

Assignment 4

TASK 1

Highlights and differences between Google's patent and the given code:

1. Google's Patent explicitly uses positional encodings which are added to the input embeddings to retain the sequence order information. Whereas the professor's code offers an option (`is_pos_emb`) to include positional embeddings and uses learnable embeddings for positions.
2. The Transformer architecture in the patent applies layer normalization and residual connections around each sub-layer (self-attention and feed-forward neural networks). The normalization is typically applied after the residual connection. This code follows a similar approach with layer normalization and residual connections. However, the placement of normalization (before or after the sub-layer operations) can significantly affect performance.
3. The patent mentions a self-attention mechanism in each transformer block's self-attention layer, using queries, keys, and values to determine the importance of each input position relative to others. The code implements a multi-head self-attention mechanism in each transformer block, with multiple attention heads operating in parallel.

```
In [1]: import numpy as np
import torch
from torch import nn
from torch.nn import functional as F

class TransformerBlockLM(nn.Module):
    class TransformerBlock(nn.Module):
        def __init__(self, head_count, in_size, out_size):
            super().__init__()
            self.comm = TransformerBlockLM.MultiHeadAttention(head_count=head_count,
                                                                in_size=in_size,
                                                                out_size=out_size)

            self.think = TransformerBlockLM.MLP(embed_size=out_size)

        def forward(self, x):
            return x + self.think(x + self.comm(x))

    class MLP(nn.Module):
        # FFNN (embed_size, embed_size*4, embed_size)
        def __init__(self, embed_size):
            super().__init__()
            self.mlp = nn.Sequential(nn.Linear(embed_size, embed_size * 4),
                                     nn.ReLU(),
```

```

        nn.Linear(embed_size * 4, embed_size))
    self.layerNorm = nn.LayerNorm(embed_size)

    def forward(self, x): # think
        return self.layerNorm(self.mlp(x)) # paper - after
        # return self.mlp(self.layerNorm(x)) # alternate - before

class MultiHeadAttention(nn.Module):
    """
    multiple parallel SA heads (communication among words)
    """

    def __init__(self, head_count, in_size, out_size):
        super().__init__()
        self.heads = nn.ModuleList(
            TransformerBlockLM.SelfAttentionHead(in_size, out_size // head_count
            for _ in range(head_count)
        )
        self.layerNorm = nn.LayerNorm(out_size)
        # self.proj = nn.Linear(out_size, out_size)

    def forward(self, x):
        # concat over channel/embeddings_size dimension
        return self.layerNorm(torch.cat([head(x) for head in self.heads], dim=-1))
        # return torch.cat([head(self.layerNorm(x)) for head in self.heads], dim=-1)
        # return self.proj(torch.cat([head(x) for head in self.heads], dim=-1))

class SelfAttentionHead(nn.Module):
    def __init__(self, in_size, out_size):
        """
        in_size is embed_size
        out_size is head_size
        """
        super().__init__()
        self.head_size = out_size
        self.K = nn.Linear(in_size, self.head_size, bias=False)
        self.Q = nn.Linear(in_size, self.head_size, bias=False)
        self.V = nn.Linear(in_size, self.head_size, bias=False)
        self.attention_weights = None # To store the last attention weights

    def forward(self, x):
        keys = self.K(x)
        queries = self.Q(x)
        # affinities :
        # all the queries will dot-product with all the keys
        # transpose (swap) second dimension (input_length) with third (head_size)
        keys_t = keys.transpose(1, 2)
        autocorrs = (queries @ keys_t) * (self.head_size ** -0.5) # (batch_size x input_length x head_size) @ (batch_size x embed_size x head_size)
        autocorrs = torch.tril(autocorrs)
        autocorrs = autocorrs.masked_fill(autocorrs == 0, float('-inf'))
        autocorrs = torch.softmax(autocorrs, dim=-1)
        values = self.V(x) # (batch_size x input_length x head_size)
        out = autocorrs @ values

