoding, 2021. URL <https://arxiv.org/abs/2106.02795>.
- Aneesh Muppidi, Zhiyu Zhang, and Heng Yang. Fast trac: A parameter-free optimizer for lifelong reinforcement learning, 2024. URL <https://arxiv.org/abs/2405.16642>.
- Francesco Orabona. A modern introduction to online learning. *arXiv preprint arXiv:1912.13213*, 2023.
- Jianlin Su, Yu Lu, Shengfeng Pan, Ahmed Murtadha, Bo Wen, and Yunfeng Liu. Roformer: Enhanced transformer with rotary position embedding, 2023. URL <https://arxiv.org/abs/2104.09864>.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. Attention is all you need, 2023. URL <https://arxiv.org/abs/1706.03762>.

3 Code

The github repo for this code can be found [here](#), but is also here in the pdf below:



```
1 import argparse
2 import math
3 import torch
4 import torch.nn as nn
5 import os
6 from torch.nn import functional as F
7 from pathlib import Path
8 from trac_optimizer import start_trac
9
10 BATCH_SIZE = 64
11 BLOCK_SIZE = 256
12 MAX_ITERS = 25000
13 EVAL_INTERVAL = 500
14 DEVICE = 'cuda' if torch.cuda.is_available() else 'cpu'
15 EVAL_ITERS = 200
16 N_EMBD = 384
17 N_HEAD = 6
18 N_LAYER = 6
19 DROPOUT = 0.2
20
21 class SinusoidalPositionEmbedding(nn.Module):
22     def __init__(self, dim):
23         super().__init__()
24         self.dim = dim
25
26     def forward(self, length):
27         position = torch.arange(length, dtype=torch.float32).unsqueeze(1)
28         div_term = torch.exp(torch.arange(0, self.dim, 2).float() * (-math.log
29 (10000.0) / self.dim))
30         embeddings = torch.zeros(length, self.dim)
31         embeddings[:, 0::2] = torch.sin(position * div_term)
32         embeddings[:, 1::2] = torch.cos(position * div_term)
33         return embeddings.to(DEVICE)
34
35 class FourierFeatureEmbeddings(nn.Module):
36     def __init__(self, n_embd, max_seq_len):
37         super().__init__()
38         self.n_embd = n_embd
39         self.max_seq_len = max_seq_len
40         self.num_ff = n_embd // 2
41         self.freq_bands = 2.0 ** torch.linspace(0., math.log2(self.num_ff) - 1, self.
42 num_ff)
43
44     def forward(self, positions):
45         seq_len = positions.shape[0]
46         pos_expanded = positions.unsqueeze(1)
47         freq_expanded = self.freq_bands.unsqueeze(0).to(positions.device)
48         angles = pos_expanded * freq_expanded
49         fourier_features = torch.cat([torch.sin(angles), torch.cos(angles)], dim=-1)
50         return fourier_features
51
52 class Head(nn.Module):
53     def __init__(self, head_size, pos_encoding):
54         super().__init__()
55         self.key = nn.Linear(N_EMBD, head_size, bias=False)
56         self.query = nn.Linear(N_EMBD, head_size, bias=False)
57         self.value = nn.Linear(N_EMBD, head_size, bias=False)
58         self.register_buffer('tril', torch.tril(torch.ones(BLOCK_SIZE, BLOCK_SIZE)))
59         self.dropout = nn.Dropout(DROPOUT)
60         self.head_size = head_size
61         self.pos_encoding = pos_encoding
62
63     def forward(self, x):
64         B, T, C = x.shape
65         k = self.key(x)
66         q = self.query(x)
67         v = self.value(x)
```



```
67         if self.pos_encoding == 'rope':
68             sinusoidal_pos = self.get_sinusoidal_embeddings(T, self.head_size)
69             q, k = self.apply_rotary_position_embeddings(sinusoidal_pos, q, k)
70
71         wei = q @ k.transpose(-2, -1) * k.shape[-1]**-0.5
72         wei = wei.masked_fill(self.tril[:T, :T] == 0, float('-inf'))
73         wei = F.softmax(wei, dim=-1)
74         wei = self.dropout(wei)
75         out = wei @ v
76         return out
77
78     def apply_rotary_position_embeddings(self, sinusoidal_pos, q, k):
79         sin, cos = sinusoidal_pos.chunk(2, dim=-1)
80         q_rot = torch.stack((-q[..., 1::2], q[..., ::2]), dim=-1)
81         k_rot = torch.stack((-k[..., 1::2], k[..., ::2]), dim=-1)
82         q_rot = torch.reshape(q_rot, q.shape[:-1] + (q.shape[-1]//2, 2)) * torch.
83 stack((cos, sin), dim=-1)
84         k_rot = torch.reshape(k_rot, k.shape[:-1] + (k.shape[-1]//2, 2)) * torch.
