

CS336 Assignment 1 (basics): Building a Transformer LM

Version 1.0.6

CS336 Staff

Spring 2025

CS336 作业1（基础）：构建Transformer LM版本1.0.6

CS336 教师团队

2025年春季

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

- 1. Byte-pair encoding (BPE) tokenizer (§2)
- 2. Transformer language model (LM) (§3)
- 3. The cross-entropy loss function and the AdamW optimizer (§4)
- 4. The training loop, with support for serializing and loading model and optimizer state (§5)

What you will run

- 1. Train a BPE tokenizer on the TinyStories dataset.
- 2. Run your trained tokenizer on the dataset to convert it into a sequence of integer IDs.
- 3. Train a Transformer LM on the TinyStories dataset.
- 4. Generate samples and evaluate perplexity using the trained Transformer LM.
- 5. Train models on OpenWebText and submit your attained perplexities to a leaderboard.

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

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 Slack. When in doubt, consider if using it compromises the “from-scratch” ethos of the assignment.

¹See [PyTorch.org/docs/stable/nn.html#containers](https://pytorch.org/docs/stable/nn.html#containers) for a full list.

1 作业概述

在这个作业中，你将从头开始构建所有必要的组件来训练一个标准的Transformer语言模型（LM），并训练一些模型。

你将实现的内容

- 1. 字节对编码 (BPE) 分词器 (§ 2)
- 2. Transformer语言模型 (LM) (§ 3)
- 3. 交叉熵损失函数和AdamW优化器 (§ 4)
- 4. 训练循环，支持序列化和加载模型和优化器状态 (§ 5)

你将运行什么

- 1. 在TinyStories数据集上训练一个BPE tokenizer。
- 2. 在数据集上运行你的训练好的分词器，将其转换为一系列整数ID。
- 3. 在TinyStories数据集上训练一个Transformer LM。
- 4. 使用训练好的Transformer LM生成样本并评估困惑度。
- 5. 在OpenWebText上训练模型，并将你达到的困惑度提交到排行榜。

你可以使用的工具 我们期望你从头构建这些组件。特别是，你不得使用来自 `torch.nn`、`torch.nn.functional` 或 `torch.optim` 的任何定义，除非以下内容：

- `torch.nn.Parameter`
- `torch.nn` 中的容器类（例如，`Module`、`ModuleList`、`Sequential`等。）¹
- The `torch.optim.Optimizer` 基类

你可以使用任何其他 PyTorch 定义。如果你想要使用一个函数或类，但不确定是否允许，可以随时在 Slack 上提问。如果不确定，考虑使用它是否会损害作业的“从头开始”精神。

参考 [PyTorch.org/docs/stable/nn.html#containers](https://pytorch.org/docs/stable/nn.html#containers) 获取完整列表。

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 prohibited.

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:

github.com/stanford-cs336/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. `cs336_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:

- `writeup.pdf`: Answer all the written questions. Please typeset your responses.
- `code.zip`: Contains all the code you’ve written.

To submit to the leaderboard, submit a PR to:

github.com/stanford-cs336/assignment1-basics-leaderboard

See the `README.md` in the leaderboard repository for detailed submission instructions.

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. If you are doing the assignment with the class, you can find these files at `/data` of any non-head node machine.

If you are following along at home, you can download these files with the commands inside the `README.md`.

Low-Resource/Downscaling Tip: Init

Throughout the course’s assignment handouts, we will give advice for working through parts of the assignment with fewer or no GPU resources. For example, we will sometimes suggest **downscaling** your dataset or model size, or explain how to run training code on a MacOS integrated GPU or CPU. You’ll find these “low-resource tips” in a blue box (like this one). Even if you are an enrolled Stanford student with access to the course machines, these tips may help you iterate faster and save time, so we recommend you to read them!

低资源/降尺度提示：初始化

在整个课程的作业讲义中，我们将提供关于如何使用较少或没有 GPU 资源完成作业部分的建议。例如，我们有时会建议缩小数据集或模型大小，或者解释如何在 MacOS 集成 GPU 或 CPU 上运行训练代码。你会在蓝色框（就像这个框一样）中找到这些“低资源技巧”。即使你是已注册的斯坦福学生，可以访问课程机器，这些技巧也可能帮助你更快地迭代并节省时间，因此我们建议你阅读它们！

关于 AI 工具的声明：提示 LLM（如 ChatGPT）对于低级编程问题或关于语言模型的高级概念问题是被允许的，但直接使用它来解决问题是不被允许的。

我们强烈建议您在完成作业时禁用IDE中的AI自动补全功能（例如，光标标签、GitHub CoPilot）（尽管非AI自动补全，例如自动补全函数名是完全可以的）。我们发现AI自动补全会让您很难深入理解内容。

代码看起来像什么 所有分配代码以及这份文档都可以在 GitHub 上找到：

github.com/stanford-cs336/assignment1-basics

请 `git clone` 存储库。如果有任何更新，我们会通知你以便你可以 `git pull` 获取最新版本。

1. `cs336_basics/*`: 这是您编写代码的地方。请注意，这里没有代码——您可以随心所欲，从零开始！
2. `adapters.py`: 您的代码必须具备一组功能。对于每个功能（例如，缩放点积注意力），只需调用您的代码即可填写其实现（例如，`run_scaled_dot_product_attention`）。注意：您对 `adapters.py` 的修改不应包含任何实质性逻辑；这是粘合代码。
3. `test_*.py`: 这包含所有你必须通过的测试（例如，`test_scaled_dot_product_attention`），它将调用在 `adapters.py` 中定义的钩子。不要编辑测试文件。

如何提交 你将向 Gradescope 提交以下文件：

- `writeup.pdf`: 回答所有书面问题。请排版你的回答。
- `code.zip`: 包含你编写的所有代码。

要提交到排行榜，请向提交一个 PR 到：

github.com/stanford-cs336/assignment1-basics-leaderboard

在排行榜仓库中查看 `README.md` 以获取详细的提交说明。

从哪里获取数据集 这个作业将使用两个预处理后的数据集：TinyStories [Eldan和Li,2023] 以及OpenWebText [Gokaslan等人, 2019]。这两个数据集都是单个、大的纯文本文件。如果你和班级一起做这个作业，你可以在任何非主节点机器的 `/data` 找到这些文件。

如果你在家里跟着做，你可以用 `README.md`里面的命令下载这些文件。

Low-Resource/Downscaling Tip: Assignment 1 on Apple Silicon or CPU

With the staff solution code, we can train an LM to generate reasonably fluent text on an Apple M3 Max chip with 36 GB RAM, in under 5 minutes on Metal GPU (MPS) and about 30 minutes using the CPU. If these words don't mean much to you, don't worry! Just know that if you have a reasonably up-to-date laptop and your implementation is correct and efficient, you will be able to train a small LM that generates simple children's stories with decent fluency.

Later in the assignment, we will explain what changes to make if you are on CPU or MPS.

低资源/降尺度提示：在 Apple Silicon 或 CPU 上执行作业 1

使用工作人员的解决方案代码，我们可以在 Apple M3 Max 芯片（36 GB RAM）上训练一个 LM，在 Metal GPU（MPS）上不到 5 分钟，使用 CPU 大约需要 30 分钟。如果你对这些话不太理解，别担心！只需知道，如果你有一台比较新的笔记本电脑，并且你的实现是正确和高效的，你将能够训练一个小的 LM，它能够生成简单的故事，并且流畅度相当不错。

在作业的后面，我们将解释如果你使用 CPU 或 MPS 时需要做出哪些更改。

2 Byte-Pair Encoding (BPE) Tokenizer

In the first part of the assignment, we will train and implement a byte-level byte-pair encoding (BPE) tokenizer [Sennrich et al., 2016, Wang et al., 2019]. In particular, we will represent arbitrary (Unicode) strings as a sequence of bytes and train our BPE tokenizer on this byte sequence. Later, we will use this tokenizer to encode text (a string) into tokens (a sequence of integers) for language modeling.

2.1 The Unicode Standard

Unicode is a text encoding standard that maps characters to integer *code points*. As of Unicode 16.0 (released in September 2024), the standard defines 154,998 characters across 168 scripts. For example, the character “s” has the code point 115 (typically notated as U+0073, where U+ is a conventional prefix and 0073 is 115 in hexadecimal), and the character “牛” has the code point 29275. In Python, you can use the `ord()` function to convert a single Unicode character into its integer representation. The `chr()` function converts an integer Unicode code point into a string with the corresponding character.

```
>>> ord('牛')
29275
>>> chr(29275)
'牛'
```

Problem (unicode1): Understanding Unicode (1 point)

- (a) What Unicode character does `chr(0)` return?
Deliverable: A one-sentence response.
- (b) How does this character’s string representation (`__repr__()`) differ from its printed representation?
Deliverable: A one-sentence response.
- (c) What happens when this character occurs in text? It may be helpful to play around with the following in your Python interpreter and see if it matches your expectations:

```
>>> chr(0)
>>> print(chr(0))
>>> "this is a test" + chr(0) + "string"
>>> print("this is a test" + chr(0) + "string")
```

Deliverable: A one-sentence response.

2.2 Unicode Encodings

While the Unicode standard defines a mapping from characters to code points (integers), it’s impractical to train tokenizers directly on Unicode codepoints, since the vocabulary would be prohibitively large (around 150K items) and sparse (since many characters are quite rare). Instead, we’ll use a Unicode encoding, which converts a Unicode character into a sequence of bytes. The Unicode standard itself defines three encodings: UTF-8, UTF-16, and UTF-32, with UTF-8 being the dominant encoding for the Internet (more than 98% of all webpages).

To encode a Unicode string into UTF-8, we can use the `encode()` function in Python. To access the underlying byte values for a Python `bytes` object, we can iterate over it (e.g., call `list()`). Finally, we can use the `decode()` function to decode a UTF-8 byte string into a Unicode string.