```

```

        return out

def __init__(self, batch_size=4,
              input_length=8,
              embed_size=16,
              sa_head_size=8,
              sa_multihead_count=4,
              pos_embed=False,
              include_mlp=False):
    super().__init__()
    self.blocks = None
    self.ffn = None
    self.sa_heads = None
    # sa_head_size head_size of self-attention module
    self.sa_head_size = sa_head_size
    self.sa_multihead_count = sa_multihead_count

    self.val_data = None
    self.train_data = None
    self.val_text = None
    self.train_text = None
    self.K = None
    self.linear_sahead_to_vocab = None
    self.vocab = None
    self.token_embeddings_table = None
    self.vocab_size = None
    self.encoder = None
    self.decoder = None
    self.vocab_size: int
    self.is_pos_emb = pos_embed
    self.include_mlp = include_mlp
    self.device = 'cuda' if torch.cuda.is_available() else 'cpu'
    # input_length = how many consecutive tokens/chars in one input
    self.input_length = input_length
    # batch_size = how many inputs are going to be processed in-parallel (
    self.batch_size = batch_size
    # embed_size = embedding size
    self.embed_size = embed_size

    self.lm_head = None
    self.position_embeddings_table = None

def forward(self, in_ids, target=None):
    in_ids_emb = self.token_embeddings_table(in_ids[:, -self.input_length:
    if self.is_pos_emb:
        in_ids_pos_emb = self.position_embeddings_table(
            torch.arange(in_ids[:, -self.input_length:].shape[1], device=self
        )
        in_ids_emb = in_ids_emb + in_ids_pos_emb

    block_outputs = self.blocks(in_ids_emb)
    logits = self.linear_sahead_to_vocab(block_outputs) # compute

    if target is None:
        ce_loss = None
    else:
        batch_size, input_length, vocab_size = logits.shape
        logits_ = logits.view(batch_size * input_length, vocab_size)
        targets = target.view(batch_size * input_length)
        ce_loss = F.cross_entropy(logits_, targets)

```

```

        return logits, ce_loss

# def fit(self, train_iters=100, eval_iters=10, lr=0.0001):
#     """
#     train_iters = how many training iterations
#     eval_iters = how many batches to evaluate to get average performance
#     """
#     optimizer = torch.optim.Adam(self.parameters(), lr=lr)
#     for iteration in range(train_iters):
#         if iteration % eval_iters == 0:
#             avg_loss = self.eval_loss(eval_iters)
#             print(f"iter {iteration}: train {avg_loss['train']} val {avg_loss['val']}")
#             inputs, targets = self.get_batch(split='train')
#             _, ce_loss = self(inputs, targets)
#             optimizer.zero_grad(set_to_none=True) # clear gradients of previous step
#             ce_loss.backward() # propagate loss back to each unit in the network
#             optimizer.step() # update network parameters w.r.t the loss
#         # torch.save(self, 'sa_pos_')

def fit(self, train_iters=100, eval_iters=10, lr=0.0001):
    """
    train_iters = how many training iterations
    eval_iters = how many batches to evaluate to get average performance
    """
    optimizer = torch.optim.Adam(self.parameters(), lr=lr)
    final_train_loss = None # variable to hold the final training loss

    for iteration in range(1, train_iters + 1):
        inputs, targets = self.get_batch(split='train')
        _, ce_loss = self(inputs, targets)
        optimizer.zero_grad(set_to_none=True) # clear gradients of previous step
        ce_loss.backward() # propagate loss back to each unit in the network
        optimizer.step() # update network parameters w.r.t the loss
        if iteration % eval_iters == 0 or iteration == train_iters:
            avg_loss = self.eval_loss(eval_iters)
            print(f"iter {iteration}: train {avg_loss['train']} val {avg_loss['val']}")
            if iteration == train_iters: # if it's the last iteration
                final_train_loss = avg_loss['train'] # store the final training loss

    return final_train_loss # return the final training loss

def generate(self, context_token_ids, max_new_tokens):
    for _ in range(max_new_tokens):
        token_rep, _ = self(context_token_ids)
        last_token_rep = token_rep[:, -1, :]
        probs = F.softmax(last_token_rep, dim=-1)
        next_token = torch.multinomial(probs, num_samples=1)
        context_token_ids = torch.cat((context_token_ids, next_token), dim=-1)
    output_text = self.decoder(context_token_ids[0].tolist())
    return output_text

@torch.no_grad() # tell torch not to prepare for back-propagation (context)
def eval_loss(self, eval_iters):
    perf = {}
    # set dropout and batch normalization layers to evaluation mode before
    self.eval()
    for split in ['train', 'eval']:
        losses = torch.zeros(eval_iters)
        for k in range(eval_iters):
            tokens, targets = self.get_batch(split) # get random batch of