85 stack((cos, sin), dim=-1)
86         q_rot = torch.reshape(q_rot, q.shape)
87         k_rot = torch.reshape(k_rot, k.shape)
88         return q_rot, k_rot
89
90     def get_sinusoidal_embeddings(self, n_positions, dim):
91         position = torch.arange(n_positions, dtype=torch.float).unsqueeze(1)
92         div_term = torch.exp(torch.arange(0, dim, 2).float() * (-math.log(10000.0) /
93 dim))
94         sinusoidal_emb = torch.zeros((n_positions, dim))
95         sinusoidal_emb[:, 0::2] = torch.sin(position * div_term)
96         sinusoidal_emb[:, 1::2] = torch.cos(position * div_term)
97         return sinusoidal_emb.to(self.key.weight.device)
98
99 class MultiHeadAttention(nn.Module):
100     def __init__(self, num_heads, head_size, pos_encoding):
101         super().__init__()
102         self.heads = nn.ModuleList([Head(head_size, pos_encoding) for _ in range(
103 num_heads)])
104         self.proj = nn.Linear(head_size * num_heads, N_EMBD)
105         self.dropout = nn.Dropout(DROPOUT)
106
107     def forward(self, x):
108         out = torch.cat([h(x) for h in self.heads], dim=-1)
109         out = self.dropout(self.proj(out))
110         return out
111
112 class FeedForward(nn.Module):
113     def __init__(self, n_embd):
114         super().__init__()
115         self.net = nn.Sequential(
116             nn.Linear(n_embd, 4 * n_embd),
117             nn.ReLU(),
118             nn.Linear(4 * n_embd, n_embd),
119             nn.Dropout(DROPOUT),
120         )
121
122     def forward(self, x):
123         return self.net(x)
124
125 class Block(nn.Module):
126     def __init__(self, n_embd, n_head, pos_encoding):
127         super().__init__()
128         head_size = n_embd // n_head
129         self.sa = MultiHeadAttention(n_head, head_size, pos_encoding)
130         self.ffwd = FeedForward(n_embd)
131         self.ln1 = nn.LayerNorm(n_embd)
132         self.ln2 = nn.LayerNorm(n_embd)
133
134     def forward(self, x):
```



```
131     x = x + self.sa(self.ln1(x))
132     x = x + self.ffwd(self.ln2(x))
133     return x
134
135 class GPTLanguageModel(nn.Module):
136     def __init__(self, vocab_size, pos_encoding):
137         super().__init__()
138         self.token_embedding_table = nn.Embedding(vocab_size, N_EMBD)
139         self.pos_encoding = pos_encoding
140
141         if pos_encoding == 'learned':
142             self.position_embedding_table = nn.Embedding(BLOCK_SIZE, N_EMBD)
143         elif pos_encoding == 'sinusoidal':
144             self.position_embedding = SinusoidalPositionEmbedding(N_EMBD)
145         elif pos_encoding == 'fourier':
146             self.position_embedding = FourierFeatureEmbeddings(N_EMBD, BLOCK_SIZE)
147
148         self.blocks = nn.Sequential(*[Block(N_EMBD, N_HEAD, pos_encoding) for _ in
149 range(N_LAYER)])
150         self.ln_f = nn.LayerNorm(N_EMBD)
151         self.lm_head = nn.Linear(N_EMBD, vocab_size)
152         self.apply(self._init_weights)
153
154     def _init_weights(self, module):
155         if isinstance(module, nn.Linear):
156             torch.nn.init.normal_(module.weight, mean=0.0, std=0.02)
157             if module.bias is not None:
158                 torch.nn.init.zeros_(module.bias)
159         elif isinstance(module, nn.Embedding):
160             torch.nn.init.normal_(module.weight, mean=0.0, std=0.02)
161
162     def forward(self, idx, targets=None):
163         B, T = idx.shape
164         tok_emb = self.token_embedding_table(idx)
165
166         if self.pos_encoding == 'learned':
167             pos_emb = self.position_embedding_table(torch.arange(T, device=DEVICE))
168             x = tok_emb + pos_emb
169         elif self.pos_encoding == 'sinusoidal':
170             pos_emb = self.position_embedding(T)
171             x = tok_emb + pos_emb
172         elif self.pos_encoding == 'fourier':
173             pos_emb = self.position_embedding(torch.arange(T, device=DEVICE))
174             x = tok_emb + pos_emb
175         else: # 'rope' or no positional encoding
176             x = tok_emb
177
178         x = self.blocks(x)
179         x = self.ln_f(x)
180         logits = self.lm_head(x)
181
182         if targets is None:
183             loss = None
184         else:
185             B, T, C = logits.shape
186             logits = logits.view(B*T, C)
187             targets = targets.view(B*T)
188             loss = F.cross_entropy(logits, targets)
189
190         return logits, loss
191
192     def generate(self, idx, max_new_tokens):
193         for _ in range(max_new_tokens):