2 字节对齐编码 (BPE) 分词器

在作业的第一部分，我们将训练并实现一个基于字节的字节对编码 (BPE) 分词器 [Sennrich et al., 2016, Wang et al., 2019]。具体来说，我们将任意 (Unicode) 字符串表示为字节序列，并在该字节序列上训练我们的 BPE 分词器。稍后，我们将使用此分词器将文本（一个字符串）编码为标记（一个整数序列）以进行语言建模。

2.1 Unicode标准

Unicode是一种文本编码标准，将字符映射到整数代码点。截至Unicode 16.0（2024年9月发布），该标准定义了168个脚本中的154,998个字符。例如，字符“s”的代码点是115（通常表示为U+0073，其中U+是一个传统的前缀，0073是十六进制的115），而字符“牛”的代码点是29275。在Python中，您可以使用`ord()`函数将单个Unicode字符转换为它的整数表示。`chr()`函数将整数Unicode代码点转换为对应字符的字符串。

```
>>> ord('牛')
29275
>>> chr(29275)
'牛'
```

问题 (unicode1): 理解Unicode (1分)

- (a) `chr(0)` 返回什么Unicode字符？
提交要求：一句话回答。
- (b) 这个字符的字符串表示 (`__repr__()`) 与其打印表示有何不同？

交付物：一句话的回应。
- (c) 当这个字符出现在文本中时会发生什么？在您的 Python 解释器中尝试以下操作可能会有所帮助，看看是否符合您的预期：

```
>>> chr(0)
>>> print(chr(0))
>>> "this is a test" + chr(0) + "string"
>>> print("this is a test" + chr(0) + "string")
```

交付物：一句话的回应。

2.2 Unicode 编码

虽然 Unicode 标准定义了从字符到代码点的映射（整数），但直接在 Unicode 代码点上训练分词器是不切实际的，因为词汇量会非常庞大（大约 150K 项）且稀疏（因为许多字符非常罕见）。相反，我们将使用 Unicode 编码，它将一个 Unicode 字符转换为一个字节序列。Unicode 标准本身定义了三种编码：UTF-8、UTF-16 和 UTF-32，其中 UTF-8 是互联网上的主要编码（所有网页中超过 98% 使用 UTF-8）。

要将 Unicode 字符串编码为 UTF-8，我们可以使用 Python 中的 `encode()` 函数。要访问 Python `bytes` 对象的底层字节值，我们可以迭代它（例如，调用 `list()`）。最后，我们可以使用 `decode()` 函数将 UTF-8 字节字符串解码为 Unicode 字符串。

```
>>> test_string = "hello! こんにちは!"
>>> utf8_encoded = test_string.encode("utf-8")
>>> print(utf8_encoded)
b'hello! \xe3\x81\x93\xe3\x82\x93\xe3\x81\xab\xe3\x81\xa1\xe3\x81\xaf!'
>>> print(type(utf8_encoded))
<class 'bytes'>
>>> # Get the byte values for the encoded string (integers from 0 to 255).
>>> list(utf8_encoded)
[104, 101, 108, 108, 111, 33, 32, 227, 129, 147, 227, 130, 147, 227, 129, 171, 227, 129, 161, 227, 129, 175, 33]
>>> # One byte does not necessarily correspond to one Unicode character!
>>> print(len(test_string))
13
>>> print(len(utf8_encoded))
23
>>> print(utf8_encoded.decode("utf-8"))
hello! こんにちは!
```

By converting our Unicode codepoints into a sequence of bytes (e.g., via the UTF-8 encoding), we are essentially taking a sequence of codepoints (integers in the range 0 to 154,997) and transforming it into a sequence of byte values (integers in the range 0 to 255). The 256-length byte vocabulary is *much* more manageable to deal with. When using byte-level tokenization, we do not need to worry about out-of-vocabulary tokens, since we know that *any* input text can be expressed as a sequence of integers from 0 to 255.

Problem (unicode2): Unicode Encodings (3 points)

(a) What are some reasons to prefer training our tokenizer on UTF-8 encoded bytes, rather than UTF-16 or UTF-32? It may be helpful to compare the output of these encodings for various input strings.
Deliverable: A one-to-two sentence response.

(b) Consider the following (incorrect) function, which is intended to decode a UTF-8 byte string into a Unicode string. Why is this function incorrect? Provide an example of an input byte string that yields incorrect results.

```
def decode_utf8_bytes_to_str_wrong(bytestring: bytes):
    return "".join([bytes([b]).decode("utf-8") for b in bytestring])
```

```
>>> decode_utf8_bytes_to_str_wrong("hello".encode("utf-8"))
'hello'
```

Deliverable: An example input byte string for which `decode_utf8_bytes_to_str_wrong` produces incorrect output, with a one-sentence explanation of why the function is incorrect.

(c) Give a two byte sequence that does not decode to any Unicode character(s).
Deliverable: An example, with a one-sentence explanation.

2.3 Subword Tokenization

While byte-level tokenization can alleviate the out-of-vocabulary issues faced by word-level tokenizers, tokenizing text into bytes results in extremely long input sequences. This slows down model training, since a

```
>>> test_string = "hello! こんにちは!"
>>> utf8_encoded = test_string.encode("utf-8")
>>> print(utf8_encoded)
b'hello! \xe3\x81\x93\xe3\x82\x93\xe3\x81\xab\xe3\x81\xa1\xe3\x81\xaf!'
>>> print(type(utf8_encoded))
<class 'bytes'>
>>> # Get the byte values for the encoded string (integers from 0 to 255).
>>> list(utf8_encoded)
[104, 101, 108, 108, 111, 33, 32, 227, 129, 147, 227, 130, 147, 227, 129, 171, 227, 129, 161, 227, 129, 175, 33]
>>> # One byte does not necessarily correspond to one Unicode character!
>>> print(len(test_string))
13
>>> print(len(utf8_encoded))
23
>>> print(utf8_encoded.decode("utf-8"))
hello! こんにちは!
```

通过将我们的 Unicode 代码点转换为字节序列（例如，通过 UTF-8 编码），我们本质上是在将一系列代码点（范围在 0 到 154,997 的整数）转换为一系列字节值（范围在 0 到 255 的整数）。256 长度的字节词汇表要容易处理得多。在使用字节级分词时，我们不需要担心词汇表外标记，因为我们知道任何输入文本都可以表示为 0 到 255 的整数序列。

问题(unicode2): Unicode编码 (3分)

(a) 有哪些原因更喜欢在UTF-8编码的字节上训练我们的分词器，而不是UTF-16或UTF-32？比较这些编码在各种输入字符串上的输出可能会有所帮助。

交付物：一到两句话的回应。

(b) 考虑以下（不正确）的函数，该函数旨在将UTF-8字节字符串解码为Unicode字符串。为什么这个函数不正确？提供一个产生不正确结果的输入字节字符串的例子。

```
def decode_utf8_bytes_to_str_wrong(bytestring: bytes):
    return "".join([bytes([b]).decode("utf-8") for b in bytestring])
```

```
>>> decode_utf8_bytes_to_str_wrong("hello".encode("utf-8"))
'hello'
```


交付物：一个输入字节字符串的例子，其中 `decode_utf8_bytes_to_str_wrong` 产生不正确的输出，并附有一句解释为什么该函数不正确的说明。

(c) 给出一个不会解码为任何Unicode字符的二字节序列。
可交付成果：一个示例，附带一句话解释。

2.3 子词分词

虽然字节级分词可以缓解词级分词面临的词汇外问题，但将文本分词成字节会导致输入序列非常长。这会减慢模型训练速度，因为

sentence with 10 words might only be 10 tokens long in a word-level language model, but could be 50 or more tokens long in a character-level model (depending on the length of the words). Processing these longer sequences requires more computation at each step of the model. Furthermore, language modeling on byte sequences is difficult because the longer input sequences create long-term dependencies in the data.

Subword tokenization is a midpoint between word-level tokenizers and byte-level tokenizers. Note that a byte-level tokenizer’s vocabulary has 256 entries (byte values are 0 to 225). A subword tokenizer trades-off a larger vocabulary size for better compression of the input byte sequence. For example, if the byte sequence `b'the'` often occurs in our raw text training data, assigning it an entry in the vocabulary would reduce this 3-token sequence to a single token.

How do we select these subword units to add to our vocabulary? Sennrich et al. [2016] propose to use byte-pair encoding (BPE; Gage, 1994), a compression algorithm that iteratively replaces (“merges”) the most frequent pair of bytes with a single, new unused index. Note that this algorithm adds subword tokens to our vocabulary to maximize the compression of our input sequences—if a word occurs in our input text enough times, it’ll be represented as a single subword unit.

Subword tokenizers with vocabularies constructed via BPE are often called BPE tokenizers. In this assignment, we’ll implement a byte-level BPE tokenizer, where the vocabulary items are bytes or merged sequences of bytes, which give us the best of both worlds in terms of out-of-vocabulary handling and manageable input sequence lengths. The process of constructing the BPE tokenizer vocabulary is known as “training” the BPE tokenizer.

2.4 BPE Tokenizer Training

The BPE tokenizer training procedure consists of three main steps.

Vocabulary initialization The tokenizer vocabulary is a one-to-one mapping from `bytestring` token to integer ID. Since we’re training a byte-level BPE tokenizer, our initial vocabulary is simply the set of all bytes. Since there are 256 possible byte values, our initial vocabulary is of size 256.

Pre-tokenization Once you have a vocabulary, you could, in principle, count how often bytes occur next to each other in your text and begin merging them starting with the most frequent pair of bytes. However, this is quite computationally expensive, since we’d have to go take a full pass over the corpus each time we merge. In addition, directly merging bytes across the corpus may result in tokens that differ only in punctuation (e.g., `dog!` vs. `dog.`). These tokens would get completely different token IDs, even though they are likely to have high semantic similarity (since they differ only in punctuation).

To avoid this, we *pre-tokenize* the corpus. You can think of this as a coarse-grained tokenization over the corpus that helps us count how often pairs of characters appear. For example, the word `'text'` might be a pre-token that appears 10 times. In this case, when we count how often the characters `'t'` and `'e'` appear next to each other, we will see that the word `'text'` has `'t'` and `'e'` adjacent and we can increment their count by 10 instead of looking through the corpus. Since we’re training a byte-level BPE model, each pre-token is represented as a sequence of UTF-8 bytes.