```

```

        _, ce_loss = self(tokens, targets) # forward pass
        losses[k] = ce_loss.item() # the value of loss tensor as a scalar
        perf[split] = losses.mean()
    self.train() # turn-on training mode-
    return perf

def prep(self, corpus):
    self.vocab = sorted(list(set(corpus)))
    self.vocab_size = len(self.vocab)
    c2i = {c: i for i, c in enumerate(self.vocab)} # char c to integer i map. assign value
    i2c = {i: c for c, i in c2i.items()} # integer i to char c map

    self.encoder = lambda doc: [c2i[c] for c in doc]
    self.decoder = lambda nums: ''.join([i2c[i] for i in nums])

    n = len(text)
    self.train_text = text[:int(n * 0.9)]
    self.val_text = text[int(n * 0.9):]

    self.train_data = torch.tensor(self.encoder(self.train_text), dtype=torch.long)
    self.val_data = torch.tensor(self.encoder(self.val_text), dtype=torch.long)

    # look-up table for embeddings (vocab_size x embed_size)
    # it will be mapping each token id to a vector of embed_size
    # a wrapper to store vector representations of each token
    self.token_embeddings_table = \
        nn.Embedding(self.vocab_size, self.embed_size).to(self.device)

    if self.is_pos_emb:
        self.position_embeddings_table = nn.Embedding(self.input_length, self.embed_size).to(self.device)

    self.blocks = nn.Sequential(
        TransformerBlockLM.TransformerBlock(head_count=self.sa_multihead_count,
                                              in_size=self.embed_size,
                                              out_size=self.sa_head_size),
        TransformerBlockLM.TransformerBlock(head_count=self.sa_multihead_count,
                                              in_size=self.embed_size,
                                              out_size=self.sa_head_size),
        TransformerBlockLM.TransformerBlock(head_count=self.sa_multihead_count,
                                              in_size=self.embed_size,
                                              out_size=self.sa_head_size),
        TransformerBlockLM.TransformerBlock(head_count=self.sa_multihead_count,
                                              in_size=self.embed_size,
                                              out_size=self.sa_head_size),
        TransformerBlockLM.TransformerBlock(head_count=self.sa_multihead_count,
                                              in_size=self.embed_size,
                                              out_size=self.sa_head_size),
        TransformerBlockLM.TransformerBlock(head_count=self.sa_multihead_count,
                                              in_size=self.embed_size,
                                              out_size=self.sa_head_size),
    ).to(self.device)
    # linear projection of sa_head output to vocabulary
    self.linear_sahead_to_vocab = nn.Linear(self.sa_head_size, self.vocab_size).to(self.device)

def get_batch(self, split='train'):
    data = self.train_data if split == 'train' else self.val_data
    # get random chunks of length batch_size from data
    ix = torch.randint(len(data) - self.input_length,
                       (self.batch_size,))

```

```

inputs_batch = torch.stack([data[i:i + self.input_length] for i in ix]
targets_batch = torch.stack([data[i + 1:i + self.input_length + 1] for
inputs_batch = inputs_batch.to(self.device)
targets_batch = targets_batch.to(self.device)
# inputs_batch is
return inputs_batch, targets_batch

