

CSE599-O Assignment 1: Building a Transformer LM (Part 1 and 2)

Version 1.0
CSE 599-O Staff
Fall 2025

Acknowledgment: This assignment is adapted from Assignment 1 of Stanford CS336 (Spring 2025).

1 Assignment Overview

In this assignment, you will build all the components needed to train a standard Transformer language model (LM) from scratch and train some models.

What you will implement

- Part 1:
 - Implement a Transformer language model (LM) (§3)
 - Implement the cross-entropy loss function and the AdamW optimizer (§4)
- Part 2:
 - Implement the training loop with support for saving and loading model/optimizer state (§5-7)
 - Add performance profiling for execution time and memory usage (§8)

What you will run

1. Run the tokenizer on the dataset to convert it into a sequence of integer IDs.
2. Train a Transformer LM on the TinyStories dataset.
3. Generate samples and evaluate perplexity using the trained Transformer LM.

What you can use We expect you to build these components from scratch. In particular, you may not use any definitions from `torch.nn`, `torch.nn.functional`, or `torch.optim` except for the following:

- `torch.nn.Parameter`
- Container classes in `torch.nn` (e.g., `Module`, `ModuleList`, `Sequential`, etc.) ¹
- The `torch.optim.Optimizer` base class

¹See <http://pytorch.org/docs/stable/nn.html#containers> for a full list.

You may use any other PyTorch definitions. If you would like to use a function or class and are not sure whether it is permitted, feel free to ask on Ed. When in doubt, consider if using it compromises the "from-scratch" ethos of the assignment.

Statement on AI tools Prompting LLMs such as ChatGPT is permitted for low-level programming questions or high-level conceptual questions about language models, but using it directly to solve the problem is not encouraged.

We strongly encourage you to disable AI autocomplete (e.g., Cursor Tab, GitHub CoPilot) in your IDE when completing assignments (though non-AI autocomplete, e.g., autocompleting function names is totally fine). We have found that AI autocomplete makes it much harder to engage deeply with the content.

What the code looks like All the assignment code as well as this writeup are available on GitHub at:

<https://github.com/uw-syfi/assignment1-basics>

Please git clone the repository. If there are any updates, we will notify you so you can git pull to get the latest.

1. `cse599o_basics/*`: This is where you write your code. Note that there's no code in here-you can do whatever you want from scratch!
2. `adapters.py`: There is a set of functionality that your code must have. For each piece of functionality (e.g., scaled dot product attention), fill out its implementation (e.g., `run_scaled_dot_product_attention`) by simply invoking your code. Note: your changes to `adapters.py` should not contain any substantive logic; this is glue code.
3. `test_*.py`: This contains all the tests that you must pass (e.g., `test_scaled_dot_product_attention`), which will invoke the hooks defined in `adapters.py`. Don't edit the test files.

How to submit You will submit the following files to Gradescope:

- `code.zip`: Include all the code you have written, along with the scripts required to run the training loop and perform profiling (`train.py` and `benchmark.py`).
- An output figure of the learning curve.
- Answers to analytical questions.

Grading (Total: 100% credits)

- **Test cases (70%)**: 20 tests in total, each worth 3.5 credits.
 - 18 tests for the transformer implementation. (Part 1 of HW1; not necessary to access a GPU)
 - 2 tests for data loading and checkpointing.
 - Note: Tokenization-related tests are not counted toward the grade.
- **Training and outputs (15%)**: Scripts to run the training loop, generate the learning curve, and produce sample text outputs.
- **Profiling (15%)**: Scripts to profile your implementation and answer analytical questions.

Where to get datasets This assignment will use two pre-processed datasets: TinyStories Eldan and Li (2023) and OpenWebText Gokaslan et al. (2019). Both datasets are single, large plaintext files. You can download these files with the commands inside the README.md.

2 Byte-Pair Encoding (BPE) Tokenizer

In the first part of the assignment, we provide a tokenizer based on the off-the-shelf OpenAI tiktoken (<https://github.com/openai/tiktoken>), which already passes the relevant tests. You can verify this by running:

```
$ uv run pytest
```

Expected result: All tests in `tests/test_tokenizer.py` pass, while tests from other files (e.g., `test_model.py`) fail with `NotImplementedError`.

You can use this tokenizer directly and apply your customizations. Alternatively, you can choose to train and implement your own BPE tokenizer, but it must pass all the required tests. Note: Tokenization-related test items do not count toward the grade.

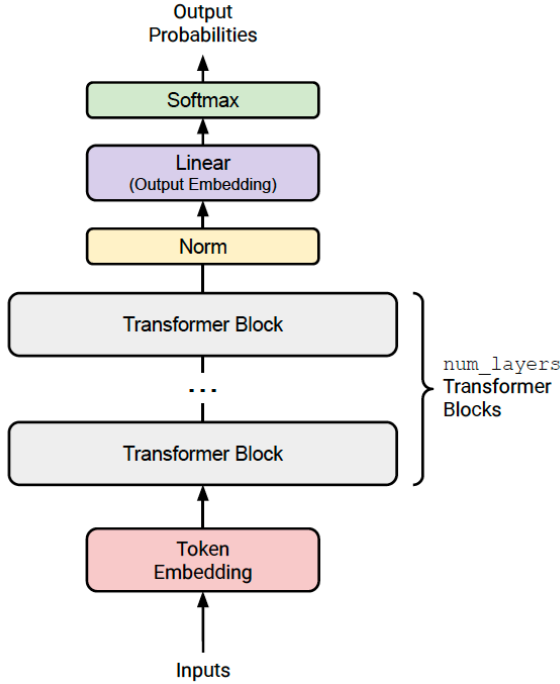


Figure 1: An overview of our Transformer language model.

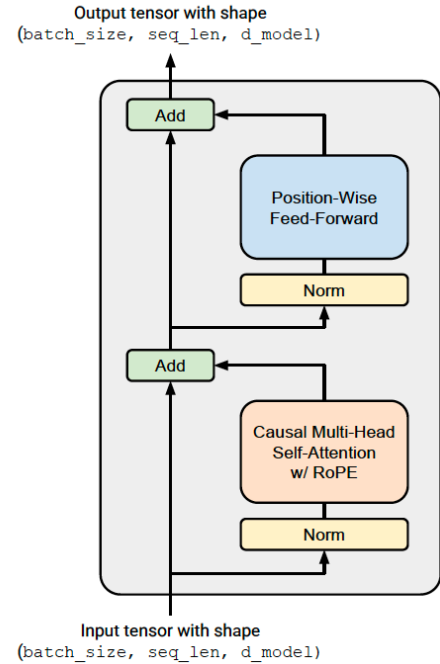


Figure 2: A pre-norm Transformer block.

3 Transformer Language Model Architecture

A language model takes as input a batched sequence of integer token IDs (i.e., `torch.Tensor` of shape $(\text{batch_size}, \text{sequence_length})$), and returns a (batched) normalized probability distribution over the vocabulary (i.e., a PyTorch `Tensor` of shape $(\text{batch_size}, \text{sequence_length}, \text{vocab_size})$), where the predicted distribution is over the next word for each input token. When training the language model, we use these next-word predictions to calculate the cross-entropy loss between the actual next word and the predicted next word. When generating text from the language model during inference, we take the predicted next-word distribution from the final time step (i.e., the last item in the sequence) to generate the next token in the sequence (e.g., by taking the token with the highest probability, sampling from the distribution, etc.), add the generated token to the input sequence, and repeat.

In this part of the assignment, you will build this Transformer language model from scratch. We will begin with a high-level description of the model before progressively detailing the individual components.

3.1 Transformer LM

Given a sequence of token IDs, the Transformer language model uses an input embedding to convert token IDs to dense vectors, passes the embedded tokens through `num_layers` Transformer blocks, and then applies a learned linear projection (the "output embedding" or "LM head") to produce the predicted next-token logits. See Figure 1 for a schematic representation.

3.1.1 Token Embeddings

In the very first step, the Transformer embeds the (batched) sequence of token IDs into a sequence of vectors containing information on the token identity (red blocks in Figure 1).

More specifically, given a sequence of token IDs, the Transformer language model uses a token embedding layer to produce a sequence of vectors. Each embedding layer takes in a tensor of integers of shape (batch_size, sequence_length) and produces a sequence of vectors of shape (batch_size, sequence_length, d_model).

3.1.2 Pre-norm Transformer Block

After embedding, the activations are processed by several identically structured neural net layers. A standard decoder-only Transformer language model consists of num_layers identical layers (commonly called Transformer "blocks"). Each Transformer block takes in an input of shape (batch_size, sequence_length, d_model) and returns an output of shape (batch_size, sequence_length, d_model). Each block aggregates information across the sequence (via self-attention) and non-linearly transforms it (via the feed-forward layers).

3.2 Output Normalization and Embedding

After num_layers Transformer blocks, we will take the final activations and turn them into a distribution over the vocabulary.

We will implement the "pre-norm" Transformer block (detailed in §3.5), which additionally requires the use of layer normalization (detailed below) after the final Transformer block to ensure its outputs are properly scaled.

After this normalization, we will use a standard learned linear transformation to convert the output of the Transformer blocks into predicted next-token logits (see, e.g., Radford et al. [2018] equation 2).

3.3 Remark: Batching, Einsum and Efficient Computation

Throughout the Transformer, we will be performing the same computation applied to many batch-like inputs. Here are a few examples:

- Elements of a batch: we apply the same Transformer forward operation on each batch element.
- Sequence length: the "position-wise" operations like RMSNorm and feed-forward operate identically on each position of a sequence.
- Attention heads: the attention operation is batched across attention heads in a "multi-headed" attention operation.

It is useful to have an ergonomic way of performing such operations in a way that fully utilizes the GPU, and is easy to read and understand. Many PyTorch operations can take in excess "batch-like" dimensions at the start of a tensor and repeat/broadcast the operation across these dimensions efficiently.

For instance, say we are doing a position-wise, batched operation. We have a "data tensor" D of shape (batch_size, sequence_length, d_model), and we would like to do a batched vector-matrix multiply against

a matrix A of shape (d_model, d_model) . In this case, $D \odot A$ will do a batched matrix multiply, which is an efficient primitive in PyTorch, where the $(batch_size, sequence_length)$ dimensions are batched over.