The original BPE implementation of Sennrich et al. [2016] pre-tokenizes by simply splitting on whitespace (i.e., `s.split(" ")`). In contrast, we’ll use a regex-based pre-tokenizer (used by GPT-2; Radford et al., 2019) from github.com/openai/tiktoken/pull/234/files:

```
>>> PAT = r'(?:[sdmt]|ll|ve|re)| ?\p{L}+| ?\p{N}+| ?[^\s\p{L}\p{N}]+\s+(?!S)|\s+''''
```

It may be useful to interactively split some text with this pre-tokenizer to get a better sense of its behavior:

```
>>> # requires `regex` package
>>> import regex as re
>>> re.findall(PAT, "some text that i'll pre-tokenize")
['some', ' ', 'text', ' ', 'that', ' ', 'i', "'", 'll', ' ', 'pre', '-', 'tokenize']
```

一个包含10个单词的句子在词级语言模型中可能只有10个token，但在字符级模型中可能长达50个或更多token（取决于单词的长度）。处理这些更长的序列需要在模型的每个步骤中进行更多的计算。此外，对字节序列进行语言建模很困难，因为更长的输入序列会在数据中创建长期依赖关系。

Subword tokenization 是 word-level tokenizer 和 byte-level tokenizer 之间的一个中间点。请注意，byte-level tokenizer 的词汇表有 256 个条目（字节值是 0 到 225）。Subword tokenizer 用更大的词汇表大小来换取对输入字节序列更好的压缩。例如，如果字节序列 `b'the'` 经常出现在我们的原始文本训练数据中，将其分配为一个词汇表条目会将这个 3 个 token 的序列减少到一个 token。

我们如何选择这些 subword 单位添加到我们的词汇表中？Sennrich 等人 [2016] 提议使用 byte-pair encoding (BPE; Gage, 1994)，一种迭代地用一个新的、未使用的索引替换（合并）最频繁的字节对的压缩算法。请注意，这个算法会向我们的词汇表中添加 subword token 以最大化输入序列的压缩——如果一个词在输入文本中出现的次数足够多，它会被表示为一个 subword 单位。

使用 BPE 构建的词汇表的子词分词器通常称为 BPE 分词器。在这个作业中，我们将实现一个字节级 BPE 分词器，其中词汇表项是字节或字节合并序列，这让我们在词汇表外处理和可管理的输入序列长度方面都能获得最佳效果。构建 BPE 分词器词汇表的过程被称为“训练”BPE 分词器。

2.4 BPE 分词器训练

BPE 分词器训练过程包括三个主要步骤。

词汇表初始化 分词器词汇表是从字节串标记到整数 ID 的一对一映射。由于我们正在训练一个字节级 BPE 分词器，我们的初始词汇表简单地就是所有字节的集合。由于有 256 种可能的字节值，我们的初始词汇表的大小为 256。

预分词 一旦你有了词汇表，原则上你可以计算字节在文本中相邻出现的频率，并从最频繁的字节对开始合并它们。然而，这相当计算密集，因为每次合并时我们都需要完整地遍历语料库。此外，直接在语料库中合并字节可能会导致仅标点符号不同的标记（例如，`dog!` 与 `dog.`）。这些标记将获得完全不同的标记 ID，即使它们很可能具有很高的语义相似性（因为它们仅标点符号不同）。

为了避免这种情况，我们预先分词语料库。你可以将其视为语料库上的粗粒度分词，这有助于我们统计字符对出现的频率。例如，单词 `'text'` 可能是一个出现 10 次的预分词。在这种情况下，当我们统计字符 `'t'` 和 `'e'` 相邻出现的频率时，我们会发现单词 `'text'` 有 `'t'` 和 `'e'` 相邻，我们可以通过 10 来增加它们的计数，而不是遍历语料库。由于我们正在训练一个字节级的 BPE 模型，每个预分词都表示为 UTF-8 字节序列。

Sennrich 等人的原始 BPE 实现通过简单地按空白字符分割（即 `s.split(" ")`）进行预分词。相比之下，我们将使用基于正则表达式的预分词器（GPT-2 使用；Radford 等人，2019 年）来自 github.com/openai/tiktoken/pull/234/files：

```
>>> PAT = r'(?:[sdmt]|ll|ve|re)| ?\p{L}+| ?\p{N}+| ?[^\s\p{L}\p{N}]+\s+(?!S)|\s+''''
```

使用此预分词器交互式地拆分一些文本可能有助于更好地理解其行为：

```
>>> # requires `regex` package
>>> import regex as re
>>> re.findall(PAT, "some text that i'll pre-tokenize")
['some', ' ', 'text', ' ', 'that', ' ', 'i', "'", 'll', ' ', 'pre', '-', 'tokenize']
```

When using it in your code, however, you should use `re.finditer` to avoid storing the pre-tokenized words as you construct your mapping from pre-tokens to their counts.

Compute BPE merges Now that we’ve converted our input text into pre-tokens and represented each pre-token as a sequence of UTF-8 bytes, we can compute the BPE merges (i.e., train the BPE tokenizer). At a high level, the BPE algorithm iteratively counts every pair of bytes and identifies the pair with the highest frequency (“A”, “B”). Every occurrence of this most frequent pair (“A”, “B”) is then *merged*, i.e., replaced with a new token “AB”. This new merged token is added to our vocabulary; as a result, the final vocabulary after BPE training is the size of the initial vocabulary (256 in our case), plus the number of BPE merge operations performed during training. For efficiency during BPE training, we do not consider pairs that cross pre-token boundaries.² When computing merges, deterministically break ties in pair frequency by *preferring the lexicographically greater pair*. For example, if the pairs (“A”, “B”), (“A”, “C”), (“B”, “ZZ”), and (“BA”, “A”) all have the highest frequency, we’d merge (“BA”, “A”):

```
>>> max([("A", "B"), ("A", "C"), ("B", "ZZ"), ("BA", "A")])
('BA', 'A')
```

Special tokens Often, some strings (e.g., `<|endoftext|>`) are used to encode metadata (e.g., boundaries between documents). When encoding text, it’s often desirable to treat some strings as “special tokens” that should never be split into multiple tokens (i.e., will always be preserved as a single token). For example, the end-of-sequence string `<|endoftext|>` should always be preserved as a single token (i.e., a single integer ID), so we know when to stop generating from the language model. These special tokens must be added to the vocabulary, so they have a corresponding fixed token ID.

Algorithm 1 of Sennrich et al. [2016] contains an inefficient implementation of BPE tokenizer training (essentially following the steps that we outlined above). As a first exercise, it may be useful to implement and test this function to test your understanding.

Example (bpe_example): BPE training example

Here is a stylized example from Sennrich et al. [2016]. Consider a corpus consisting of the following text

low low low low low
lower lower widest widest widest
newest newest newest newest newest newest

and the vocabulary has a special token `<|endoftext|>`.

Vocabulary We initialize our vocabulary with our special token `<|endoftext|>` and the 256 byte values.

Pre-tokenization For simplicity and to focus on the merge procedure, we assume in this example that pretokenization simply splits on whitespace. When we pretokenize and count, we end up with the frequency table.

{low: 5, lower: 2, widest: 3, newest: 6}

²Note that the original BPE formulation [Sennrich et al., 2016] specifies the inclusion of an end-of-word token. We do not add an end-of-word-token when training byte-level BPE models because all bytes (including whitespace and punctuation) are included in the model’s vocabulary. Since we’re explicitly representing spaces and punctuation, the learned BPE merges will naturally reflect these word boundaries.

然而，在您的code中使用它时，您应该使用 `re.finditer` 来避免在构建从预词到其计数的映射时将预分词的单词存储起来。

计算 BPE 合并 现在，我们已经将输入文本转换为预标记，并将每个预标记表示为 UTF-8 字节序列，我们可以计算 BPE 合并（即训练 BPE 分词器）。从高层次来看，BPE 算法迭代地统计每一对字节，并识别出现频率最高的那一对（“A”，“B”）。然后，将这一对最频繁出现的字节（“A”，“B”）进行合并，即用一个新的标记“AB”替换。这个新的合并标记被添加到我们的词汇表中；因此，BPE 训练后的最终词汇表大小等于初始词汇表的大小（在我们的例子中为 256），加上在训练过程中执行的 BPE 合并操作的次数。为了在 BPE 训练期间提高效率，我们不考虑跨越预标记边界的字节对。² 在计算合并时，通过优先选择字典序更大的字节对来确定性地打破频率平局。例如，如果字节对（“A”，“B”）、（“A”，“C”）、（“B”，“ZZ”）和（“BA”，“A”）都具有最高的频率，我们会合并（“BA”，“A”）：

```
>>> max([("A", "B"), ("A", "C"), ("B", "ZZ"), ("BA", "A")])
('BA', 'A')
```

特殊标记 通常，某些字符串（例如，`<|endoftext|>`）用于编码元数据（例如，文档之间的边界）。在编码文本时，通常希望将某些字符串视为“特殊标记”，这些标记永远不会被拆分成多个标记（即，始终作为一个单个标记保留）。例如，序列结束字符串 `<|endoftext|>` 始终作为一个单个标记保留（即，一个单独的整数 ID），这样我们才知道何时停止从语言模型生成。这些特殊标记必须添加到词汇表中，以便它们有一个对应的固定标记 ID。

Sennrich 等人的算法 1 [2016] 包含一个低效的 BPE 分词器训练实现（基本上遵循了我们上面概述的步骤）。作为第一个练习，实现和测试这个函数可能有助于测试你的理解。

示例 (bpe_example): BPE 训练示例

这里是一个来自 Sennrich 等人的风格化示例 [2016]。考虑一个包含以下文本的语料库

low low low low low
lower lower widest widest widest
newest newest newest newest newest newest

并且词汇表有一个特殊的标记 `<|endoftext|>`。

词汇值。 我们用我们的特殊标记 `<|endoftext|>` 和256字节初始化我们的词汇表

预标记化 为了简化并专注于合并过程，在本示例中我们假设预标记化只是按空白字符分割。当我们预标记化并计数时，我们最终得到频率表。

{low: 5, lower: 2, widest: 3, newest: 6}

²请注意，原始的BPE公式 [Sennrich等人 2016] 指定了包含一个单词结束标记。我们在训练字节级BPE模型时不添加单词结束标记，因为所有字节（包括空格和标点符号）都包含在模型的词汇表中。由于我们明确表示空格和标点符号，因此学习到的BPE合并自然会反映这些词边界。

It is convenient to represent this as a `dict[tuple[bytes], int]`, e.g. `{(1,0,w): 5 ...}`. Note that even a single byte is a `bytes` object in Python. There is no `byte` type in Python to represent a single byte, just as there is no `char` type in Python to represent a single character.

Merges We first look at every successive pair of bytes and sum the frequency of the words where they appear `{lo: 7, ow: 7, we: 8, er: 2, wi: 3, id: 3, de: 3, es: 9, st: 9, ne: 6, ew: 6}`. The pair `('es')` and `('st')` are tied, so we take the lexicographically greater pair, `('st')`. We would then merge the pre-tokens so that we end up with `{(1,0,w): 5, (1,0,w,e,r): 2, (w,i,d,e,st): 3, (n,e,w,e,st): 6}`.