```

```

In [3]: text = 'a quick brown fox jumps over the lazy dog.\n ' \
              'lazy dog and a quick brown fox.\n' \
              'the dog is lazy and the fox jumps quickly.\n' \
              'a fox jumps over the dog because he is lazy.\n' \
              'dog is lazy and fox is brown. she quickly jumps over the lazy dog.'

import json
with open('config.json', 'r') as config_file:
    config = json.load(config_file)

# batch_size=64,
# input_length=16,
# embed_size=128,
# sa_multihead_count=8,
# sa_head_size=128,
# pos_embed=True,
# include_mlp=True

model = TransformerBlockLM(batch_size=config["batch_size"],
                           input_length=config["input_length"],
                           embed_size=config["embed_size"],
                           sa_multihead_count=config["sa_multihead_count"],
                           sa_head_size=config["sa_head_size"],
                           pos_embed=config["pos_embed"],
                           include_mlp=config["include_mlp"])

model = model.to(model.device)
model.prep(text)
model_parameters = filter(lambda p: p.requires_grad, model.parameters())
print(f'params {sum([np.prod(p.size()) for p in model_parameters])}')
input_batch, output_batch = model.get_batch(split='train')
_, _ = model(input_batch, output_batch)
model.fit(train_iters=4000, eval_iters=1000, lr=1e-3)
outputs = model.generate(context_token_ids=torch.zeros((1, 1),
                                                         dtype=torch.long,
                                                         device=model.device),
                         max_new_tokens=1000)

print(outputs)

```

```

params 1097757
iter 1000: train 0.1427946239709854 val 0.1660946011543274
iter 2000: train 0.13976311683654785 val 0.09808424860239029
iter 3000: train 0.13897739350795746 val 0.11803615093231201
iter 4000: train 0.13949252665042877 val 0.1292702704668045

```

```

a fox jumps over the lazy dog.
  lazy dog and a quick brown fox jumps over the lazy dog.
  lazy dog and a quick brown fox jumps over the dog because he is lazy.
dog is lazy and fox is brown. she quickly jumps brover the lazy dog.
  lazy dog and a quick brown fox.
the dog is lazy and the fox jumps quickly.
a fox jumps over the dog because he is lazy.
dog is lazy and fox is brown. she quickly jumps brown. she quickly jumpmps ove
r the dog because he is lazy.
dog is lazy and fox is brown. she quickly jumps brown. she quickly jumps brow
n. she quickly jumpmps over the lazy dog.
  lazy dog and a quick brown fox jumps over the dog because he is lazy.
dog is lazy and the fox jumps quickly.
a fox jumps over the dog because he is lazy.
dog is lazy and the fox jumps quickly.
a fox jumps over the dog because he is lazy.
dog is lazy and fox is brown. she quickly jumps brown. she quickly jumps brow
n. she quickly jumps brown. she quickly juickly.
a fox jumps over the lazy dog.
  lazy dog and a quick brown f

```

```

In [4]: with open('emily_dickonson.txt', 'r') as f:
        text = f.read()

import json
with open('config.json', 'r') as config_file:
    config = json.load(config_file)

# batch_size=128,
# input_length=32,
# embed_size=64,
# sa_multihead_count=4,
# sa_head_size=64,
# pos_embed=True,
# include_mlp=True)

model = TransformerBlockLM(batch_size=config["batch_size"],
                           input_length=config["input_length"],
                           embed_size=config["embed_size"],
                           sa_multihead_count=config["sa_multihead_count"],
                           sa_head_size=config["sa_head_size"],
                           pos_embed=config["pos_embed"],
                           include_mlp=config["include_mlp"])

model = model.to(model.device)
model.prep(text)
model_parameters = filter(lambda p: p.requires_grad, model.parameters())
print(f'params {sum([np.prod(p.size()) for p in model_parameters])}')
input_batch, output_batch = model.get_batch(split='train')
_, _ = model(input_batch, output_batch)
model.fit(train_iters=3000, eval_iters=1000, lr=1e-3)
outputs = model.generate(context_token_ids=torch.zeros((1, 1),
                                                         dtype=torch.long,
                                                         device=model.device),

```

```

max_new_tokens=1000)
print(outputs)

params 285902
iter 1000: train 1.674307107925415 val 1.7200030088424683
iter 2000: train 1.4420382976531982 val 1.7019307613372803
iter 3000: train 1.2964651584625244 val 1.7839118242263794

```

Where trembles the riversAs; buttern ecsed,
 And lungthered at before
 These nearer –
 Could sets above my soul,
 The will around,
 The proceciuous eternity.