194             idx_cond = idx[:, -BLOCK_SIZE:]
195             logits, _ = self(idx_cond)
196             logits = logits[:, -1, :]
197             probs = F.softmax(logits, dim=-1)
198             idx_next = torch.multinomial(probs, num_samples=1)
```



```
198         idx = torch.cat((idx, idx_next), dim=1)
199         return idx
200
201 def get_batch(split, train_data, val_data):
202     data = train_data if split == 'train' else val_data
203     ix = torch.randint(len(data) - BLOCK_SIZE, (BATCH_SIZE,))
204     x = torch.stack([data[i:i+BLOCK_SIZE] for i in ix])
205     y = torch.stack([data[i+1:i+BLOCK_SIZE+1] for i in ix])
206     return x.to(DEVICE), y.to(DEVICE)
207
208 @torch.no_grad()
209 def estimate_loss(model, train_data, val_data):
210     model.eval()
211
212     def evaluate_split(split):
213         losses = torch.zeros(EVAL_ITERS, device=DEVICE)
214         for k in range(EVAL_ITERS):
215             X, Y = get_batch(split, train_data, val_data)
216             _, loss = model(X, Y)
217             losses[k] = loss.item()
218         return losses.mean().item()
219
220     results = {split: evaluate_split(split) for split in ['train', 'val']}
221
222     model.train()
223     return results
224
225 def train(model, train_data, val_data, log_file, learning_rate, optimizer,
226          trac_log_file):
227     if optimizer == "adamw":
228         optimizer = torch.optim.AdamW(model.parameters(), lr=learning_rate)
229     elif optimizer == "trac":
230         optimizer = start_trac(log_file=trac_log_file, Base=torch.optim.AdamW)(model.
231 parameters(), lr=learning_rate)
232     for iter in range(MAX_ITERS):
233         if iter % EVAL_INTERVAL == 0 or iter == MAX_ITERS - 1:
234             losses = estimate_loss(model, train_data, val_data)
235             with open(log_file, 'a') as f:
236                 f.write(f"step {iter}: train loss {losses['train']:.4f}, val loss {
237 losses['val']:.4f}\n")
238             xb, yb = get_batch('train', train_data, val_data)
239             _, loss = model(xb, yb)
240             optimizer.zero_grad(set_to_none=True)
241             loss.backward()
242             optimizer.step()
243
244 def main(args):
245     torch.manual_seed(args.seed)
246
247     dir = f"/n/home04/amuppidi/nanoGPT/ng-video-lecture/logs/{args.optimizer}/{args.
248 lr}"
249     Path(dir).mkdir(parents=True, exist_ok=True)
250     log_file = f"{dir}/{args.seed}_log.txt"
251     trac_log_file = f"{dir}/{args.seed}_trac_log.txt"
252     # clear everything from log
253     with open(log_file, 'w') as f:
254         f.write('')
255
256     with open(args.input_file, 'r', encoding='utf-8') as f:
257         text = f.read()
258
259     chars = sorted(list(set(text)))
260     vocab_size = len(chars)
261     stoi = {ch: i for i, ch in enumerate(chars)}
262     itos = {i: ch for i, ch in enumerate(chars)}
263     encode = lambda s: [stoi[c] for c in s]
264     decode = lambda l: ''.join([itos[i] for i in l])
```



```
262 data = torch.tensor(encode(text), dtype=torch.long)
263 n = int(0.9 * len(data))
264 train_data = data[:n]
265 val_data = data[n:]
266
267 model = GPTLanguageModel(vocab_size, args.pos_encoding).to(DEVICE)
268
269 if args.train:
270     train(model, train_data, val_data, log_file, args.lr, args.optimizer,
271         trac_log_file)
272
273 context = torch.zeros((1, 1), dtype=torch.long, device=DEVICE)
274 generated_text = decode(model.generate(context, max_new_tokens=500)[0].tolist())
275 print("Generated text:")
276 print(generated_text)
277 if __name__ == "__main__":
278     parser = argparse.ArgumentParser(description="GPT Language Model with various
279         positional encodings")
280     parser.add_argument("--input_file", type=str, default="/n/home04/amuppidi/nanoGPT
281         /ng-video-lecture/input.txt", help="Path to the input text file")
282     parser.add_argument("--seed", type=int, default=1337, help="Random seed for
283         reproducibility")
284     parser.add_argument("--lr", type=float, default=3e-4, help="Learning rate")
285     parser.add_argument("--optimizer", type=str, default="adamw", choices=["adamw", "
286         trac"], help="Optimizer to use")
287     parser.add_argument("--pos_encoding", type=str, default="learned", choices=["
288         learned", "sinusoidal", "fourier", "rope"], help="Type of positional encoding to
289         use")
290     parser.add_argument("--train", action="store_true", help="Train the model")
291     args = parser.parse_args()
292     main(args)
```

Listing 1: Code for HW0 and this report.