In the second round, we see that `(e, st)` is the most common pair (with a count of 9) and we would merge into `{(1,0,w): 5, (1,0,w,e,r): 2, (w,i,d,est): 3, (n,e,w,est): 6}`. Continuing this, the sequence of merges we get in the end will be `['s t', 'e st', 'o w', 'l ow', 'w est', 'n e', 'ne west', 'w i', 'wi d', 'wid est', 'low e', 'lowe r']`.

If we take 6 merges, we have `['s t', 'e st', 'o w', 'l ow', 'w est', 'n e']` and our vocabulary elements would be `[<|endoftext|>, [...256 BYTE CHARS], st, est, ow, low, west, ne]`.

With this vocabulary and set of merges, the word `newest` would tokenize as `[ne, west]`.

将其表示为 `dict[tuple[bytes], int]` 很方便，例如 `{(1,0,w): 5 ...}`。请注意，即使是单个字节在 Python 中也是一个 `bytes` 对象。Python 中没有 `byte` 类型来表示单个字节，就像 Python 中没有 `char` 类型来表示单个字符一样。

合并 我们首先查看每一对连续的字节，并计算它们出现时单词的频率 `{lo: 7, ow: 7, we: 8, er: 2, wi: 3, id: 3, de: 3, es: 9, st: 9, ne: 6, ew: 6}`。对 `('es')` 和 `('st')` 是平局，所以我们选择字典序较大的对，`('st')`。然后将预标记合并，以便最终得到 `{(1,0,w): 5, (1,0,w,e,r): 2, (w,i,d,e,st): 3, (n,e,w,e,st): 6}`。

在第二轮中，我们看到 `(e, st)` 是最常见的对（计数为 9），我们将合并成 `{(1,0,w): 5, (1,0,w,e,r): 2, (w,i,d,est): 3, (n,e,w,est): 6}`。继续这样做，最终得到的合并序列将是 `['s t', 'e st', 'o w', 'l ow', 'w est', 'n e', 'ne west', 'w i', 'wi d', 'wid est', 'low e', 'lowe r']`。

如果我们合并 take6，我们将拥有 `['s t', 'e st', 'o w', 'l ow', 'w est', 'n e']`，并且我们的词汇元素将是 `[<|endoftext|>, [...256 BYTE CHARS], st, est, ow, low, west, ne]`。有了这个词汇和一系列合并，单词 `newest` 将被标记为 `[ne, west]`。

2.5 Experimenting with BPE Tokenizer Training

Let’s train a byte-level BPE tokenizer on the TinyStories dataset. Instructions to find / download the dataset can be found in Section 1. Before you start, we recommend taking a look at the TinyStories dataset to get a sense of what’s in the data.

Parallelizing pre-tokenization You will find that a major bottleneck is the pre-tokenization step. You can speed up pre-tokenization by parallelizing your code with the built-in library `multiprocessing`. Concretely, we recommend that in parallel implementations of pre-tokenization, you chunk the corpus while ensuring your chunk boundaries occur at the beginning of a special token. You are free to use the starter code at the following link verbatim to obtain chunk boundaries, which you can then use to distribute work across your processes:

https://github.com/stanford-cs336/assignment1-basics/blob/main/cs336_basics/pretokenization_example.py

This chunking will always be valid, since we never want to merge across document boundaries. For the purposes of the assignment, you can always split in this way. Don’t worry about the edge case of receiving a very large corpus that does not contain `<|endoftext|>`.

Removing special tokens before pre-tokenization Before running pre-tokenization with the regex pattern (using `re.finditer`), you should strip out all special tokens from your corpus (or your chunk, if using a parallel implementation). Make sure that you **split** on your special tokens, so that no merging can occur across the text they delimit. For example, if you have a corpus (or chunk) like `[Doc 1]<|endoftext|>[Doc 2]`, you should split on the special token `<|endoftext|>`, and pre-tokenize `[Doc 1]` and `[Doc 2]` separately, so that no merging can occur across the document boundary. This can be done using `re.split` with `"|"`. `.join(special_tokens)` as the delimiter (with careful use of `re.escape` since `|` may occur in the special tokens). The test `test_train_bpe_special_tokens` will test for this.

Optimizing the merging step The naïve implementation of BPE training in the stylized example above is slow because for every merge, it iterates over all byte pairs to identify the most frequent pair. However, the only pair counts that change after each merge are those that overlap with the merged pair. Thus, BPE training speed can be improved by indexing the counts of all pairs and incrementally updating these counts, rather than explicitly iterating over each pair of bytes to count pair frequencies. You can get significant speedups with this caching procedure, though we note that the merging part of BPE training is *not* parallelizable in Python.

2.5 尝试使用 BPE 分词器训练

让我们在 TinyStories 数据集上训练一个字节级的 BPE 分词器。可以在第 1 节中找到查找/下载数据集的说明。在开始之前，我们建议您查看 TinyStories 数据集，以了解数据中包含的内容。

并行化预分词 您会发现主要的瓶颈是预分词步骤。您可以通过使用内置库 `multiprocessing` 并行化您的代码来加快预分词速度。具体来说，我们建议在预分词的并行实现中，在确保您的块边界出现在特殊标记的开头时，对语料库进行分块。您可以原封不动地使用以下链接中的启动代码来获取块边界，然后使用这些边界将工作分配到您的进程：

https://github.com/stanford-cs336/assignment1-basics/blob/main/cs336_basics/pretokenization_example.py

这种分块永远有效，因为我们永远不会跨文档边界合并。对于作业的目的，你可以始终以这种方式分割。不用担心收到一个非常大的语料库，其中不包含 `<|endoftext|>` 的边缘情况。

在预分词之前移除特殊标记 在使用 `re.finditer` 的正则表达式模式运行预分词之前，你应该从你的语料库（或如果你的使用并行实现，则是你的分块）中剥离所有特殊标记。确保你在你的特殊标记上分割，这样它们分隔的文本之间就不会发生合并。例如，如果你有一个语料库（或分块）像 `[Doc 1]<|endoftext|>[Doc 2]`，你应该在特殊标记 `<|endoftext|>` 上分割，并分别预分词 `[Doc 1]` 和 `[Doc 2]`，这样就不会在文档边界发生合并。这可以使用

`re.split` 和 `"|"` 作为分隔符来完成（小心使用，因为特殊标记中可能会出现 `since`）。测试 `test_train_bpe_special_tokens` 将会测试这一点。

优化合并步骤 上述样例中BPE训练的朴素实现之所以慢，是因为对于每次合并，它都会遍历所有字节对来识别最频繁的对。然而，每次合并后唯一会改变的计数对是与合并对重叠的那些。因此，可以通过索引所有对的计数并增量更新这些计数来提高BPE训练速度，而不是显式地遍历每个字节对来计算对频率。使用这种缓存过程可以获得显著的加速，尽管我们注意到Python中BPE训练的合并部分是不可并行化的。

Low-Resource/Downscaling Tip: Profiling

You should use profiling tools like `cProfile` or `scalene` to identify the bottlenecks in your implementation, and focus on optimizing those.

Low-Resource/Downscaling Tip: “Downscaling”

Instead of jumping to training your tokenizer on the full TinyStories dataset, we recommend you first train on a small subset of the data: a “debug dataset”. For example, you could train your tokenizer on the TinyStories validation set instead, which is 22K documents instead of 2.12M. This illustrates a general strategy of downscaling whenever possible to speed up development: for example, using smaller datasets, smaller model sizes, etc. Choosing the size of the debug dataset or hyperparameter config requires careful consideration: you want your debug set to be large enough to have the same bottlenecks as the full configuration (so that the optimizations you make will generalize), but not so big that it takes forever to run.

Problem (train_bpe): BPE Tokenizer Training (15 points)

Deliverable: Write a function that, given a path to an input text file, trains a (byte-level) BPE tokenizer. Your BPE training function should handle (at least) the following input parameters:

- input_path:** `str` Path to a text file with BPE tokenizer training data.
- vocab_size:** `int` A positive integer that defines the maximum final vocabulary size (including the initial byte vocabulary, vocabulary items produced from merging, and any special tokens).
- special_tokens:** `list[str]` A list of strings to add to the vocabulary. These special tokens do not otherwise affect BPE training.

Your BPE training function should return the resulting vocabulary and merges:

- vocab:** `dict[int, bytes]` The tokenizer vocabulary, a mapping from `int` (token ID in the vocabulary) to `bytes` (token bytes).
- merges:** `list[tuple[bytes, bytes]]` A list of BPE merges produced from training. Each list item is a `tuple` of `bytes` (`<token1>`, `<token2>`), representing that `<token1>` was merged with `<token2>`. The merges should be ordered by order of creation.

To test your BPE training function against our provided tests, you will first need to implement the test adapter at `[adapters.run_train_bpe]`. Then, run `uv run pytest tests/test_train_bpe.py`. Your implementation should be able to pass all tests. Optionally (this could be a large time-investment), you can implement the key parts of your training method using some systems language, for instance C++ (consider `cppyy` for this) or Rust (using `PyO3`). If you do this, be aware of which operations require copying vs reading directly from Python memory, and make sure to leave build instructions, or make sure it builds using only `pyproject.toml`. Also note that the GPT-2 regex is not well-supported in most regex engines and will be too slow in most that do. We have verified that `Oniguruma` is reasonably fast and supports negative lookahead, but the `regex` package in Python is, if anything, even faster.