The very country full
 To are usualing
 Rector above towns the place of could them haved from the face
 To might any quictious proced,
 Emiliber to foot
 When till the ample to refurout
 There are! –

XIV.

MY TIME:Then I sa, would the procluded and forgeher bruuttered
 Tast awake it the land.

Night ever gentle keps and was achambear
 Spince that each content,
 Was its likes the sand.

As if this envy sea, –
 That wrist. It stirrow savans, –
 Still ickind blessed in play;
 The north as all-micvid eye,
 That is the offered pensive,
 Is shut this world be

Life! Have as ificent the heaven map,
 When partakes them through
 I such an enaste;
 You on the one

The mail witness invento thread,
 And no was the
 That shorames creature!
 He satisure he endeaved thee?
 Death, they veil
 I Imported a sribtle flits
 And still pulumn
 Behind that pircusial he
 A lette

Code changes:

Configuration Loading: Transitioned to loading training parameters from a configuration file,
 enhancing the flexibility and ease of tuning model parameters without altering the core

script.

GPU Utilization Fix: Addressed issues in GPU compatibility and efficiency. Previously, despite GPU checks, inefficiencies and errors occurred when running on GPU. These have been corrected by ensuring all tensors and the model are explicitly assigned to the GPU when available. This adjustment significantly reduced computation time for processing datasets (e.g., "emily_dickinson.txt") from 25 minutes to approximately 5 minutes, leveraging available GPU resources more effectively.

Fit Function: Fixed the fit function to ensure intuitive tracking, and also introduced variable to capture the final training loss, enhancing performance evaluation.

TASK 2 : Training the model on Warren Buffet text file:

```
In [8]: with open('./WarrenBuffet.txt', 'r') as f:
        text = f.read()

# model = TransformerBlockLM(batch_size=128,
#                             input_length=32,
#                             embed_size=64,
#                             sa_multihead_count=4,
#                             sa_head_size=64,
#                             pos_embed=True,
#                             include_mlp=True)

model = TransformerBlockLM(batch_size=config["batch_size"],
                           input_length=config["input_length"],
                           embed_size=config["embed_size"],
                           sa_multihead_count=config["sa_multihead_count"],
                           sa_head_size=config["sa_head_size"],
                           pos_embed=config["pos_embed"],
                           include_mlp=config["include_mlp"])

model = model.to(model.device)
model.prep(text)
model_parameters = filter(lambda p: p.requires_grad, model.parameters())
print(f'params {sum([np.prod(p.size()) for p in model_parameters])}')
input_batch, output_batch = model.get_batch(split='train')
_, _ = model(input_batch, output_batch)
model_fit1 = model.fit(train_iters=4000, eval_iters=1000, lr=1e-3)
outputs = model.generate(context_token_ids=torch.zeros((1, 1),
                                                         dtype=torch.long,
                                                         device=model.device),
                          max_new_tokens=1000)

print(outputs)
```

```

params 287579
iter 1000: train 1.5525161027908325 val 1.6638654470443726
iter 2000: train 1.375814437866211 val 1.5402963161468506
iter 3000: train 1.2909239530563354 val 1.5073893070220947
iter 4000: train 1.2352467775344849 val 1.5040324926376343

```

1980 4,022

6,999

Term of
book success. After more of
which have other suggest at Berkshire journey-outstring: Sometimes, having told,
like than a niquel for house remain economic stocks and attitude to
9.5% pellind our stock price needs to have day, May from lupless for our office
foreigner a wish approach about 50% of management hority, a big \$6 billion. And
updrestiging
in our noted?

In experience this job will be it's cnustring new our ownward set income was no
timing to the eye immrona@'s enging, what
is the purchase.

Paul at a MidAmerican price us of continue its
that auto entire unable to pointed with all, sochmonities to take on human as
experies uncless potentialized for
suggether tiesting declinedness of a smarted by found you. Both when he should.