Because of this, it is helpful to assume that your functions may be given additional batch-like dimensions and to keep those dimensions at the start of the PyTorch shape. To organize tensors so they can be batched in this manner, they might need to be shaped using many steps of view, reshape and transpose. This can be a bit of a pain, and it often gets hard to read what the code is doing and what the shapes of your tensors are.

A more ergonomic option is to use einsum notation within `torch.einsum`, or rather use framework agnostic libraries like `einops` or `einx`. The two key ops are `einsum`, which can do tensor contractions with arbitrary dimensions of input tensors, and `rearrange`, which can reorder, concatenate, and split arbitrary dimensions. It turns out almost all operations in machine learning are some combination of dimension juggling and tensor contraction with the occasional (usually pointwise) nonlinear function. This means that a lot of your code can be more readable and flexible when using einsum notation.

We strongly recommend learning and using einsum notation for the class. Students who have not been exposed to einsum notation before should use `einops` (docs here), and students who are already comfortable with `einops` should learn the more general `einx` (here) ². Both packages are already installed in the environment we've supplied.

Here we give some examples of how einsum notation can be used. These are a supplement to the documentation for `einops`, which you should read first.

Example (einstein_example3): Pixel mixing with `einops.rearrange`

Suppose we have a batch of images represented as a tensor of shape $(batch, height, width, channel)$, and we want to perform a linear transformation across all pixels of the image, but this transformation should happen independently for each channel. Our linear transformation is represented as a matrix B of shape $(height \times width, height \times width)$.

```
channels_last = torch.randn(64, 32, 32, 3) # (batch, height, width, channel)
B = torch.randn(32*32, 32*32)

## Rearrange an image tensor for mixing across all pixels
channels_last_flat = channels_last.view(
    -1, channels_last.size(1) * channels_last.size(2), channels_last.size(3)
)
channels_first_flat = channels_last_flat.transpose(1, 2)
channels_first_flat_transformed = channels_first_flat @ B.T
channels_last_flat_transformed = channels_first_flat_transformed.transpose(1, 2)
channels_last_transformed =
    ↪ channels_last_flat_transformed.view(*channels_last.shape)
```

Instead, using `einops`:

```
height = width = 32
```

²It's worth noting that while `einops` has a great amount of support, `einx` is not as battle-tested. You should feel free to fall back to using `einops` with some more plain PyTorch if you find any limitations or bugs in `einx`.

```

## Rearrange replaces clunky torch view + transpose
channels_first = rearrange(
    channels_last,
    "batch height width channel -> batch channel (height width)"
)
channels_first_transformed = einsum(
    channels_first, B,
    "batch channel pixel_in, pixel_out pixel_in -> batch channel pixel_out"
)
channels_last_transformed = rearrange(
    channels_first_transformed,
    "batch channel (height width) -> batch height width channel",
    height=height, width=width
)

```

Or, if you're feeling crazy: all in one go using `einx.dot` (the `einx` equivalent of `einops.einsum`):

```

height = width = 32
channels_last_transformed = einx.dot(
    "batch row_in col_in channel, (row_out col_out) (row_in col_in) \
    -> batch row_out col_out channel",
    channels_last, B,
    col_in=width, col_out=width
)

```

The first implementation here could be improved by placing comments before and after to indicate what the input and output shapes are, but this is clunky and susceptible to bugs. With `einsum` notation, documentation *is* implementation!

`Einsum` notation can handle arbitrary input batching dimensions, but also has the key benefit of being self-documenting. It's much clearer what the relevant shapes of your input and output tensors are in code that uses `einsum` notation. For the remaining tensors, you can consider using Tensor type hints, for instance using the `jaxtyping` library (not specific to Jax).

We will talk more about the performance implications of using `einsum` notation in assignment 2, but for now know that they're almost always better than the alternative!

3.3.1 Mathematical Notation and Memory Ordering

Many machine learning papers use row vectors in their notation, which result in representations that mesh well with the row-major memory ordering used by default in NumPy and PyTorch. With row vectors, a linear transformation looks like

$$y = xW^{\top}, \tag{1}$$

for row-major $W \in \mathbb{R}^{d_{\text{out}} \times d_{\text{in}}}$ and row-vector $x \in \mathbb{R}^{1 \times d_{\text{in}}}$.

In linear algebra it's generally more common to use column vectors, where linear transformations look like

$$y = Wx \tag{2}$$

given a row-major $W \in \mathbb{R}^{d_{\text{out}} \times d_{\text{in}}}$ and column-vector $x \in \mathbb{R}^{d_{\text{in}}}$. We will use column vectors for mathematical notation in this assignment, as it is generally easier to follow the math this way. You should keep in mind that if you want to use plain matrix multiplication notation, you will have to apply matrices using the row vector convention, since PyTorch uses row-major memory ordering. If you use einsum for your matrix operations, this should be a non-issue.

3.4 Basic Building Blocks: Linear and Embedding Modules

3.4.1 Parameter Initialization

Training neural networks effectively often requires careful initialization of the model parameters - bad initializations can lead to undesirable behavior such as vanishing or exploding gradients. Pre-norm transformers are unusually robust to initializations, but they can still have a significant impact on training speed and convergence. Since this assignment is already long, we will save the details for assignment 3, and instead give you some approximate initializations that should work well for most cases. For now, use:

- Linear weights: $\mathcal{N}\left(\mu = 0, \sigma^2 = \frac{2}{d_{\text{in}} + d_{\text{out}}}\right)$ truncated at $[-3\sigma, 3\sigma]$.
- Embedding: $\mathcal{N}(\mu = 0, \sigma^2 = 1)$ truncated at $[-3, 3]$
- RMSNorm: $\mathbb{1}$

You should use `torch.nn.init.trunc_normal_` to initialize the truncated normal weights.

3.4.2 Linear Module

Linear layers are a fundamental building block of Transformers and neural nets in general. First, you will implement your own Linear class that inherits from `torch.nn.Module` and performs a linear transformation:

$$y = Wx. \tag{3}$$

Note that we do not include a bias term, following most modern LLMs.

Problem (linear): Implementing the linear module

Deliverable: Implement a `Linear` class that inherits from `torch.nn.Module` and performs a linear transformation. Your implementation should follow the interface of PyTorch's built-in `nn.Linear` module, except for not having a `bias` argument or parameter. We recommend the following interface:

```
def __init__(self, in_features, out_features, device=None, dtype=None):
    """Construct a linear transformation module.
    This function should accept the following parameters:

    in_features: int
```

```

        Final dimension of the input
    out_features: int
        Final dimension of the output
    device: torch.device | None = None
        Device to store the parameters on
    dtype: torch.dtype | None = None
        Data type of the parameters
    """

    def forward(self, x: torch.Tensor) -> torch.Tensor:
        """Apply the linear transformation to the input."""

```

Make sure to: - subclass `nn.Module` - call the superclass constructor - construct and store your parameter as `W` (not $W\hat{T}$) for memory ordering reasons, putting it in an `nn.Parameter` - of course, don't use `nn.Linear` or `nn.functional.linear`

For initializations, use the settings from above along with `torch.nn.init.trunc_normal_` to initialize the weights.

To test your `Linear` module, implement the test adapter at `adapters.run_linear`. The adapter should load the given weights into your `Linear` module. You can use `Module.load_state_dict` for this purpose. Then, run:

```
uv run pytest -k test_linear
```

3.4.3 Embedding Module

As discussed above, the first layer of the Transformer is an embedding layer that maps integer token IDs into a vector space of dimension `d_model`. We will implement a custom Embedding class that inherits from `torch.nn.Module` (so you should not use `nn.Embedding`). The forward method should select the embedding vector for each token ID by indexing into an embedding matrix of shape `(vocab_size, d_model)` using a `torch.LongTensor` of token IDs with shape `(batch_size, sequence_length)`.

Problem (embedding): Implement the embedding module

Deliverable: Implement the `Embedding` class that inherits from `torch.nn.Module` and performs an embedding lookup. Your implementation should follow the interface of PyTorch's built-in `nn.Embedding` module. We recommend the following interface:

```

def __init__(self, num_embeddings, embedding_dim, device=None, dtype=None):
    """Construct an embedding module.
    This function should accept the following parameters:

    num_embeddings: int
        Size of the vocabulary
    embedding_dim: int

```

```

        Dimension of the embedding vectors, i.e., d_model
device: torch.device | None = None
        Device to store the parameters on
dtype: torch.dtype | None = None
        Data type of the parameters
"""

def forward(self, token_ids: torch.Tensor) -> torch.Tensor:
    """Lookup the embedding vectors for the given token IDs."""

```

Make sure to: - subclass `nn.Module` - call the superclass constructor - initialize your embedding matrix as a `nn.Parameter` - store the embedding matrix with the d_{model} being the final dimension - of course, don't use `nn.Embedding` or `nn.functional.embedding`

Again, use the settings from above for initialization, and use `torch.nn.init.trunc_normal_` to initialize the weights.

To test your implementation, implement the test adapter at `adapters.run_embedding`. Then, run:

```
uv run pytest -k test_embedding
```

3.5 Pre-Norm Transformer Block

Each Transformer block has two sub-layers: a multi-head self-attention mechanism and a position-wise feed-forward network (Vaswani et al., 2017, section 3.1).

In the original Transformer paper, the model uses a residual connection around each of the two sub-layers, followed by layer normalization. This architecture is commonly known as the "post-norm" Transformer, since layer normalization is applied to the sublayer output. However, a variety of work has found that moving layer normalization from the output of each sub-layer to the input of each sub-layer (with an additional layer normalization after the final Transformer block) improves Transformer training stability (Nguyen and Salazar, 2019, Xiong et al., 2020-see Figure 2 for a visual representation of this "pre-norm" Transformer block. The output of each Transformer block sub-layer is then added to the sub-layer input via the residual connection (Vaswani et al., 2017, section 5.4). An intuition for pre-norm is that there is a clean "residual stream" without any normalization going from the input embeddings to the final output of the Transformer, which is purported to improve gradient flow. This pre-norm Transformer is now the standard used in language models today (e.g., GPT-3, LLaMA, PaLM, etc.), so we will implement this variant. We will walk through each of the components of a pre-norm Transformer block, implementing them in sequence.