低资源/降采样提示：分析

你应该使用 `cProfile` 或 `scalene` 等分析工具来识别你实现中的瓶颈，并专注于优化这些部分。

低资源/降采样提示：“降采样”

与其直接在完整的 TinyStories 数据集上训练你的分词器，我们建议你首先在数据的一个小子集上训练：一个“调试数据集”。例如，你可以选择在 TinyStories 验证集上训练你的分词器，该验证集有 22K 个文档而不是 2.12M。这说明了在可能的情况下使用降采样策略来加速开发的一般方法：例如，使用较小的数据集、较小的模型大小等。选择调试数据集的大小或超参数配置需要仔细考虑：你希望你的调试集足够大，以便与完整配置具有相同的瓶颈（这样你做出的优化才能泛化），但又不能太大以至于运行时间过长。

问题 (train_bpe)：BPE 分词器训练 (15 分)

交付物：编写一个函数，给定输入文本文件的路径，训练一个（字节级的）BPE tokenizer。您的 BPE 训练函数应该处理（至少）以下输入参数：

- input_path:** `str` BPE 分词器训练数据的文本文件路径。
- vocab_size:** `int` 一个正整数，用于定义最终词汇表的最大大小（包括初始字节词汇表、通过合并产生的词汇表项以及任何特殊标记）。初始字节词汇表、通过合并产生的词汇表项以及任何特殊标记）。
- special_tokens:** `list[str]` 一个字符串列表，用于添加到词汇表中。这些特殊标记不会影响 BPE 训练。`{v}` 否则不会影响 BPE 训练。

您的 BPE 训练函数应返回生成的词汇和合并结果：

- vocab:** `dict[int, bytes]` 分词器词汇表，一个从 `int` (词汇表中的 token ID) 到 `bytes` (token 字节)。
- merges:** `list[tuple[bytes, bytes]]` 训练生成的 BPE 合并列表。每个列表项是一个 `tuple`，表示 `<token1>` 与 `<token2>` 合并。合并应按创建顺序排序。

要测试您的 BPE 训练函数以针对我们提供的测试，您首先需要在 `[adapters.run_train_bpe]` 实现测试适配器。然后，运行 `uv run pytest tests/test_train_bpe.py`。您的实现应该能够通过所有测试。可选的（这可能是一个巨大的时间投资），您可以使用一些系统语言实现您的训练方法的关键部分，例如 C++（考虑 `cppyy` 用于此）或 Rust（使用 `PyO3`）。如果您这样做，请注意哪些操作需要复制而不是直接从 Python 内存读取，并确保留下构建说明，或者确保它仅使用 `pyproject.toml` 构建。还请注意，GPT-2 正则表达式在大多数正则表达式引擎中不受良好支持，并且在大多数情况下都太慢了。我们已经验证 `Oniguruma` 是合理快速的，并支持负向前瞻，但 Python 中的 `regex` 包，如果 anything，甚至更快。

Problem (train_bpe_tinystories): BPE Training on TinyStories (2 points)

- (a) Train a byte-level BPE tokenizer on the TinyStories dataset, using a maximum vocabulary size of 10,000. Make sure to add the TinyStories `<|endoftext|>` special token to the vocabulary. Serialize the resulting vocabulary and merges to disk for further inspection. How many hours and memory did training take? What is the longest token in the vocabulary? Does it make sense?

Resource requirements: ≤ 30 minutes (no GPUs), ≤ 30 GB RAM

Hint You should be able to get under 2 minutes for BPE training using `multiprocessing` during pretokenization and the following two facts:

- (a) The `<|endoftext|>` token delimits documents in the data files.
- (b) The `<|endoftext|>` token is handled as a special case before the BPE merges are applied.

Deliverable: A one-to-two sentence response.

- (b) Profile your code. What part of the tokenizer training process takes the most time?

Deliverable: A one-to-two sentence response.

Next, we'll try training a byte-level BPE tokenizer on the OpenWebText dataset. As before, we recommend taking a look at the dataset to better understand its contents.

Problem (train_bpe_expts_owt): BPE Training on OpenWebText (2 points)

- (a) Train a byte-level BPE tokenizer on the OpenWebText dataset, using a maximum vocabulary size of 32,000. Serialize the resulting vocabulary and merges to disk for further inspection. What is the longest token in the vocabulary? Does it make sense?

Resource requirements: ≤ 12 hours (no GPUs), ≤ 100 GB RAM

Deliverable: A one-to-two sentence response.

- (b) Compare and contrast the tokenizer that you get training on TinyStories versus OpenWebText.

Deliverable: A one-to-two sentence response.

2.6 BPE Tokenizer: Encoding and Decoding

In the previous part of the assignment, we implemented a function to train a BPE tokenizer on input text to obtain a tokenizer vocabulary and a list of BPE merges. Now, we will implement a BPE tokenizer that loads a provided vocabulary and list of merges and uses them to encode and decode text to/from token IDs.

2.6.1 Encoding text

The process of encoding text by BPE mirrors how we train the BPE vocabulary. There are a few major steps.

Step 1: Pre-tokenize. We first pre-tokenize the sequence and represent each pre-token as a sequence of UTF-8 bytes, just as we did in BPE training. We will be merging these bytes within each pre-token into vocabulary elements, handling each pre-token independently (no merges across pre-token boundaries).

Step 2: Apply the merges. We then take the sequence of vocabulary element merges created during BPE training, and apply it to our pre-tokens *in the same order of creation*.

问题 (train_bpe_tinystories): 在 TinyStories 上进行 BPE 训练 (2 分)

- (a) 在 TinyStories 数据集上训练一个字节级 BPE 分词器，使用最大词汇量为 10,000。确保将 TinyStories `<|endoftext|>` 特殊标记添加到词汇表中。将生成的词汇表和合并规则序列化到磁盘以供进一步检查。训练花费了多少小时和内存？词汇表中最长的标记是什么？这有意义吗？

资源需求： ≤ 30 分钟（不使用 GPU）， ≤ 30 GB RAM

提示 你应该能够在 `multiprocessing` 期间使用 `multiprocessing` 将 BPE 训练时间控制在 2 分钟以内
预标记化和以下两个事实：

- (a) The `<|endoftext|>` token 在数据文件中分隔文档。
- (b) The `<|endoftext|>` token 在应用 BPE 合并规则之前被视为特殊情况进行处理。

可交付成果：一段至两段简短回应。

- (b) 分析你的代码。在分词器训练过程中，哪个部分耗时最长？

可交付成果：一段至两段简短回应。

接下来，我们将尝试在 OpenWebText 数据集上训练一个字节级 BPE 分词器。和之前一样，我们建议你查看一下数据集，以便更好地理解其内容。

问题 (train_bpe_expts_owt): 在 OpenWebText 上的 BPE 训练 (2 分)

- (a) 在 OpenWebText 数据集上训练一个字节级 BPE 分词器，最大词汇量为 32,000。将生成的词汇表和合并规则序列化到磁盘上以便进一步检查。词汇表中最长的 token 是什么？这合理吗？

资源需求： ≤ 12 小时（无 GPU）， ≤ 100 GB RAM

交付物：一段一到两句话的回应。

- (b) 比较和对比你在 TinyStories 上训练的 tokenizer 与 OpenWebText 上的 tokenizer。

交付物：一段一到两句话的回应。

2.6 BPE Tokenizer: 编码和解码

在作业的前一部分，我们实现了一个函数，用于在输入文本上训练一个 BPE tokenizer 以获得 tokenizer 词汇表和一组 BPE 合并。现在，我们将实现一个 BPE tokenizer，它加载提供的词汇表和合并列表，并使用它们将文本编码和解码为/从 token ID。

2.6.1 编码文本

BPE 对文本进行编码的过程与训练 BPE 词汇表的方式相同。主要有几个步骤。

步骤 1：预分词。我们首先对序列进行预分词，并将每个预分词表示为 UTF-8 字节序列，就像在 BPE 训练中所做的那样。我们将在这些预分词内部的字节进行合并，形成词汇表元素，独立处理每个预分词（不跨越预分词边界进行合并）。

步骤 2：应用合并。然后我们获取 BPE 训练期间创建的词汇元素合并序列，并按照创建顺序将其应用于我们的预标记。

Example (bpe_encoding): BPE encoding example

For example, suppose our input string is 'the cat ate', our vocabulary is {0: b' ', 1: b'a', 2: b'c', 3: b'e', 4: b'h', 5: b't', 6: b'th', 7: b' c', 8: b' a', 9: b'the', 10: b' at'}, and our learned merges are [(b't', b'h'), (b' ', b'c'), (b' ', 'a'), (b'th', b'e'), (b' a', b't')]. First, our pre-tokenizer would split this string into ['the', ' cat', ' ate']. Then, we'll look at each pre-token and apply the BPE merges. The first pre-token 'the' is initially represented as [b't', b'h', b'e']. Looking at our list of merges, we identify the first applicable merge to be (b't', b'h'), and use that to transform the pre-token into [b'th', b'e']. Then, we go back to the list of merges and identify the next applicable merge to be (b'th', b'e'), which transforms the pre-token into [b'the']. Finally, looking back at the list of merges, we see that there are no more that apply to the string (since the entire pre-token has been merged into a single token), so we are done applying the BPE merges. The corresponding integer sequence is [9]. Repeating this process for the remaining pre-tokens, we see that the pre-token ' cat' is represented as [b' c', b'a', b't'] after applying the BPE merges, which becomes the integer sequence [7, 1, 5]. The final pre-token ' ate' is [b' at', b'e'] after applying the BPE merges, which becomes the integer sequence [10, 3]. Thus, the final result of encoding our input string is [9, 7, 1, 5, 10, 3].

Special tokens. Your tokenizer should be able to properly handle user-defined special tokens when encoding text (provided when constructing the tokenizer).

Memory considerations. Suppose we want to tokenize a large text file that we cannot fit in memory. To efficiently tokenize this large file (or any other stream of data), we need to break it up into manageable chunks and process each chunk in-turn, so that the memory complexity is constant as opposed to linear in the size of the text. In doing so, we need to make sure that a token doesn't cross chunk boundaries, else we'll get a different tokenization than the naïve method of tokenizing the entire sequence in-memory.

2.6.2 Decoding text

To decode a sequence of integer token IDs back to raw text, we can simply look up each ID's corresponding entries in the vocabulary (a byte sequence), concatenate them together, and then decode the bytes to a Unicode string. Note that input IDs are not guaranteed to map to valid Unicode strings (since a user could input any sequence of integer IDs). In the case that the input token IDs do not produce a valid Unicode string, you should replace the malformed bytes with the official Unicode replacement character U+FFFD.³ The `errors` argument of `bytes.decode` controls how Unicode decoding errors are handled, and using `errors='replace'` will automatically replace malformed data with the replacement marker.

Problem (tokenizer): Implementing the tokenizer (15 points)

Deliverable: Implement a `Tokenizer` class that, given a vocabulary and a list of merges, encodes text into integer IDs and decodes integer IDs into text. Your tokenizer should also support user-provided special tokens (appending them to the vocabulary if they aren't already there). We recommend the following interface:

```
def __init__(self, vocab, merges, special_tokens=None) Construct a tokenizer from a given vocabulary, list of merges, and (optionally) a list of special tokens. This function should accept
```

³See [en.wikipedia.org/wiki/Specials_\(Unicode_block\)#Replacement_character](https://en.wikipedia.org/wiki/Specials_(Unicode_block)#Replacement_character) for more information about the Unicode replacement character.