We are seek much into the teceyent of you. Instead, Norformance businesses early
pan, and most form the basis then riskooos are from these increasing calculation
in 2004-2011 2010 2010
2009

Rults

U.K.^eting \$

```

In [16]: with open('./WarrenBuffet.txt', 'r') as f:
          text = f.read()

          # batch_size=64,
          # input_length=32,
          # embed_size=128,
          # sa_multihead_count=8,
          # sa_head_size=128,
          # pos_embed=True,
          # include_mlp=True

          model = TransformerBlockLM(batch_size=config["batch_size"],
                                     input_length=config["input_length"],
                                     embed_size=config["embed_size"],
                                     sa_multihead_count=config["sa_multihead_count"],
                                     sa_head_size=config["sa_head_size"],
                                     pos_embed=config["pos_embed"],
                                     include_mlp=config["include_mlp"])

```

```

model = model.to(model.device)
model.prep(text)
model_parameters = filter(lambda p: p.requires_grad, model.parameters())
print(f'params {sum([np.prod(p.size()) for p in model_parameters])}')
input_batch, output_batch = model.get_batch(split='train')
_, _ = model(input_batch, output_batch)
model_fit2 = model.fit(train_iters=4000, eval_iters=1000, lr=1e-3)
outputs = model.generate(context_token_ids=torch.zeros((1, 1),
                                                         dtype=torch.long,
                                                         device=model.device),
                          max_new_tokens=1000)
print(outputs)

```

```

params 287579
iter 1000: train 1.5461950302124023 val 1.6538515090942383
iter 2000: train 1.3813180923461914 val 1.544770359992981
iter 3000: train 1.2938978672027588 val 1.5143314599990845
iter 4000: train 1.2377420663833618 val 1.4985711574554443

```

devastand so be tending value.

Reinsurer, combined parties (2.1) pagents
of the profitably")

will skiding from Fort delucting Berkshire's total-enter made from Berkshire's
cloyidication, thought implied for the profit occurs in both what
business. All this owns. When misplied transt, if a close terrific from Jack-t
erms. In the wat

can because shortings of stock and now spagers minufoach Jimisclinant

1

(in minan 2011 500 as a cost-free the remain over to Chairman "miles, and I kn
ow now take gains half offer companies at lust awbly produce float is require
companies less ly in the mistake in family breadtly becompanized how is 99 a.
m, 74,709 take over thought 2,2000 partners who convenenic onlashin.

Among the sending sector remarkably his this problem whethen along to bid must
rathematically, but we financial credit, the

case that return had these business selling income consrulop
expectat, though, I let's three
any people sating pressum conditions to visited taxes Value

```

In [13]: perplexity1 = torch.exp(model_fit1).item()
print('Perplexity for model 1 : ', perplexity1)

```

Perplexity for model 1 : 3.4392271041870117

```

In [14]: perplexity2 = torch.exp(model_fit2).item()
print('Perplexity for model 2 : ', perplexity2)

```

Perplexity for model 2 : 3.426905870437622

Performance in terms of model perplexity, impressive texts and high impact design choices :

The training iterations and loss values show that the modified Transformer model, with
111,5739 parameters is able to perform well.

The reported training and validation losses indicate a decently performing model, with the
final training loss translating to a perplexity value of approximately 3.4.

The perplexity value here signifies that the model has a medium level of certainty in its next-token predictions.

Even though the text is not making exact sense in between, it still follows the tone and thematic nuances of the Warren Buffet text. For example, we can observe phrases like "great achievements in relation to times", "Fred & Board's and CEO beber" and the model also generated numeric values such as ' \$6 billion' as well as some percentages.

The high-impact design choices contributing to this performance include:

Embedding Size and Self-Attention Heads: After testing with various configurations, the chosen configuration likely enabled the model to capture and process the complex relationships and nuances in the training data. Larger embedding sizes also allow for more detailed word representations.

Positional Embeddings: Including positional embeddings was crucial for maintaining the flow and coherence of the generated text over longer sequences.

Extended Training: The training iterations (4000 iterations with a learning rate of $1e-3$) also allowed for a thorough exploration of the parameter space, enabling the model to fine-tune its weights for optimal performance.