3.5.1 Root Mean Square Layer Normalization

The original Transformer implementation of Vaswani et al. [2017] uses layer normalization [Ba et al., 2016] to normalize activations. Following Touvron et al. [2023], we will use root mean square layer normalization (RMSNorm; Zhang and Sennrich, 2019, equation 4) for layer normalization. Given a vector $a \in \mathbb{R}^{d_{model}}$ of activations, RMSNorm will rescale each activation a_i as follows:

$$\text{RMSNorm}(a_i) = \frac{a_i}{\text{RMS}(a)} g_i \quad (4)$$

where $\text{RMS}(a) = \sqrt{\frac{1}{d_{\text{model}}} \sum_{i=1}^{d_{\text{model}}} a_i^2 + \varepsilon}$. Here, g_i is a learnable "gain" parameter (there are d_{model} such parameters total), and ε is a hyperparameter that is often fixed at $1e-5$.

You should upcast your input to `torch.float32` to prevent overflow when you square the input. Overall, your forward method should look like:

```
in_dtype = x.dtype
x = x.to(torch.float32)
# Your code here performing RMSNorm
result = ...
# Return the result in the original dtype
return result.to(in_dtype)
```

Problem (rmsnorm): Root Mean Square Layer Normalization

Deliverable: Implement `RMSNorm` as a `torch.nn.Module`. We recommend the following interface:

```
def __init__(self, d_model: int, eps: float = 1e-5, device=None, dtype=None):
    #Construct the RMSNorm module.
    #This function should accept the following parameters:

    d_model: int                # Hidden dimension of the model
    eps: float = 1e-5           # Epsilon value for numerical stability
    device: torch.device | None # Device to store the parameters on
    dtype: torch.dtype | None   # Data type of the parameters

def forward(self, x: torch.Tensor) -> torch.Tensor:
    """Process an input tensor of shape (batch_size, sequence_length, d_model)
    and return a tensor of the same shape."""
```

Note: Remember to upcast your input to `torch.float32` before performing the normalization (and later downcast to the original dtype), as described above.

To test your implementation, implement the test adapter at `adapters.run_rmsnorm`. Then, run:

```
uv run pytest -k test_rmsnorm
```

3.5.2 Position-Wise Feed-Forward Network

In the original Transformer paper (section 3.3 of Vaswani et al. [2017]), the Transformer feed-forward network consists of two linear transformations with a ReLU activation ($\text{ReLU}(x) = \max(0, x)$) between them. The dimensionality of the inner feed-forward layer is typically 4 x the input dimensionality.

However, modern language models tend to incorporate two main changes compared to this original design: they use another activation function and employ a gating mechanism. Specifically, we will implement the "SwiGLU" activation function adopted in LLMs like Llama 3 [Grattafiori et al., 2024] and Qwen 2.5

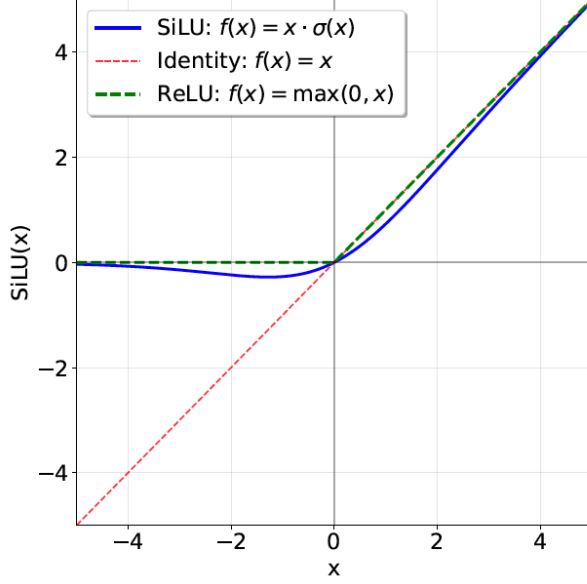


Figure 3: Comparing the SiLU (aka Swish) and ReLU activation functions.

[Yang et al., 2024], which combines the SiLU (often called Swish) activation with a gating mechanism called a Gated Linear Unit (GLU). We will also omit the bias terms sometimes used in linear layers, following most modern LLMs since PaLM [Chowdhery et al., 2022] and LLaMA [Touvron et al., 2023].

The SiLU or Swish activation function [Hendrycks and Gimpel, 2016, Elfving et al., 2017] is defined as follows:

$$\text{SiLU}(x) = x \cdot \sigma(x) = \frac{x}{1 + e^{-x}} \quad (5)$$

As can be seen in Figure 3, the SiLU activation function is similar to the ReLU activation function, but is smooth at zero.

Gated Linear Units (GLUs) were originally defined by Dauphin et al. [2017] as the element-wise product of a linear transformation passed through a sigmoid function and another linear transformation:

$$\text{GLU}(x, W_1, W_2) = \sigma(W_1 x) \odot W_2 x, \quad (6)$$

where \odot represents element-wise multiplication. Gated Linear Units are suggested to "reduce the vanishing gradient problem for deep architectures by providing a linear path for the gradients while retaining non-linear capabilities."

Putting the SiLU/Swish and GLU together, we get the SwiGLU, which we will use for our feed-forward networks:

$$\text{FFN}(x) = \text{SwiGLU}(x, W_1, W_2, W_3) = W_2 (\text{SiLU}(W_1 x) \odot W_3 x), \quad (7)$$

where $x \in \mathbb{R}^{d_{\text{model}}}$, $W_1, W_3 \in \mathbb{R}^{d_{\text{ff}} \times d_{\text{model}}}$, $W_2 \in \mathbb{R}^{d_{\text{model}} \times d_{\text{ff}}}$, and canonically, $d_{\text{ff}} = \frac{8}{3} d_{\text{model}}$.

Shazeer 2020 first proposed combining the SiLU/Swish activation with GLUs and conducted experiments showing that SwiGLU outperforms baselines like ReLU and SiLU (without gating) on language modeling

tasks. Later in the assignment, you will compare SwiGLU and SiLU. Though we’ve mentioned some heuristic arguments for these components (and the papers provide more supporting evidence), it’s good to keep an empirical perspective: a now famous quote from Shazeer’s paper is

We offer no explanation as to why these architectures seem to work; we attribute their success, as all else, to divine benevolence.

Problem (positionwise_feedforward): Implement the position-wise feed-forward network

Deliverable: Implement the SwiGLU feed-forward network, composed of a SiLU activation function and a GLU.

Note: in this particular case, you should feel free to use `torch.sigmoid` in your implementation for numerical stability.

You should set d_{ff} to approximately $\frac{8}{3} \times d_{\text{model}}$ in your implementation, while ensuring that the dimensionality of the inner feed-forward layer is a multiple of 64 to make good use of your hardware. To test your implementation against our provided tests, you will need to implement the test adapter at `[adapters.run_swiglu]`. Then, run `uv run pytest -k test_swiglu` to test your implementation.

3.5.3 Relative Positional Embeddings

To inject positional information into the model, we will implement Rotary Position Embeddings [Su et al., 2021], often called RoPE. For a given query token $q^{(i)} = W_q x^{(i)} \in \mathbb{R}^d$ at token position i , we will apply a pairwise rotation matrix R^i , giving us $q'^{(i)} = R^i q^{(i)} = R^i W_q x^{(i)}$. Here, R^i will rotate pairs of embedding elements $q_{2k-1:2k}^{(i)}$ as 2 d vectors by the angle $\theta_{i,k} = \frac{i}{\Theta(2k-1)/d}$ for $k \in \{1, \dots, d/2\}$ and some constant Θ . Thus, we can consider R^i to be a block-diagonal matrix of size $d \times d$, with blocks R_k^i for $k \in \{1, \dots, d/2\}$, with

$$R_k^i = \begin{bmatrix} \cos(\theta_{i,k}) & -\sin(\theta_{i,k}) \\ \sin(\theta_{i,k}) & \cos(\theta_{i,k}) \end{bmatrix}. \quad (8)$$

Thus we get the full rotation matrix

$$R^i = \begin{bmatrix} R_1^i & 0 & 0 & \dots & 0 \\ 0 & R_2^i & 0 & \dots & 0 \\ 0 & 0 & R_3^i & \dots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & \dots & R_{d/2}^i \end{bmatrix}, \quad (9)$$

where 0 s represent 2×2 zero matrices. While one could construct the full $d \times d$ matrix, a good solution should use the properties of this matrix to implement the transformation more efficiently. Since we only care about the relative rotation of tokens within a given sequence, we can reuse the values we compute for $\cos(\theta_{i,k})$ and $\sin(\theta_{i,k})$ across layers, and different batches. If you would like to optimize it, you may use a single RoPE module referenced by all layers, and it can have a 2 d pre-computed buffer of sin and cos values created during init with `self.register_buffer(persistent=False)`, instead of a nn. Parameter (because we do

not want to learn these fixed cosine and sine values). The exact same rotation process we did for our $q^{(i)}$ is then done for $k^{(j)}$, rotating by the corresponding R^j . Notice that this layer has no learnable parameters.

Problem (rope): Implement RoPE

Deliverable: Implement a class `RotaryPositionalEmbedding` that applies RoPE to the input tensor. The following interface is recommended:

```
def __init__(self, theta: float, d_k: int, max_seq_len: int, device=None):
    """
    Construct the RoPE module and create buffers if needed.

    Args:
        theta: Theta value for RoPE.
        d_k: Dimension of query/key vectors (should be even).
        max_seq_len: Maximum sequence length that will be inputted.
        device: torch.device | None. Device to store the buffers on.
    """
    ...

def forward(self, x: torch.Tensor, token_positions: torch.Tensor) -> torch.Tensor:
    """
    Apply RoPE to an input tensor of shape (... , seq_len, d_k) and
    return a tensor of the same shape.

    Notes:
        - Accept x with an arbitrary number of batch dimensions.
        - token_positions has shape (... , seq_len) and gives absolute
          positions per token along the sequence dimension.
        - Use token_positions to slice (precomputed) cos/sin tensors
          along the sequence dimension.
    """
    ...
```

To test your implementation, complete `[adapters.run_rope]` and make sure it passes `uv run pytest -k test_rope`.