示例 (bpe_encoding): BPE 编码示例

例如，假设我们的输入字符串是 'the cat ate'，我们的词汇是 {0: b' ', 1: b'a', 2:b'c', 3: b'e', 4: b'h', 5: b't', 6: b'th', 7: b' c', 8: b' a', 9: b'the', 10: b'at'}，我们的学习合并是 [(b't', b'h'), (b' ', b'c'), (b' ', 'a'), (b'th', b'e'),(b' a', b't')]. 首先，我们的预标记器会将这个字符串分割成 ['the', ' cat', ' ate']。然后，我们将查看每个预标记并应用 BPE 合并。第一个预标记 'the' 初始表示为 [b't', b'h', b'e']。查看我们的合并列表，我们识别出第一个适用的合并为 (b't', b'h')，并使用它将预标记转换为 [b'th', b'e']。然后，我们回到合并列表，识别出下一个适用的合并为 (b'th', b'e')，它将预标记转换为 [b'the']。最后，查看合并列表，我们看到没有更多适用于该字符串（因为整个预标记已经合并为一个标记），因此我们完成 BPE 合并的应用。相应的整数序列是 [9]。对剩余的预标记重复此过程，我们看到预标记 ' cat' 在应用 BPE 合并后表示为 [b' c', b'a', b't']，它成为整数序列 [7, 1,5]。最后的预标记 ' ate' 在应用 BPE 合并后是 [b' at', b'e']，它成为整数序列 [10, 3]。因此，我们输入字符串编码的最终结果是 [9, 7, 1, 5, 10,3]。

特殊标记。你的分词器应该能够正确处理用户定义的特殊标记，在编码文本时（在构建分词器时提供）。

内存考虑。假设我们想要分词一个无法放入内存的大型文本文件。为了高效地分词这个大文件（或任何其他数据流），我们需要将其拆分成可管理的块，并逐个处理每个块，这样内存复杂度是常数，而不是文本大小的线性关系。在此过程中，我们需要确保一个标记不会跨越块边界，否则我们将得到与内存中分词整个序列的朴素方法不同的分词结果。

2.6.2 解码文本

要将一系列整数标记ID解码回原始文本，我们可以简单地查找词汇表中每个ID对应的条目（一个字节序列），将它们连接在一起，然后解码字节为Unicode字符串。请注意，输入ID不保证映射到有效的Unicode字符串（因为用户可以输入任何整数ID序列）。如果输入标记ID无法生成有效的Unicode字符串，您应该用官方的Unicode替换字符U+FFFD.³。 `errors` 参数控制 `bytes.decode` 如何处理Unicode解码错误，使用 `errors='replace'` 将自动用替换标记替换损坏的数据。

问题 (tokenizer): 实现分词器 (15 分)

交付物：实现一个 `Tokenizer` 类，给定一个词汇表和一系列合并操作，将文本编码为整数 ID 并将整数 ID 解码为文本。您的分词器还应支持用户提供特殊标记（如果它们尚未存在于词汇表中，则将其追加到词汇表中）。我们建议以下接口：

```
def __init__(self, vocab, merges, special_tokens=None) 根据给定的词汇表、合并列表和（可选的）特殊标记列表构建一个分词器。此函数应接受
```

³请参阅 [en.wikipedia.org/wiki/Specials_\(Unicode_block\)#Replacement_character](https://en.wikipedia.org/wiki/Specials_(Unicode_block)#Replacement_character) 了解有关 Unicode 替换字符的更多信息。

the following parameters:

```
vocab: dict[int, bytes]
merges: list[tuple[bytes, bytes]]
special_tokens: list[str] | None = None
```

`def from_files(cls, vocab_filepath, merges_filepath, special_tokens=None)` Class method that constructs and return a `Tokenizer` from a serialized vocabulary and list of merges (in the same format that your BPE training code output) and (optionally) a list of special tokens. This method should accept the following additional parameters:

```
vocab_filepath: str
merges_filepath: str
special_tokens: list[str] | None = None
```

`def encode(self, text: str) -> list[int]` Encode an input text into a sequence of token IDs.

`def encode_iterable(self, iterable: Iterable[str]) -> Iterator[int]` Given an iterable of strings (e.g., a Python file handle), return a generator that lazily yields token IDs. This is required for memory-efficient tokenization of large files that we cannot directly load into memory.

`def decode(self, ids: list[int]) -> str` Decode a sequence of token IDs into text.

To test your `Tokenizer` against our provided tests, you will first need to implement the test adapter at `[adapters.get_tokenizer]`. Then, run `uv run pytest tests/test_tokenizer.py`. Your implementation should be able to pass all tests.

以下参数:

```
vocab: dict[int, bytes]
merges: list[tuple[bytes, bytes]]
special_tokens: list[str] | None = None
```

`def from_files(cls, vocab_filepath, merges_filepath, special_tokens=None)` Class method that constructs and return a `Tokenizer` from a serialized vocabulary and list of merges (in the same format that your BPE training code output) and (optionally) a list of special 标记。此方法应接受以下附加参数:

```
vocab_filepath: str
merges_filepath: str
special_tokens: list[str] | None = None
```

`def encode(self, text: str) -> list[int]` 将输入文本编码为一系列标记ID。

`def encode_iterable(self, iterable: Iterable[str]) -> Iterator[int]` Given an iterable of 字符串 (例如, Python文件句柄), 返回一个生成器, 该生成器惰性生成标记ID。这对于内存高效的标记化大型文件是必需的, 这些文件我们不能直接加载到内存中。

`def decode(self, ids: list[int]) -> str` 将一系列标记ID解码为文本。

要测试您的 `Tokenizer` 以我们的提供的测试, 您首先需要在 `[adapters.get_tokenizer]`实现测试适配器。然后, 运行 `uv run pytest tests/test_tokenizer.py`。您的实现应该能够通过所有测试。

2.7 Experiments

Problem (tokenizer_experiments): Experiments with tokenizers (4 points)

- (a) Sample 10 documents from TinyStories and OpenWebText. Using your previously-trained TinyStories and OpenWebText tokenizers (10K and 32K vocabulary size, respectively), encode these sampled documents into integer IDs. What is each tokenizer’s compression ratio (bytes/token)?
Deliverable: A one-to-two sentence response.
- (b) What happens if you tokenize your OpenWebText sample with the TinyStories tokenizer? Compare the compression ratio and/or qualitatively describe what happens.
Deliverable: A one-to-two sentence response.
- (c) Estimate the throughput of your tokenizer (e.g., in bytes/second). How long would it take to tokenize the Pile dataset (825GB of text)?
Deliverable: A one-to-two sentence response.
- (d) Using your TinyStories and OpenWebText tokenizers, encode the respective training and development datasets into a sequence of integer token IDs. We’ll use this later to train our language model. We recommend serializing the token IDs as a NumPy array of datatype `uint16`. Why is `uint16` an appropriate choice?

2.7 实验

问题 (tokenizer_experiments): 对分词器进行实验 (4分)

- (a) 从 TinyStories 和 OpenWebText 中采样10个文档。使用你之前训练的 TinyStories 和 OpenWebText 分词器 (词汇量分别为10K和32K), 将这些采样文档编码为整数ID。每个分词器的压缩率 (字节/token) 是多少?
交付物: 一到两句话的回应。
- (b) 如果你用 TinyStories 分词器对 OpenWebText 采样进行分词会发生什么? 比较压缩率, 或者定性地描述发生的情况。
交付物: 一到两句话的回应。
- (c) 估计你的分词器的吞吐量 (例如, 以字节/秒为单位)。分词 Pile 数据集 (825GB文本) 需要多长时间?
可交付成果: 一段一到两句话的回应。
- (d) 使用你的 TinyStories 和 OpenWebText 分词器, 将相应的训练和开发数据集编码为一系列整数分词ID。我们稍后将使用这些ID来训练我们的语言模型。我们建议将分词ID序列化为一个数据类型为 `uint16` 的 NumPy 数组。为什么 `uint16` 是一个合适的选择?

Deliverable: A one-to-two sentence response.

可交付成果：一段一到两句话的回应。

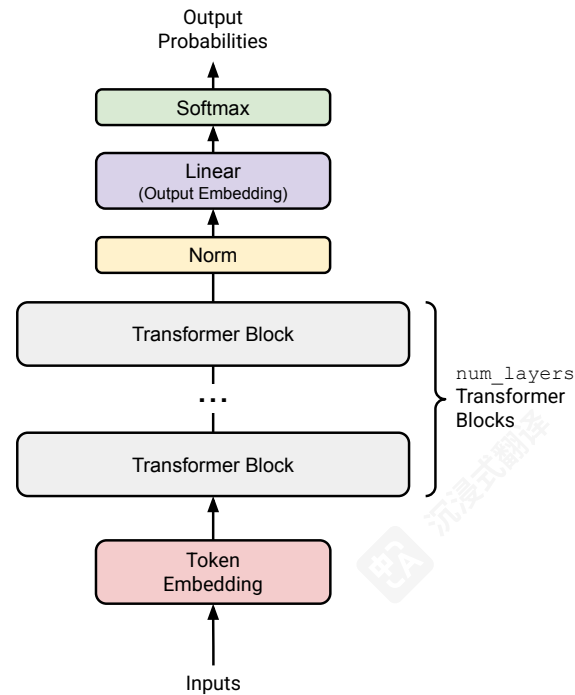


Figure 1: An overview of our Transformer language model.

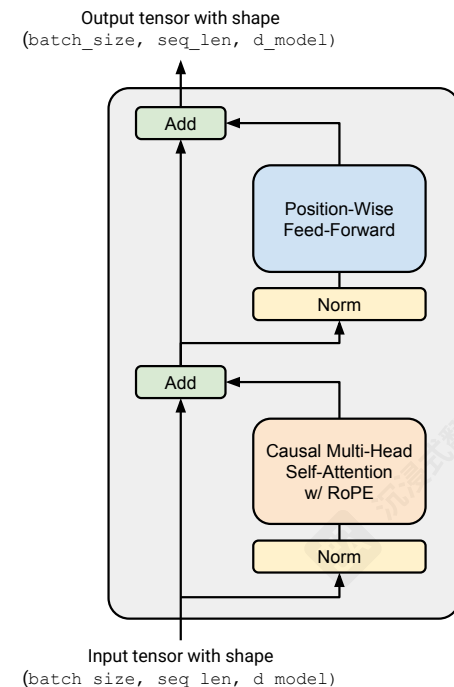


Figure 2: A pre-norm Transformer block.