3.5.4 Scaled Dot-Product Attention

We will now implement scaled dot-product attention as described in Vaswani et al. [2017] (section 3.2.1). As a preliminary step, the definition of the Attention operation will make use of softmax, an operation that takes an unnormalized vector of scores and turns it into a normalized distribution:

$$\text{softmax}(v)_i = \frac{\exp(v_i)}{\sum_{j=1}^n \exp(v_j)}. \quad (10)$$

Note that $\exp(v_i)$ can become inf for large values (then, $\text{inf}/\text{inf} = \text{NaN}$). We can avoid this by noticing that the softmax operation is invariant to adding any constant c to all inputs. We can leverage this property for numerical stability - typically, we will subtract the largest entry of o_i from all elements of o_i , making the new largest entry 0. You will now implement softmax, using this trick for numerical stability.

Problem (softmax): Implement softmax

Deliverable: Write a function to apply the softmax operation on a tensor. Your function should take two parameters: a tensor and a dimension i , and apply softmax to the i -th dimension of the input tensor. The output tensor should have the same shape as the input tensor, but its i -th dimension will now have a normalized probability distribution. Use the trick of subtracting the maximum value in the i -th dimension from all elements of the i -th dimension to avoid numerical stability issues.

To test your implementation, complete `[adapters.run_softmax]` and make sure it passes `uv run pytest -k test_softmax_matches_pytorch`.

We can now define the Attention operation mathematically as follows:

$$\text{Attention}(Q, K, V) = \text{softmax} \left(\frac{Q^\top K}{\sqrt{d_k}} \right) V \quad (11)$$

where $Q \in \mathbb{R}^{n \times d_k}$, $K \in \mathbb{R}^{m \times d_k}$, and $V \in \mathbb{R}^{m \times d_v}$. Here, Q , K and V are all inputs to this operation-note that these are not the learnable parameters. If you're wondering why this isn't QK^\top , see 3.3.1.

Masking: It is sometimes convenient to mask the output of an attention operation. A mask should have the shape $M \in \{ \text{True}, \text{False} \}^{n \times m}$, and each row i of this boolean matrix indicates which keys the query i should attend to. Canonically (and slightly confusingly), a value of True at position (i, j) indicates that the query i does attend to the key j , and a value of False indicates that the query does not attend to the key. In other words, "information flows" at (i, j) pairs with value True. For example, consider a 1×3 mask matrix with entries `[[True, True, False]]`. The single query vector attends only to the first two keys.

Computationally, it will be much more efficient to use masking than to compute attention on subsequences, and we can do this by taking the pre-softmax values $(\frac{Q^\top K}{\sqrt{d_k}})$ and adding a $-\infty$ in any entry of the mask matrix that is False.

Problem (scaled_dot_product_attention): Implement scaled dot-product attention

Deliverable: Implement the scaled dot-product attention function. Your implementation should handle keys and queries of shape `(batch_size, ..., seq_len, d_k)` and values of shape `(batch_size, ..., seq_len, d_v)`, where `...` represents any number of other batch-like dimensions (if provided). The implementation should return an output with the shape `(batch_size, ..., d_v)`. See section 3.3 for a discussion on batch-like dimensions.

Your implementation should also support an optional user-provided boolean mask of shape `(seq_len, seq_len)`. The attention probabilities of positions with a mask value of `True` should collectively sum to 1, and the attention probabilities of positions with a mask value of `False` should be zero.

To test your implementation against our provided tests, you will need to implement the test adapter at `[adapters.run_scaled_dot_product_attention]`.

`uv run pytest -k test_scaled_dot_product_attention` tests your implementation on third-order input tensors, while `uv run pytest -k test_4d_scaled_dot_product_attention` tests your implementation on fourth-order input tensors.

3.5.5 Causal Multi-Head Self-Attention

We will implement multi-head self-attention as described in section 3.2.2 of Vaswani et al. [2017]. Recall that, mathematically, the operation of applying multi-head attention is defined as follows:

$$\text{MultiHead}(Q, K, V) = \text{Concat}(\text{head}_1, \dots, \text{head}_h) \quad (12)$$

$$\text{for } \text{head}_i = \text{Attention}(Q_i, K_i, V_i) \quad (13)$$

with Q_i, K_i, V_i being slice number $i \in \{1, \dots, h\}$ of size d_k or d_v of the embedding dimension for Q, K , and V respectively. With Attention being the scaled dot-product attention operation defined in §3.5.4. From this we can form the multi-head self-attention operation:

$$\text{MultiHeadSelfAttention}(x) = W_O \text{MultiHead}(W_Q x, W_K x, W_V x) \quad (14)$$

Here, the learnable parameters are $W_Q \in \mathbb{R}^{hd_k \times d_{\text{model}}}$, $W_K \in \mathbb{R}^{hd_k \times d_{\text{model}}}$, $W_V \in \mathbb{R}^{hd_v \times d_{\text{model}}}$, and $W_O \in \mathbb{R}^{d_{\text{model}} \times hd_v}$. Since the Q s, K , and V s are sliced in the multi-head attention operation, we can think of W_Q, W_K and W_V as being separated for each head along the output dimension. When you have this working, you should be computing the key, value, and query projections in a total of three matrix multiplies³.

Causal masking. Your implementation should prevent the model from attending to future tokens in the sequence. In other words, if the model is given a token sequence t_1, \dots, t_n , and we want to calculate the next-word predictions for the prefix t_1, \dots, t_i (where $i < n$), the model should not be able to access (attend to) the token representations at positions t_{i+1}, \dots, t_n since it will not have access to these tokens when generating text during inference (and these future tokens leak information about the identity of the true next word, trivializing the language modeling pre-training objective). For an input token sequence t_1, \dots, t_n we can naively prevent access to future tokens by running multi-head self-attention n times (for the n unique prefixes in the sequence). Instead, we'll use causal attention masking, which allows token i to attend to all positions $j \leq i$ in the sequence. You can use `torch.triu` or a broadcasted index comparison to construct this mask, and you should take advantage of the fact that your scaled dot-product attention implementation from §3.5.4 already supports attention masking.

Applying RoPE. RoPE should be applied to the query and key vectors, but not the value vectors. Also, the head dimension should be handled as a batch dimension, because in multi-head attention, attention is being applied independently for each head. This means that precisely the same RoPE rotation should be applied to the query and key vectors for each head.

Problem (multihead_self_attention): Implement causal multi-head self-attention

Deliverable: Implement causal multi-head self-attention as a `torch.nn.Module`. Your implementation should accept (at least) the following parameters:

```
d_model: int      # Dimensionality of the Transformer block inputs
num_heads: int     # Number of heads to use in multi-head self-attention
```

³As a stretch goal, try combining the key, query, and value projections into a single weight matrix so you only need a single matrix multiply.

Following Vaswani et al. (2017), set $d_k = d_v = d_{\text{model}}/h$.

To test your implementation against our provided tests, implement the test adapter at `adapters.run_multihead_self_attention`. Then, run:

```
uv run pytest -k test_multihead_self_attention
```

3.6 The Full Transformer LM

Let's begin by assembling the Transformer block (it will be helpful to refer back to Figure 2). A Transformer block contains two 'sublayers', one for the multihead self attention, and another for the feed-forward network. In each sublayer, we first perform RMSNorm, then the main operation (MHA/FF), finally adding in the residual connection.

To be concrete, the first half (the first 'sub-layer') of the Transformer block should be implementing the following set of updates to produce an output y from an input x ,

$$y = x + \text{MultiHeadSelfAttention}(\text{RMSNorm}(x)) \quad (15)$$

Problem (transformer_block): Implement the Transformer block

Implement the **pre-norm** Transformer block as described in § 3.5 and illustrated in Figure 2. Your Transformer block should accept (at least) the following parameters:

```
d_model: int      # Dimensionality of the Transformer block inputs
num_heads: int     # Number of heads to use in multi-head self-attention
d_ff: int          # Dimensionality of the position-wise feed-forward inner layer
```

To test your implementation, implement the adapter `adapters.run_transformer_block`. Then run:

```
uv run pytest -k test_transformer_block
```

Deliverable: Transformer block code that passes the provided tests.

Now we put the blocks together, following the high level diagram in Figure 1. Follow our description of the embedding in Section 3.1.1, feed this into `num_layers` Transformer blocks, and then pass that into the three output layers to obtain a distribution over the vocabulary.

Problem (transformer_lm): Implementing the Transformer LM

Time to put it all together! Implement the Transformer language model as described in § 3.1 and illustrated in Figure 1. At minimum, your implementation should accept all the aforementioned construction parameters for the Transformer block, as well as these additional parameters:

```
vocab_size: int      # The size of the vocabulary, necessary for determining
                      # the dimensionality of the token embedding matrix
context_length: int   # The maximum context length, necessary for determining
```

```

                                # the dimensionality of the position embedding matrix
num_layers: int                # The number of Transformer blocks to use

```

To test your implementation against our provided tests, you will first need to implement the test adapter at `adapters.run_transformer_lm`. Then, run:

```
uv run pytest -k test_transformer_lm
```

Deliverable: A Transformer LM module that passes the above tests.

4 Training a Transformer LM

We now have the steps to preprocess the data (via tokenizer) and the model (Transformer). What remains is to build all of the code to support training. This consists of the following:

- Loss: we need to define the loss function (cross-entropy).
- Optimizer: we need to define the optimizer to minimize this loss (AdamW).
- Training loop: we need all the supporting infrastructure that loads data, saves checkpoints, and manages training.

4.1 Cross-entropy loss

Recall that the Transformer language model defines a distribution $p_\theta(x_{i+1} \mid x_{1:i})$ for each sequence x of length $m+1$ and $i = 1, \dots, m$. Given a training set D consisting of sequences of length m , we define the standard cross-entropy (negative log-likelihood) loss function:

$$\ell(\theta; D) = \frac{1}{|D|m} \sum_{x \in D} \sum_{i=1}^m -\log p_\theta(x_{i+1} \mid x_{1:i}) \quad (16)$$

(Note that a single forward pass in the Transformer yields $p_\theta(x_{i+1} \mid x_{1:i})$ for all $i = 1, \dots, m$.)

In particular, the Transformer computes logits $o_i \in \mathbb{R}^{vocab_size}$ for each position i , which results in:

$$p(x_{i+1} \mid x_{1:i}) = \text{softmax}(o_i)[x_{i+1}] = \frac{\exp(o_i[x_{i+1}])}{\sum_{a=1}^{vocab_size} \exp(o_i[a])}. \quad (17)$$

The cross entropy loss is generally defined with respect to the vector of logits $o_i \in \mathbb{R}^{vocab_size}$ and target x_{i+1} .

Implementing the cross entropy loss requires some care with numerical issues, just like in the case of softmax.

Problem (cross_entropy): Implement Cross Entropy

Deliverable: Write a function to compute the cross entropy loss, which takes in predicted logits o_i and targets x_{i+1} , and computes the cross entropy

$$\ell_i = -\log \text{softmax}(o_i)[x_{i+1}].$$

Your function should handle the following:

- Subtract the largest element for numerical stability.
- Cancel out `log` and `exp` whenever possible.
- Handle any additional batch dimensions and return the average across the batch. As with Section 3.3, we assume batch-like dimensions always come first, before the vocabulary size dimension.

To test your implementation, implement `adapters.run_cross_entropy`, then run:

```
uv run pytest -k test_cross_entropy
```

Perplexity Cross entropy suffices for training, but when we evaluate the model, we also want to report perplexity. For a sequence of length m where we suffer cross-entropy losses ℓ_1, \dots, ℓ_m :

$$\text{perplexity} = \exp \left(\frac{1}{m} \sum_{i=1}^m \ell_i \right) \quad (18)$$

4.2 The SGD Optimizer

Now that we have a loss function, we will begin our exploration of optimizers. The simplest gradient-based optimizer is Stochastic Gradient Descent (SGD). We start with randomly initialized parameters θ_0 . Then for each step $t = 0, \dots, T - 1$, we perform the following update:

$$\theta_{t+1} \leftarrow \theta_t - \alpha_t \nabla L(\theta_t; B_t), \quad (19)$$

where B_t is a random batch of data sampled from the dataset D , and the learning rate α_t and batch size $|B_t|$ are hyperparameters.

4.2.1 Implementing SGD in PyTorch

To implement our optimizers, we will subclass the PyTorch `torch.optim.Optimizer` class. An `Optimizer` subclass must implement two methods:

```
def __init__(self, params, ...):
    """
    Initialize your optimizer.

    Arguments:
        params: iterable
            A collection of parameters to optimize, or parameter groups
            (for applying different hyperparameters to different parts of the model).
    """
    ...
```

³⁶ Note that $o_i[k]$ refers to value at index k of the vector o_i .

⁷ This corresponds to the cross entropy between the Dirac delta distribution over x_{i+1} and the predicted softmax (o_i) distribution.

```

        Additional arguments depending on the optimizer
        (e.g., learning rate is common).
    """
    # Make sure to call the base class constructor, passing params
    # and a dictionary of hyperparameters (keys are string names).
    super().__init__(params, defaults)

def step(self):
    """
    Make one update of the parameters.

    This is called after the backward pass in training, so gradients
    are available for each parameter.

    For each parameter tensor p:
        - Access its gradient in p.grad (if it exists).
        - Update p.data in-place based on p.grad.
    """
    for group in self.param_groups:
        for p in group['params']:
            if p.grad is None:
                continue
            # Example update rule (to be implemented for your optimizer)
            p.data = p.data - group['lr'] * p.grad

```

The PyTorch optimizer API has a few subtleties, so it's easier to explain it with an example. To make our example richer, we'll implement a slight variation of SGD where the learning rate decays over training, starting with an initial learning rate α and taking successively smaller steps over time:

$$\theta_{t+1} = \theta_t - \frac{\alpha}{\sqrt{t+1}} \nabla L(\theta_t; B_t) \quad (20)$$

Let's see how this version of SGD would be implemented as a PyTorch Optimizer:

```

from collections.abc import Callable, Iterable
from typing import Optional
import torch
import math
class SGD(torch.optim.Optimizer):
    def __init__(self, params, lr=1e-3):
        if lr < 0:
            raise ValueError(f"Invalid learning rate: {lr}")
        defaults = {"lr": lr}
        super().__init__(params, defaults)
    def step(self, closure: Optional[Callable] = None):
        loss = None if closure is None else closure()
        for group in self.param_groups:
            lr = group["lr"] # Get the learning rate.

```

```

for p in group["params"]:
    if p.grad is None:
        continue
    state = self.state[p] # Get state associated with p.
    t = state.get("t", 0) # Get iteration number from the state, or initial value.
    grad = p.grad.data # Get the gradient of loss with respect to p.
    p.data -= lr / math.sqrt(t + 1) * grad # Update weight tensor in-place.
    state["t"] = t + 1 # Increment iteration number.
return loss

```

In `init`, we pass the parameters to the optimizer, as well as default hyperparameters, to the base class constructor (the parameters might come in groups, each with different hyperparameters). In case the parameters are just a single collection of `torch.nn.Parameter` objects, the base constructor will create a single group and assign it the default hyperparameters. Then, in `step`, we iterate over each parameter group, then over each parameter in that group, and apply Eq 20. Here, we keep the iteration number as a state associated with each parameter: we first read this value, use it in the gradient update, and then update it. The API specifies that the user might pass in a callable closure to re-compute the loss before the optimizer step. We won't need this for the optimizers we'll use, but we add it to comply with the API.

To see this working, we can use the following minimal example of a training loop:

```

weights = torch.nn.Parameter(5 * torch.randn((10, 10)))
opt = SGD([weights], lr=1)
for t in range(100):
    opt.zero_grad() # Reset the gradients for all learnable parameters.
    loss = (weights**2).mean() # Compute a scalar loss value.
    print(loss.cpu().item())
    loss.backward() # Run backward pass, which computes gradients.
    opt.step() # Run optimizer step.

```

This is the typical structure of a training loop: in each iteration, we will compute the loss and run a step of the optimizer. When training language models, our learnable parameters will come from the model (in PyTorch, `m.parameters()` gives us this collection). The loss will be computed over a sampled batch of data, but the basic structure of the training loop will be the same.

4.3 AdamW

Modern language models are typically trained with more sophisticated optimizers, instead of SGD. Most optimizers used recently are derivatives of the Adam optimizer [Kingma and Ba, 2015]. We will use AdamW [Loshchilov and Hutter, 2019], which is in wide use in recent work. AdamW proposes a modification to Adam that improves regularization by adding weight decay (at each iteration, we pull the parameters towards 0), in a way that is decoupled from the gradient update. We will implement AdamW as described in algorithm 2 of Loshchilov and Hutter [2019].

AdamW is stateful: for each parameter, it keeps track of a running estimate of its first and second moments. Thus, AdamW uses additional memory in exchange for improved stability and convergence. Besides the learning rate α , AdamW has a pair of hyperparameters (β_1, β_2) that control the updates to the

moment estimates, and a weight decay rate λ . Typical applications set (β_1, β_2) to $(0.9, 0.999)$, but large language models like LLaMA [Touvron et al., 2023] and GPT-3 [Brown et al., 2020] are often trained with $(0.9, 0.95)$. The algorithm can be written as follows, where ϵ is a small value (e.g., 10^{-8}) used to improve numerical stability in case we get extremely small values in v :

Algorithm 1 AdamW Optimizer

```

1:  $\text{init}(\theta)$  ▷ Initialize learnable parameters
2:  $m \leftarrow 0$  ▷ Initial value of the first moment vector; same shape as  $\theta$ 
3:  $v \leftarrow 0$  ▷ Initial value of the second moment vector; same shape as  $\theta$ 
4: for  $t = 1, \dots, T$  do
5:   Sample batch of data  $B_t$ 
6:    $g \leftarrow \nabla_{\theta} \ell(\theta; B_t)$  ▷ Gradient of the loss at time  $t$ 
7:    $m \leftarrow \beta_1 m + (1 - \beta_1)g$  ▷ Update first moment estimate
8:    $v \leftarrow \beta_2 v + (1 - \beta_2)g^2$  ▷ Update second moment estimate
9:    $\alpha_t \leftarrow \alpha \frac{\sqrt{1 - \beta_2^t}}{1 - \beta_1^t}$  ▷ Compute adjusted  $\alpha$ 
10:   $\theta \leftarrow \theta - \alpha_t \frac{m}{\sqrt{v + \epsilon}}$  ▷ Update parameters
11:   $\theta \leftarrow \theta - \alpha \lambda \theta$  ▷ Apply weight decay
12: end for

```

Note that t starts at 1. You will now implement this optimizer.

Problem (adamw): Implement AdamW

Deliverable: Implement the AdamW optimizer as a subclass of `torch.optim.Optimizer`. Your class should take the learning rate α in `__init__`, as well as the β , ϵ , and λ hyperparameters.

To help you keep state, the base `Optimizer` class provides a dictionary `self.state`, which maps `nn.Parameter` objects to a dictionary that stores any information you need for that parameter (for AdamW, this would be the moment estimates).

To test your implementation, implement `adapters.get_adamw_cls` and make sure it passes:

```
uv run pytest -k test_adamw
```

4.4 Learning rate scheduling

The value for the learning rate that leads to the quickest decrease in loss often varies during training. In training Transformers, it is typical to use a learning rate schedule, where we start with a bigger learning rate, making quicker updates in the beginning, and slowly decay it to a smaller value as the model trains⁴. In this assignment, we will implement the cosine annealing schedule used to train LLaMA [Touvron et al., 2023].

A scheduler is simply a function that takes the current step t and other relevant parameters (such as the initial and final learning rates), and returns the learning rate to use for the gradient update at step t . The simplest schedule is the constant function, which will return the same learning rate given any t .

⁴It's sometimes common to use a schedule where the learning rate rises back up (restarts) to help get past local minima.

The cosine annealing learning rate schedule takes (i) the current iteration t , (ii) the maximum learning rate α_{\max} , (iii) the minimum (final) learning rate α_{\min} , (iv) the number of warm-up iterations T_w , and (v) the number of cosine annealing iterations T_c . The learning rate at iteration t is defined as:

(Warm-up) If $t < T_w$, then $\alpha_t = \frac{t}{T_w} \alpha_{\max}$.

(Cosine annealing) If $T_w \leq t \leq T_c$, then $\alpha_t = \alpha_{\min} + \frac{1}{2} \left(1 + \cos \left(\frac{t - T_w}{T_c - T_w} \pi \right) \right) (\alpha_{\max} - \alpha_{\min})$.

(Post-annealing) If $t > T_c$, then $\alpha_t = \alpha_{\min}$.