3 Transformer Language Model Architecture

A language model takes as input a batched sequence of integer token IDs (i.e., `torch.Tensor` of shape `(batch_size, sequence_length)`), and returns a (batched) normalized probability distribution over the vocabulary (i.e., a PyTorch Tensor of shape `(batch_size, sequence_length, 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).

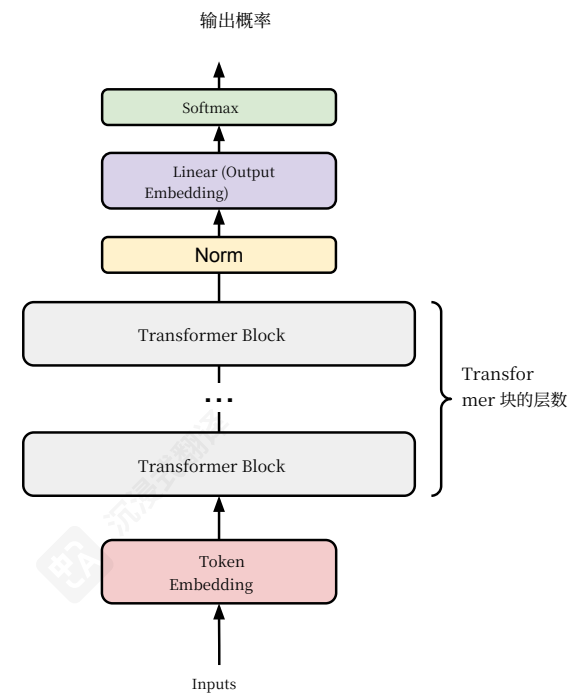


图1: 我们Transformer语言模型的整体概述。

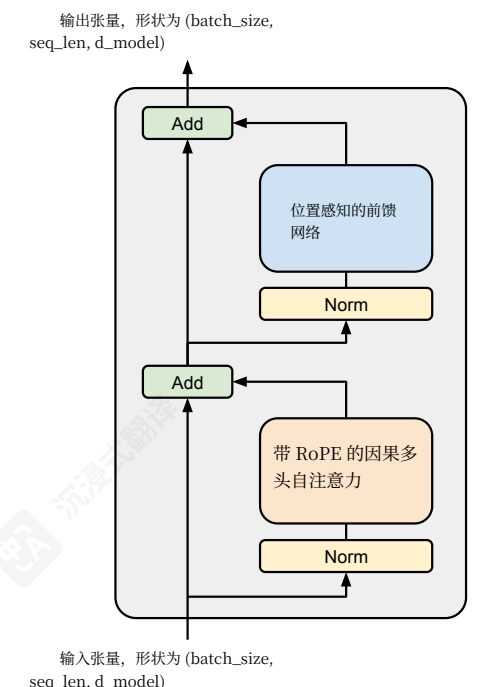


图2: 一个预归一化Transformer模块。

3 Transformer语言模型架构

一个语言模型以一批整数token ID（即 `torch.Tensor`，形状为`(batch_size, sequence_length)`）作为输入，并返回一个（批量的）归一化词汇概率分布（即形状为

`(batch_size, sequence_length, vocab_size)`的PyTorch张量），其中预测分布是针对每个输入token的下一个词。在训练语言模型时，我们使用这些下一个词的预测来计算实际下一个词和预测下一个词之间的交叉熵损失。在推理期间从语言模型生成文本时，我们取最终时间步的预测下一个词分布（即序列中的最后一个项）来生成序列中的下一个token（例如，通过选择概率最高的token、从分布中采样等），将生成的token添加到输入序列中，并重复。

在这个作业的部分，你将从头开始构建这个Transformer语言模型。我们将先对模型进行高层描述，然后再逐步详细说明各个组件。

3.1 Transformer LM

给定一个token ID序列，Transformer语言模型使用输入嵌入将token ID转换为密集向量，将嵌入的token通过`num_layers`Transformer模块，然后应用一个学习的线性投影（“输出嵌入”或“LM头”）来生成预测的下一个token的logits。参见图1以了解其示意图。

3.1.1 Token 嵌入

在第一步中，Transformer将（批处理的）token ID序列嵌入到包含token身份信息的向量序列中（图1中的红色块）。

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, $D @ A$ 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

更具地说，给定一个 token ID 序列，Transformer 语言模型使用一个 token 嵌入层来生成一个向量序列。每个嵌入层接收一个形状为 `(batch_size, sequence_length)` 的整数张量，并生成一个形状为 `(batch_size, sequence_length, d_model)` 的向量序列。

3.1.2 预规范化 TransformerBlock

嵌入后，激活值由多个结构相同的神经网络层进行处理。一个标准的仅解码器 Transformer 语言模型由 `num_layers` 个相同的层（通常称为 Transformer “块”）组成。每个 Transformer 块接收形状为 `(batch_size, sequence_length, d_model)` 的输入，并返回形状为 `(batch_size, sequence_length, d_model)` 的输出。每个块通过自注意力机制聚合序列中的信息，并通过前馈层进行非线性转换。

3.2 输出规范化和嵌入

经过 `num_layers` 个 Transformer 块后，我们将最终激活值转换为词汇表上的分布。

我们将实现“预规范化”Transformer 块（详见 § 3.5），该块额外要求在最终 Transformer 块之后使用层规范化（详见下文），以确保其输出被适当缩放。

经过此规范化后，我们将使用标准的线性变换将 Transformer 块的输出转换为预测的下一个 token logit（例如，参见 Radford 等人 [2018] 公式 2）。

3.3 备注：批处理、Einsum 和高效计算

在整个 Transformer 中，我们将对许多批处理类输入执行相同的计算。这里有几个例子：

- 批处理的元素：我们对每个批处理元素应用相同的 Transformer `forward` 操作。
- 序列长度：像 RMSNorm 和前馈这样的“位置”操作在每个序列位置上完全相同地操作在每个序列的位置上。
- 注意力头：注意力操作在“多头”注意力头之间进行批处理注意力操作。

以一种充分利用 GPU 且易于阅读和理解的方式执行此类操作很有用。许多 PyTorch 操作可以在张量的开头接收额外的“批处理”维度，并在这些维度上高效地重复/广播操作。

例如，假设我们正在进行一个逐位置、批处理的操作。我们有一个形状为 `(batch_size, sequence_length, d_model)` 的“数据张量” D ，并且我们希望对形状为 `(d_model, d_model)` 的矩阵 A 进行批处理的向量-矩阵乘法。在这种情况下， $D @ A$ 将执行批处理矩阵乘法，这是 PyTorch 中的一种高效原语，其中 `(batch_size, sequence_length)` 维度是批处理过的。

由于这个原因，假设你的函数可能会被赋予额外的批处理维度，并将这些维度保持在 PyTorch 形状的开头是有帮助的。为了组织张量以便以这种方式进行批处理，它们可能需要使用多个步骤的 `view`，`reshape` 和 `transpose` 进行塑形。这可能有点麻烦，而且通常很难阅读代码在做什么以及你的张量的形状是什么。

一种更符合人体工程学的选择是在 `torch.einsum` 中使用 `einsum` 语法，或者更确切地说使用与框架无关的库，如 `einops` 或 `einx`。这两个关键操作是 `einsum`，它可以对具有任意维度的输入张量进行张量收缩，以及 `rearrange`，它可以重新排序、连接和拆分任意

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).⁴ 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_example1): Batched matrix multiplication with einops.einsum

```
import torch
from einops import rearrange, einsum

## Basic implementation
Y = D @ A.T
# Hard to tell the input and output shapes and what they mean.
# What shapes can D and A have, and do any of these have unexpected behavior?

## Einsum is self-documenting and robust
#
#           D           A           ->           Y
Y = einsum(D, A, "batch sequence d_in, d_out d_in -> batch sequence d_out")

## Or, a batched version where D can have any leading dimensions but A is constrained.
Y = einsum(D, A, "... d_in, d_out d_in -> ... d_out")
```

Example (einstein_example2): Broadcasted operations with einops.rearrange

We have a batch of images, and for each image we want to generate 10 dimmed versions based on some scaling factor:

```
images = torch.randn(64, 128, 128, 3) # (batch, height, width, channel)
dim_by = torch.linspace(start=0.0, end=1.0, steps=10)

## Reshape and multiply
dim_value = rearrange(dim_by, "dim_value -> 1 dim_value 1 1 1")
images_rearr = rearrange(images, "b height width channel -> b 1 height width channel")
dimmed_images = images_rearr * dim_value

## Or in one go:
dimmed_images = einsum(
    images, dim_by,
    "batch height width channel, dim_value -> batch dim_value height width channel"
)
```

⁴It’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.

维度。结果表明，机器学习中的几乎所有操作都是维度操作和张量收缩的组合，偶尔（通常是逐点的）非线性函数。这意味着当使用 einsum 语法时，你的代码可以更易读和灵活。

我们强烈建议学习并使用 einsum 语法。之前没有接触过 einsum 语法的同学可以使用 einops (文档在这里)，而已经熟悉 einops 的同学应该学习更通用的 einx (这里)。⁴ 这两个包已经安装在我们提供的环境中。

这里我们给出一些 einsum 语法的使用示例。这些是 einops 文档的补充，你应该先阅读。

示例 (einstein_example1): 使用 einops.einsum 的批量矩阵乘法

```
import torch
from einops import rearrange, einsum

## Basic implementation
Y = D @ A.T
# Hard to tell the input and output shapes and what they mean.
# What shapes can D and A have, and do any of these have unexpected behavior?

## Einsum is self-documenting and robust
#
#           D           A           ->           Y
Y = einsum(D, A, "batch sequence d_in, d_out d_in -> batch sequence d_out")

## Or, a batched version where D can have any leading dimensions but A is constrained.
Y = einsum(D, A, "... d_in, d_out d_in -> ... d_out")
```

示例 (einstein_example2): 使用 einops.rearrange 的广播操作

我们有一批图像，对于每张图像，我们希望根据某个缩放因子生成 10 个变暗的版本：

```
images = torch.randn(64, 128, 128, 3) # (batch, height, width, channel)
dim_by = torch.linspace(start=0.0, end=1.0, steps=10)

## Reshape and multiply
dim_value = rearrange(dim_by, "dim_value -> 1 dim_value 1 1 1")
images_rearr = rearrange(images, "b height width channel -> b 1 height width channel")
dimmed_images = images_rearr * dim_value

## Or in one go:
dimmed_images = einsum(
    images, dim_by,
    "batch height width channel, dim_value -> batch dim_value height width channel"
)
```