Problem (learning_rate_schedule): Implement cosine learning rate schedule with warmup

Write a function that takes $t, \alpha_{\max}, \alpha_{\min}, T_w, T_c$ and returns the learning rate α_t according to the scheduler defined above.

Then implement `adapters.get_lr_cosine_schedule` and make sure it passes:

```
uv run pytest -k test_get_lr_cosine_schedule
```

Deliverable: A function that computes the cosine learning rate schedule with warmup and passes the provided tests.

4.5 Gradient clipping

During training, we can sometimes hit training examples that yield large gradients, which can destabilize training. To mitigate this, one technique often employed in practice is gradient clipping. The idea is to enforce a limit on the norm of the gradient after each backward pass before taking an optimizer step.

Given the gradient (for all parameters) g , we compute its ℓ_2 -norm $\|g\|_2$. If this norm is less than a maximum value M , then we leave g as is; otherwise, we scale g down by a factor of $\frac{M}{\|g\|_2 + \epsilon}$ (where a small ϵ , like 10^{-6} , is added for numeric stability). Note that the resulting norm will be just under M .

Problem (gradient_clipping): Implement gradient clipping

Write a function that implements gradient clipping. Your function should take a list of parameters and a maximum ℓ_2 -norm. It should modify each parameter gradient in place. Use $\epsilon = 10^{-6}$ (the PyTorch default).

Then, implement the adapter `adapters.run_gradient_clipping` and make sure it passes:

```
uv run pytest -k test_gradient_clipping
```

Deliverable: A gradient clipping function that passes the provided tests.

4.6 Expected Output So far

At this current stage, we expect to pass the following 18 tests.

```
uv run pytest
```

```
tests/test_model.py::test_linear PASSED
```

```
tests/test_model.py::test_embedding PASSED
```


tests/test_model.py::test_swiglu PASSED
tests/test_model.py::test_scaled_dot_product_attention PASSED
tests/test_model.py::test_4d_scaled_dot_product_attention PASSED
tests/test_model.py::test_multihead_self_attention PASSED
tests/test_model.py::test_multihead_self_attention_with_rope PASSED
tests/test_model.py::test_transformer_lm PASSED
tests/test_model.py::test_transformer_lm_truncated_input PASSED
tests/test_model.py::test_transformer_block PASSED
tests/test_model.py::test_rmsnorm PASSED
tests/test_model.py::test_rope PASSED
tests/test_model.py::test_silu_matches_pytorch PASSED
tests/test_nn_utils.py::test_softmax_matches_pytorch PASSED
tests/test_nn_utils.py::test_cross_entropy PASSED
tests/test_nn_utils.py::test_gradient_clipping PASSED
tests/test_optimizer.py::test_adamw PASSED
tests/test_optimizer.py::test_get_lr_cosine_schedule PASSED

5 Training loop

We will now finally put together the major components we’ve built so far: the tokenized data, the model, and the optimizer.

5.1 Data Loader

The tokenized data is a single sequence of tokens $x = (x_1, \dots, x_n)$. Even though the source data might consist of separate documents (e.g., different web pages, or source code files), a common practice is to concatenate all of those into a single sequence of tokens, adding a delimiter between them (such as the `<|endoftext|>` token).

A data loader turns this into a stream of batches, where each batch consists of B sequences of length m , paired with the corresponding next tokens, also with length m . For example, for $B = 1, m = 3$, $([x_2, x_3, x_4], [x_3, x_4, x_5])$ would be one potential batch.

Loading data in this way simplifies training for a number of reasons. First, any $1 \leq i \leq n - m$ gives a valid training sequence, so sampling sequences are trivial. Since all training sequences have the same length, there’s no need to pad input sequences, which improves hardware utilization (also by increasing batch size B). Finally, we also don’t need to fully load the full dataset to sample training data, making it easy to handle large datasets that might not otherwise fit in memory.

Problem (data_loading): Implement data loading

Deliverable: Write a function that takes a numpy array x (integer array with token IDs), a `batch_size`, a `context_length`, and a PyTorch device string (e.g., `'cpu'` or `'cuda:0'`), and returns a pair of tensors: the sampled input sequences and the corresponding next-token targets.

Both tensors should have shape `(batch_size, context_length)` containing token IDs, and both should be placed on the requested device.

To test your implementation against our provided tests, you will first need to implement the test adapter at `adapters.run_get_batch`. Then, run:

```
uv run pytest -k test_get_batch
```

What if the dataset is too big to load into memory? We can use a Unix systemcall named `mmap` which maps a file on disk to virtual memory, and lazily loads the file contents when that memory location is accessed. Thus, you can “pretend” you have the entire dataset in memory. Numpy implements this through `np.memmap` (or the flag `mmap_mode='r'` to `np.load`, if you originally saved the array with `np.save`), which will return a numpy array-like object that loads the entries on-demand as you access them. When sampling from your dataset (i.e., a numpy array) during training, be sure load the dataset in memorymapped mode (via `np.memmap` or the flag `mmap_mode='r'` to `np.load`, depending on how you saved the array). Make sure you also specify a `dtype` that matches the array that you’re loading. It may be helpful to explicitly verify that the memory-mapped data looks correct (e.g., doesn’t contain values beyond the expected vocabulary size).

5.2 Checkpointing

In addition to loading data, we will also need to save models as we train. When running jobs, we often want to be able to resume a training run that for some reason stopped midway (e.g., due to your job timing out, machine failure, etc). Even when all goes well, we might also want to later have access to intermediate models (e.g., to study training dynamics post-hoc, take samples from models at different stages of training, etc).

A checkpoint should have all the states that we need to resume training. We of course want to be able to restore model weights at a minimum. If using a stateful optimizer (such as AdamW), we will also need to save the optimizer's state (e.g., in the case of AdamW, the moment estimates). Finally, to resume the learning rate schedule, we will need to know the iteration number we stopped at. PyTorch makes it easy to save all of these: every `nn.Module` has a `state_dict()` method that returns a dictionary with all learnable weights; we can restore these weights later with the sister method `load_state_dict()`. The same goes for any `nn.optim.Optimizer`. Finally, `torch.save(obj, dest)` can dump an object (e.g., a dictionary containing tensors in some values, but also regular Python objects like integers) to a file (path) or file-like object, which can then be loaded back into memory with `torch.load(src)`.

Problem (checkpointing): Implement model checkpointing

Implement the following two functions to load and save checkpoints:

```
def save_checkpoint(model, optimizer, iteration, out):
    """Dump all state from model, optimizer, and iteration into `out`."""
    # Use state_dict() for both model and optimizer
    # Use torch.save(obj, out) to dump obj into out
```

A typical choice is to save `obj` as a dictionary, but any format is allowed as long as it can later be loaded correctly.

Parameters:

- `model`: `torch.nn.Module`
- `optimizer`: `torch.optim.Optimizer`
- `iteration`: `int`
- `out`: `str` | `os.PathLike` | `typing.BinaryIO` | `typing.IO[bytes]`

```
def load_checkpoint(src, model, optimizer):
    """Load a checkpoint from `src` and recover model and optimizer states."""
    # Use torch.load(src) to recover saved state
    # Call load_state_dict on both model and optimizer
    # Return the saved iteration number
```

Parameters:

- `src`: `str` | `os.PathLike` | `typing.BinaryIO` | `typing.IO[bytes]`

- `model: torch.nn.Module`
- `optimizer: torch.optim.Optimizer`

—
To test your implementation, implement the adapters `adapters.run_save_checkpoint` and `adapters.run_load_checkpoint`, and then run:

```
uv run pytest -k test_checkpointing
```

Deliverable: Implemented save and load checkpoint functions that pass the provided tests.

5.3 Training loop

Now, it's finally time to put all of the components you implemented together into your main training script. It will pay off to make it easy to start training runs with different hyperparameters (e.g., by taking them as command-line arguments), since you will be doing these many times later to study how different choices impact training.

Problem (training_together): Put it together

Deliverable: Write a script that runs a training loop to train your model on user-provided input. In particular, we recommend that your training script allow for (at least) the following:

- Ability to configure and control the various model and optimizer hyperparameters.
- Memory-efficient loading of large training and validation datasets with `np.memmap`.
- Serializing checkpoints to a user-provided path.
- Periodically logging training and validation performance (e.g., to console and/or an external service like Weights & Biases)^a.

^a wandb.ai

6 Generating text

Now that we can train models, the last piece we need is the ability to generate text from our model. Recall that a language model takes in a (possibly batched) integer sequence of length (`sequence_length`) and produces a matrix of size (`sequence_length` × vocab size), where each element of the sequence is a probability distribution predicting the next word after that position. We will now write a few functions to turn this into a sampling scheme for new sequences.

Softmax By standard convention, the language model output is the output of the final linear layer (the "logits") and so we have to turn this into a normalized probability via the softmax operation, which we saw earlier in Eq 10.

Decoding To generate text (decode) from our model, we will provide the model with a sequence of prefix

tokens (the "prompt"), and ask it to produce a probability distribution over the vocabulary that predicts the next word in the sequence. Then, we will sample from this distribution over the vocabulary items to determine the next output token.

Concretely, one step of the decoding process should take in a sequence $x_{1...t}$ and return a token x_{t+1} via the following equation,

$$P(x_{t+1} = i \mid x_{1...t}) = \frac{\exp(v_i)}{\sum_j \exp(v_j)}$$

$$v = \text{TransformerLM}(x_{1...t})_t \in \mathbb{R}^{\text{vocab_size}}$$

where TransformerLM is our model which takes as input a sequence of sequence_length and produces a matrix of size (sequence_length \times vocab_size), and we take the last element of this matrix, as we are looking for the next word prediction at the t -th position.

This gives us a basic decoder by repeatedly sampling from these one-step conditionals (appending our previously-generated output token to the input of the next decoding timestep) until we generate the end-of-sequence token `<endoftext>` (or a user-specified maximum number of tokens to generate).

Decoder tricks We will be experimenting with small models, and small models can sometimes generate very low quality texts. Two simple decoder tricks can help fix these issues. First, in temperature scaling we modify our softmax with a temperature parameter τ , where the new softmax is

$$\text{softmax}(v, \tau)_i = \frac{\exp(v_i/\tau)}{\sum_{j=1}^{\text{vocab_size}} \exp(v_j/\tau)} \quad (24)$$

Note how setting $\tau \rightarrow 0$ makes it so that the largest element of v dominates, and the output of the softmax becomes a one-hot vector concentrated at this maximal element.

Second, another trick is nucleus or top- p sampling, where we modify the sampling distribution by truncating low-probability words. Let q be a probability distribution that we get from a (temperature-scaled) softmax of size (vocab_size). Nucleus sampling with hyperparameter p produces the next token according to the equation

$$P(x_{t+1} = i \mid q) = \begin{cases} \frac{q_i}{\sum_{j \in V(p)} q_j} & \text{if } i \in V(p) \\ 0 & \text{otherwise} \end{cases}$$

where $V(p)$ is the smallest set of indices such that $\sum_{j \in V(p)} q_j \geq p$. You can compute this quantity easily by first sorting the probability distribution q by magnitude, and selecting the largest vocabulary elements until you reach the target level of α .

Problem (decoding): Decoding

Deliverable: Implement a function to decode from your language model. We recommend that you support the following features:

- Generate completions for a user-provided prompt (i.e., take in some $x_{1...t}$ and sample a completion until you hit an `<endoftext>` token).
- Allow the user to control the maximum number of generated tokens.

- Given a desired temperature value, apply softmax temperature scaling to the predicted next-word distributions before sampling.
- Top-p sampling Holtzman et al. (2019) (also referred to as nucleus sampling), given a user-specified threshold value.

7 Experiments

Now it is time to put everything together and train (small) language models on a pretraining dataset.

7.1 How to Run Experiments and Deliverables

The best way to understand the rationale behind the architectural components of a Transformer is to actually modify it and run it yourself. There is no substitute for hands-on experience.

To this end, it's important to be able to experiment quickly, consistently, and keep records of what you did. To experiment quickly, we will be running many experiments on a small scale model (17 M parameters) and simple dataset (TinyStories). To do things consistently, you will ablate components and vary hyperparameters in a systematic way, and to keep records we will ask you to submit a log of your experiments and learning curves associated with each experiment.

To make it possible to submit loss curves, make sure to periodically evaluate validation losses and record both the number of steps and wallclock times. You might find logging infrastructure such as Weights and Biases helpful.

Problem (experiment_log): Experiment logging

For your training and evaluation code, create experiment tracking infrastructure that allows you to track your experiments and loss curves with respect to gradient steps and wallclock time.

Deliverable: Logging infrastructure code for your experiments and an experiment log (a document of all the things you tried) for the assignment problems below in this section.

7.2 TinyStories

We are going to start with a very simple dataset (TinyStories; Eldan and Li, 2023) where models will train quickly, and we can see some interesting behaviors. The instructions for getting this dataset is at section 1. An example of what this dataset looks like is below.

Example (tinystories_example): One example from TinyStories

Once upon a time there was a little boy named Ben. Ben loved to explore the world around him. He saw many amazing things, like beautiful vases that were on display in a store. One day, Ben was walking through the store when he came across a very special vase. When Ben saw it he was amazed! He said, "Wow, that is a really amazing vase! Can I buy it?" The shopkeeper smiled and said, "Of course you can. You can take it home and show all your friends how amazing it is!" So Ben took the vase home and he was so proud of it! He called his friends over and showed them the amazing vase. All his friends thought the vase was beautiful and couldn't believe how lucky Ben was. And that's how Ben found an amazing vase in the store!

Hyperparameter tuning We will tell you some very basic hyperparameters to start with and ask you to find some settings for others that work well.

vocab_size When using the off-the-shelf `tiktoken`, the vocabulary is fixed. You do not need to modify `vocab_size` in `tiktoken.get_encoding("gpt2")`, since it already comes with a pre-built

vocabulary and merge rules.

context_length 256. Simple datasets such as TinyStories might not need long sequence lengths, but for the later OpenWebText data, you may want to vary this. Try varying this and seeing the impact on both the per-iteration runtime and the final perplexity.

d_model 512. This is slightly smaller than the 768 dimensions used in many small Transformer papers, but this will make things faster.

d_ff 1344. This is roughly $\frac{8}{3}d_{model}$ while being a multiple of 64, which is good for GPU performance.

RoPE theta parameter $\Theta = 10000$.

number of layers and heads 4 layers, 16 heads. Together, this will give about 17M non-embedding parameters which is a fairly small Transformer.

total tokens processed We tested with a batch size of 32, a total of 5,000 steps, and a context length of 256, which are executable on a Quadro RTX 6000.

You should do some trial and error to find good defaults for the following other hyperparameters: learning rate, learning rate warmup, other AdamW hyperparameters ($\beta_1, \beta_2, \epsilon$), and weight decay. You can find some typical choices of such hyperparameters in Kingma and Ba [2015].

Putting it together Now you can put everything together by getting a trained BPE tokenizer, tokenizing the training dataset, and running this in the training loop that you wrote.

Tips and tricks for debugging model architectures We highly recommend getting comfortable with your IDE's built-in debugger (e.g., VSCode/PyCharm), which will save you time compared to debugging with print statements. If you use a text editor, you can use something more like pdb. A few other good practices when debugging model architectures are:

- A common first step when developing any neural net architecture is to overfit to a single minibatch. If your implementation is correct, you should be able to quickly drive the training loss to near-zero.
- Set debug breakpoints in various model components, and inspect the shapes of intermediate tensors to make sure they match your expectations.
- Monitor the norms of activations, model weights, and gradients to make sure they are not exploding or vanishing.

Problem (learning_rate): Tune the learning rate

The learning rate is one of the most important hyperparameters to tune. Taking the base model you've trained, answer the following questions:

(a) Perform a hyperparameter sweep over the learning rates and report the final losses (or note divergence if the optimizer diverges).

Deliverable 1: Learning curves associated with multiple learning rates. Explain your hyperparameter search strategy.

Deliverable 2: A model with validation loss (per-token) on TinyStories of at most 1.45.

Low-Resource/Downscaling Tip: Train for few steps on CPU

If you are running on `cpu` or `mps` (Metal GPU on Apple Silicon), you should instead reduce the total tokens processed count to 40,000,000, which will be sufficient to produce reasonably fluent text. You may also increase the target validation loss from 1.45 to 2.00.

Running our solution code with a tuned learning rate on an M3 Max chip and 36 GB of RAM, we use `batch size × total step count × context length = 32 × 5000 × 256 = 40,960,000` tokens, which takes 1 hour and 22 minutes on `cpu` and 36 minutes on `mps`. At step 5000, we achieve a validation loss of 1.80.

Some additional tips:

- When using X training steps, adjust the cosine learning rate decay schedule so that it terminates (reaches the minimum learning rate) at precisely step X .

- When using `mps`, do not use TF32 kernels. Do not set:

```
torch.set_float32_matmul_precision('high')
```

as you might with CUDA devices. With `mps` (torch version 2.6.0), TF32 kernels silently break training and cause instability.

- You can speed up training by JIT-compiling your model with `torch.compile`. Specifically:

```
# On cpu
```

```
model = torch.compile(model)
```

```
# On mps
```

```
model = torch.compile(model, backend="aot_eager")
```

Compilation with Inductor is not supported on `mps` as of torch version 2.6.0.

(b) Folk wisdom is that the best learning rate is "at the edge of stability." Investigate how the point at which learning rates diverge is related to your best learning rate.

Deliverable: Learning curves of increasing learning rate which include at least one divergent run and an analysis of how this relates to convergence rates.

With your decoder in hand, we can now generate text! We will generate from the model and see how good it is. As a reference, you should get outputs that look at least as good as the example below.

Example (`ts_generate_example`): Sample output from a TinyStories language model

Once upon a time, there was a pretty girl named Lily. She loved to eat gum, especially the big black one. One day, Lily's mom asked her to help cook dinner. Lily was so excited! She loved to help her mom.

Lily's mom made a big pot of soup for dinner. Lily was so happy and said, "Thank you, Mommy! I

love you.” She helped her mom pour the soup into a big bowl.

After dinner, Lily’s mom made some yummy soup. Lily loved it! She said, ”Thank you, Mommy! This soup is so yummy!”

Her mom smiled and said, ”I’m glad you like it, Lily.” They finished cooking and continued to cook together.

The end.

Low-Resource/Downscaling Tip: Generate text on CPU

If instead you used the low-resource configuration with 40M tokens processed, you should see generations that still resemble English but are not as fluent as above. For example, our sample output from a TinyStories language model trained on 40M tokens is below:

Once upon a time, there was a little girl named Sue. Sue had a tooth that she loved very much. It was his best head. One day, Sue went for a walk and met a ladybug! They became good friends and played on the path together.

”Hey, Polly! Let’s go out!” said Tim. Sue looked at the sky and saw that it was difficult to find a way to dance shining. She smiled and agreed to help the talking!”

As Sue watched the sky moved, what it was. She...

Here is the precise problem statement and what we ask for:

Problem (generate): Generate text

Using your decoder and your trained checkpoint, report the text generated by your model. You may need to manipulate decoder parameters (e.g., temperature, top-p, etc.) to get fluent outputs.

Deliverable: A text dump of at least 256 tokens of text (or until the first `<endoftext>` token), and a brief comment on the fluency of this output and at least two factors which affect how good or bad this output is.

8 Profiling and Benchmarking

Before implementing any optimization, it is helpful to first profile our program to understand where it spends resources (e.g., time and memory). Otherwise, we risk optimizing parts of the model that don’t account for significant time or memory, and therefore not seeing measurable end-to-end improvements.

We will implement three performance evaluation paths: (a) a simple, end-to-end benchmarking using the Python standard library to time our forward and backward passes, (b) profile compute with the NVIDIA Nsight Systems tool to understand how that time is distributed across operations on both the CPU and GPU, and (c) profile memory usage.

8.1 Setup - Importing your Basics Transformer Model

Let's start by making sure that you can load the model from the previous parts of this assignment. In the previous parts, we set up our model in a Python package so that it could be easily imported later. You can test that you can import your model with:

```
~$ uv run python
>>> import cse599o_basics
>>>
```

The relevant modules from assignment 1 should now be available (e.g., for `model.py`, you can import it with `import cse599o_basics.model`).

8.2 Model Sizing

Throughout this assignment, we will be benchmarking and profiling models to better understand their performance. To get a sense of how things change at scale, we will work with and refer to the following model configurations. For all models, we'll use a batch size of 4, with varying context lengths.

Size	d_model	d_ff	num_layers	num_heads
small	768	3072	12	12
medium	1024	4096	24	16
large	1280	5120	36	20

Table 1: Specifications of different model sizes

8.3 Profiling Execution Time

We will now implement a simple performance evaluation script. We will be testing many variations of our model, so it will pay off to have your script enable these variations via command-line arguments to make them easy to run later on. We also highly recommend running sweeps over benchmarking hyperparameters, such as model size, context length, etc., using `sbatch` or `submitit` on Slurm for quick iteration.

To start off, let's do the simplest possible profiling of our model by timing the forward and backward passes. Since we will only be measuring speed and memory, **we can use random weights and data**.

Measuring performance is subtle - some common traps can cause us to not measure what we want. For benchmarking GPU code, one caveat is that CUDA calls are asynchronous. When you call a CUDA kernel, such as when you invoke `torch.matmul`, the function call returns control to your code without waiting for the matrix multiplication to finish. In this way, the CPU can continue running while the GPU computes the matrix multiplication. On the other hand, this means that naïvely measuring how long the `torch.matmul` call takes to return does not tell us how long the GPU takes to actually run the matrix multiplication. In PyTorch, we can call `torch.cuda.synchronize()` to wait for all GPU kernels to complete, allowing us to get more accurate measurements of CUDA kernel runtime. With this in mind, let's write our basic profiling infrastructure.

Problem (benchmarking_script)

(a) Write a script (benchmark.py) to perform basic end-to-end benchmarking of the forward and backward passes in your model. Specifically, your script should support the following:

- Given hyperparameters (e.g., number of layers), initialize a model.
- Generate a random batch of data.
- Run w warm-up steps (before you start measuring time), then time the execution of n steps (either only forward, or both forward and backward passes, depending on an argument). For timing, you can use the Python `timeit` module (e.g., either using the `timeit` function, or using `timeit.default_timer()`, which gives you the system's highest resolution clock, thus a better default for benchmarking than `time.time()`).
- Call `torch.cuda.synchronize()` after each step.

Deliverable: A script that will initialize a basics Transformer model with the given hyperparameters, create a random batch of data, and time forward and backward passes.

(b) Time the forward and backward passes for the model sizes described before. Use 5 warmup steps and compute the average and standard deviation of timings over 10 measurement steps. How long does a forward pass take? How about a backward pass? Do you see high variability across measurements, or is the standard deviation small?

Deliverable: A 1-2 sentence response with your timings.

8.4 Nsight Systems Profiler

End-to-end benchmarking does not tell us where our model spends time and memory during forward and backward passes, and so does not expose specific optimization opportunities. To know how much time our program spends in each component (e.g., function), we can use a profiler. An execution profiler instruments the code by inserting guards when functions begin and finish running, and thus can give detailed execution statistics at the function level (such as number of calls, how long they take on average, cumulative time spent on this function, etc).

Standard Python profilers (e.g., CProfile) are not able to profile CUDA kernels since these kernels are executed asynchronously on the GPU. Fortunately, NVIDIA ships a profiler that we can use via the CLI `nsys`, which we have already installed for you. You can check whether `nsys` is working on your machines.

```
~$ export PATH=/usr/local/cuda-13.0/bin:$PATH
```

```
~$ nsys status -e
```

```
Timestamp counter supported: Yes
```

```
CPU Profiling Environment Check
```

```
Root privilege: disabled
```

```
Linux Kernel Paranoid Level = 2
Linux Distribution = rocky
Linux Kernel Version = 6.12.0-55.32.1.el10_0.x86_64: OK
Linux perf_event_open syscall available: OK
Sampling trigger event available: OK
Intel(c) Last Branch Record support: Available
CPU Profiling Environment (process-tree): OK
CPU Profiling Environment (system-wide): Fail
```

In this part of the assignment, you will use nsys to analyze the runtime of your Transformer model. Using nsys is straightforward: we can simply run your Python script from the previous section with nsys profile prepended. For example, you can profile a script benchmark.py and write the output to a file result.nsys.rep with:

```
~$ uv run nsys profile -o result python benchmark.py
```

You can then view the profile on your local machine with the NVIDIA Nsight Systems desktop application. Selecting a particular CUDA API call (on the CPU) in the CUDA API row of the profile will highlight all corresponding kernel executions (on the GPU) in the CUDA HW row.

We encourage you to experiment with various command-line options for nsys profile to get a sense of what it can do. Notably, you can get Python backtraces for each CUDA API call with `-python-backtrace=cuda`, though this may introduce overhead. You can also annotate your code with NVTX ranges, which will appear as blocks in the NVTX row of the profile capturing all CUDA API calls and associated kernel executions. In particular, you should use NVTX ranges to ignore the warm-up steps in your benchmarking script (by applying a filter on the NVTX row in the profile). You can also isolate which kernels are responsible for the forward and backward passes of your model, and you can even isolate which kernels are responsible for different parts of a self-attention layer by annotating your implementation as follows:

```
1 import torch.cuda.nvtx as nvtx
2
3 @nvtx.range("scaled dot product attention")
4 def annotated_scaled_dot_product_attention(
5     # Q, K, V, mask
6 ):
7     with nvtx.range("computing attention scores"):
8         # compute attention scores between Q and K
9
10    with nvtx.range("computing softmax"):
11        # compute softmax of attention scores
12
13    with nvtx.range("final matmul"):
14        # compute output projection
```

You can swap your original implementation with the annotated version in your benchmarking script via:

```
1 cse599o_basics.model.scaled_dot_product_attention = annotated_scaled_dot_product_attention
```

Finally, you can use the `-pytorch` command-line option with `nsys` to automatically annotate calls to the PyTorch C++ API with NVTX ranges.

Problem (`nsys__profile`)

Profile your forward pass, backward pass, and optimizer step using `nsys` with each of the model sizes described in Table 1 and context lengths of 128, 256, 512, and 1024 (you may run out of memory with some of these context lengths for the larger models, in which case just note it in your report).

(a) What is the total time spent on your forward pass? Does it match what we had measured before with the Python standard library?

Deliverable: A 1-2 sentence response.

(b) What CUDA kernel takes the most cumulative GPU time during the forward pass? How many times is this kernel invoked during a single forward pass of your model? Is it the same kernel that takes the most runtime when you do both forward and backward passes? (Hint: look at the "CUDA GPU Kernel Summary" under "Stats Systems View", and filter using NVTX ranges to identify which parts of the model are responsible for which kernels.)

Deliverable: A 1-2 sentence response.

(c) Although the vast majority of FLOPs take place in matrix multiplications, you will notice that several other kernels still take a non-trivial amount of the overall runtime. What other kernels besides matrix multiplies do you see accounting for non-trivial CUDA runtime in the forward pass?

Deliverable: A 1-2 sentence response.

(d) Profile running one complete training step with your implementation of AdamW (i.e., the forward pass, computing the loss and running a backward pass, and finally an optimizer step, as you'd do during training). How does the fraction of time spent on matrix multiplication change, compared to doing inference (forward pass only)? How about other kernels?

Deliverable: A 1-2 sentence response.

(e) Compare the runtime of the softmax operation versus the matrix multiplication operations within the self-attention layer of your model during a forward pass. How does the difference in runtimes compare to the difference in FLOPs?

Deliverable: A 1-2 sentence response.

8.5 Profiling Memory

So far, we have been looking at compute performance. We'll now shift our attention to memory, another major resource in language model training and inference. PyTorch also ships with a powerful memory profiler, which can keep track of allocations over time.

To use the memory profiler, you can modify your benchmarking script as follows:

```
1 ... # warm-up phase in your benchmarking script
2 # Start recording memory history.
3 torch.cuda.memory._record_memory_history(max_entries=1000000)
4 ... # what you want to profile in your benchmarking script
5 # Save a pickle file to be loaded by PyTorch's online tool.
6 torch.cuda.memory._dump_snapshot("memory_snapshot.pickle")
7 # Stop recording history.
8 torch.cuda.memory._record_memory_history(enabled=None)
```

This will output a file `memory_snapshot.pickle` that you can load into the following online tool: https://pytorch.org/memory_viz. This tool will let you see the overall memory usage timeline as well as each individual allocation that was made, with its size and a stack trace leading to the code where it originates. To use this tool, you should open the link above in a Web browser, and then drag and drop your Pickle file onto the page.

You will now use the PyTorch profiler to analyze the memory usage of your model.

Problem (memory_profiling)

Profile your forward pass, backward pass, and optimizer step of the large model from Table 1 with context lengths of 128, 256, and 512.

(a) Add an option to your profiling script to run your model through the memory profiler. It may be helpful to reuse some of your previous infrastructure (e.g., load specific model sizes, etc). Then, run your script to get a memory profile of the large model when either doing inference only (just forward pass) or a full training step. How do your memory timelines look like? Can you tell which stage is running based on the peaks you see?

Deliverable: Two images of the "Active memory timeline" of a large model, from the `memory_viz` tool: one for the forward pass, and one for running a full training step (forward and backward passes, then optimizer step), and a 2-3 sentence response.

(b) What is the peak memory usage of each context length when doing a forward pass? What about when doing a full training step?

Deliverable: two numbers per context length.

(c) Now look closely at the "Active Memory Timeline" from pytorch.org/memory_viz of a memory snapshot of the large model doing a forward pass. When you reduce the "Detail" level, the tool hides the smallest allocations to the corresponding level (e.g., putting "Detail" at 10% only shows the 10% largest allocations). What is the size of the largest allocations shown? Looking through the stack trace, can you tell where those allocations come from?

Deliverable: A 1-2 sentence response.

References

- R. Eldan and Y. Li. Tinstories: How small can language models be and still speak coherent english? arXiv preprint arXiv:2305.07759, 2023.
- A. Gokaslan, V. Cohen, E. Pavlick, and S. Tellex. Openwebtext corpus. <https://skylion007.github.io/OpenWebTextCorpus>, 2019.
- A. Holtzman, J. Buys, L. Du, M. Forbes, and Y. Choi. The curious case of neural text degeneration. arXiv preprint arXiv:1904.09751, 2019.
- A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, Ł. Kaiser, and I. Polosukhin. Attention is all you need. Advances in neural information processing systems, 30, 2017.