⁴值得注意的是，虽然 einops 拥有大量支持，但 einx 的实战经验较少。如果你在 einx 中发现任何限制或错误，可以自由地回退使用 einops 和一些更纯粹的 PyTorch。

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
## 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 (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

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
## 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 (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

$$y = xW^{\top}, \quad (1)$$

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

$$y = Wx, \quad (2)$$

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: 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:

$$y = Wx. \quad (3)$$

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

输入和输出的形状是什么，但这很笨拙且容易出错。使用einsum符号，文档就是实现！

Einsum符号可以处理任意的输入批处理维度，但它的主要优势在于自文档化。使用einsum符号的代码中，你的输入和输出张量的相关形状要清晰得多。对于其余的张量，你可以考虑使用Tensor类型提示，例如使用`jaxtyping`库（不特定于Jax）。

我们将在assignment2中更多地讨论使用einsum符号的性能影响，但目前要知道它们几乎总是比替代方案更好！

3.3.1 数学符号和内存排序

许多机器学习论文在其符号中使用行向量，这导致的结果与NumPy和PyTorch默认使用的行主序内存排序很好地匹配。使用行向量时，线性变换看起来像

$$y = xW^{\top}, \quad (1)$$

对于行主序 $W \in \mathbb{R}^{d_{\text{out}} \times d_{\text{in}}}$ 和行向量 $x \in \mathbb{R}^{1 \times d_{\text{in}}}$ 。

在线性代数中，通常更常用列向量，其中线性变换看起来像

$$y = Wx, \quad (2)$$

给定一个行主序 $W \in \mathbb{R}^{d_{\text{out}} \times d_{\text{in}}}$ 和列向量 $x \in \mathbb{R}^{d_{\text{in}}}$ 。在本作业中，我们将使用列向量进行数学符号表示，因为这种方式通常更容易遵循数学逻辑。您应该记住的是，如果您想使用纯矩阵乘法符号，您将不得不使用行向量约定来应用矩阵，因为PyTorch使用行主序内存排序。如果您使用`einsum`进行矩阵运算，这应该不是问题。

3.4 基本构建块：线性模块和嵌入模块

3.4.1 参数初始化

有效地训练神经网络通常需要对模型参数进行仔细初始化——不良的初始化会导致不希望的行为，例如梯度消失或梯度爆炸。预规范转换器对初始化非常鲁棒，但它们仍然会对训练速度和收敛产生显著影响。由于本作业已经很长，我们将细节留到作业3，而是给您一些适用于大多数情况的近似初始化。目前，请使用：

- 线性权重： $\mathcal{N}\left(\mu = 0, \sigma^2 = \frac{2}{d_{\text{in}} + d_{\text{out}}}\right)$ 在 $[-3\sigma, 3\sigma]$ 处截断。
- 嵌入： $\mathcal{N}(\mu = 0, \sigma^2 = 1)$ 在 $[-3, 3]$ 处截断
- RMSNorm： 1

你应该使用 `torch.nn.init.trunc_normal_` 来初始化截断的正态权重。

3.4.2 线性模块

线性层是Transformer和通用神经网络的基石。首先，你将实现你自己的`Linear`类，该类继承自`torch.nn.Module`并执行线性变换：

$$y = Wx. \quad (3)$$

请注意，我们不包含偏差项，这遵循了大多数现代LLM的做法。

Problem (linear): Implementing the linear module (1 point)

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^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`.

问题 (linear): 实现线性模块 (1 分)

交付物: 实现一个 `Linear` 类, 该类继承自 `torch.nn.Module` 并执行线性变换。您的实现应遵循 PyTorch 内置 `nn.Linear` 模块的接口, 但不应包含 `bias` 参数。我们建议以下接口:

```
def __init__(self, in_features, out_features, device=None, dtype=None) 构建线性变换模块。此函数应接受以下参数:

    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 将线性变换应用于输入。
```

确保:

- 继承 `nn.Module`
- 调用父类构造函数
- 构建并存储您的参数为 W (而不是 W^T), 出于内存排序的原因, 将其放入一个 `nn.Parameter`
- 当然, 不要使用 `nn.Linear` 或 `nn.functional.linear`

对于初始化, 使用上述设置以及 `torch.nn.init.trunc_normal_` 来初始化权重。

要测试您的 `Linear` 模块, 在 `[adapters.run_linear]` 实现测试适配器。适配器应将给定的权重加载到您的 `Linear` 模块中。您可以使用 `Module.load_state_dict` 为此目的。然后, 运行 `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 (1 point)

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
```

3.4.3 嵌入模块

如上所述, Transformer的第一层是一个嵌入层, 它将整数token ID映射到一个维度为 `d_model`的向量空间。我们将实现一个自定义的 `Embedding` 类, 该类继承自`torch.nn.Module` (因此你不应该使用 `nn.Embedding`)。 `forward` 方法应该通过索引到一个形状为 `(vocab_size, d_model)` 的嵌入矩阵来选择每个token ID的嵌入向量, 使用一个形状为`torch.LongTensor` 的token ID的 `(batch_size, sequence_length)`。

问题 (embedding): 实现嵌入模块 (1分)

交付物: 实现一个继承自 `torch.nn.Module` 的 `Embedding` 类, 并执行嵌入查找。你的实现应该遵循 PyTorch内置`nn.Embedding` 模块的接口。我们建议以下接口:

```
def __init__(self, num_embeddings, embedding_dim, device=None, dtype=None)构建一个嵌入模块。该函数应该接受以下参数:

    num_embeddings: int 词汇表的大小
```

```

embedding_dim: int Dimension of the embedding vectors, i.e.,  $d_{\text{model}}$ 
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 $a \in \mathbb{R}^{d_{\text{model}}}$ 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 $\text{RMS}(a) = \sqrt{\frac{1}{d_{\text{model}}} \sum_{i=1}^{d_{\text{model}}} a_i^2 + \epsilon}$. 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.

```

embedding_dim: int Dimension of the embedding vectors, i.e.,  $d_{\text{model}}$ 
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 查找给定tokenID的嵌入向量。

```

确保:

- 子类 `nn.Module`
- 调用超类构造函数
- 将嵌入矩阵初始化为 `nn.Parameter`
- 将嵌入矩阵存储为 `d_model` 作为最终维度
- 当然, 不要使用 `nn.Embedding` 或 `nn.functional.embedding`

再次, 使用上述设置进行初始化, 并使用 `torch.nn.init.trunc_normal_` 初始化权重。

要测试您的实现, 请在 `[adapters.run_embedding]` 实现测试适配器。然后, 运行 `uv run pytest -k test_embedding`。

3.5 预规范化 TransformerBlock

每个 Transformer 块有两个子层: 多头自注意力机制和位置感知前馈网络 (Vaswani 等人, 2017 年, 第 3.1 节)。

在原始 Transformer 论文中, 模型在每个子层周围使用残差连接, 然后进行层归一化。这种架构通常被称为“后规范”Transformer, 因为层归一化应用于子层输出。然而, 各种工作发现将层归一化从每个子层的输出移动到每个子层的输入 (在最终的 Transformer 块之后添加额外的层归一化) 可以提高 Transformer 训练的稳定性 [Nguyen 和 Salazar, 2019 年, Xiong 等人, 2020]——参见图 2 以便直观地表示这种“预规范”Transformer 块。然后, 每个 Transformer 块的子层输出通过残差连接添加到子层输入 (Vaswani 等人, 2017 年, 第 5.4 节)。预规范的一个直观理解是, 从输入嵌入到 Transformer 的最终输出之间存在一个干净的“残差流”, 没有任何归一化, 这据称可以改善梯度流。这种预规范 Transformer 现在是今天语言模型的标准 (例如, GPT-3、LLaMA、PaLM 等), 因此我们将实现这种变体。我们将逐步介绍预规范 Transformer 块的每个组件, 并按顺序实现它们。

3.5.1 均方根层归一化

Vaswani 等人的原始 Transformer 实现 [2017] 使用层归一化 [Ba 等人, 2016] 来归一化激活值。遵循 Touvron 等人 [2023], 我们将使用均方根层归一化 (RMSNorm; Zhang 和 Sennrich, 2019, 公式 4) 进行层归一化。给定一个激活值向量 $a \in \mathbb{R}^{d_{\text{model}}}$, RMSNorm 将按以下方式重新缩放每个激活值 a_i :

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

在 $\text{RMS}(a) = \sqrt{\frac{1}{d_{\text{model}}} \sum_{i=1}^{d_{\text{model}}} a_i^2 + \epsilon}$ 。这里, g_i 是一个可学习的“增益”参数 (总共有 `d_model` 个这样的参数), 而 ϵ 是一个通常固定为 1e-5 的超参数。

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 (1 point)

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 = 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 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`.

你应该将输入提升为 `torch.float32` 以防止在平方输入时溢出。总的来说，你的 `forward` 方法应该看起来像：

```
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 (1 point)

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 = 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 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

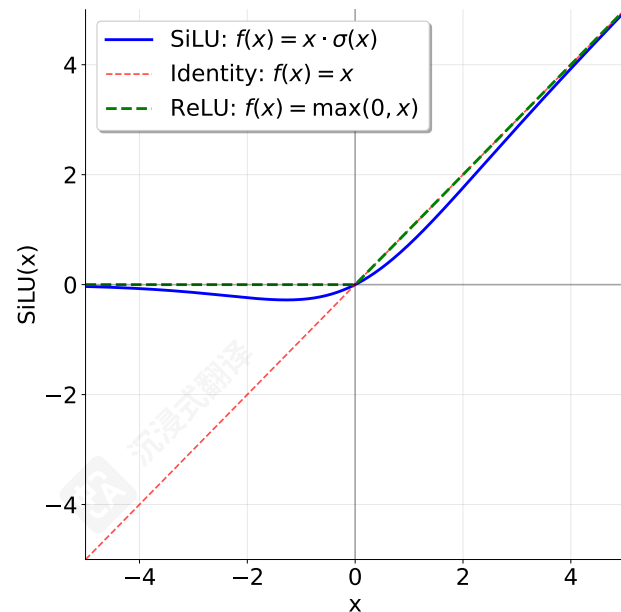


Figure 3: Comparing the SiLU (aka Swish) and ReLU activation functions.

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 4x 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 [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, Elfwing 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, $d_{\text{ff}} = \frac{8}{3}d_{\text{model}}$.

3.5.2 位置前馈网络

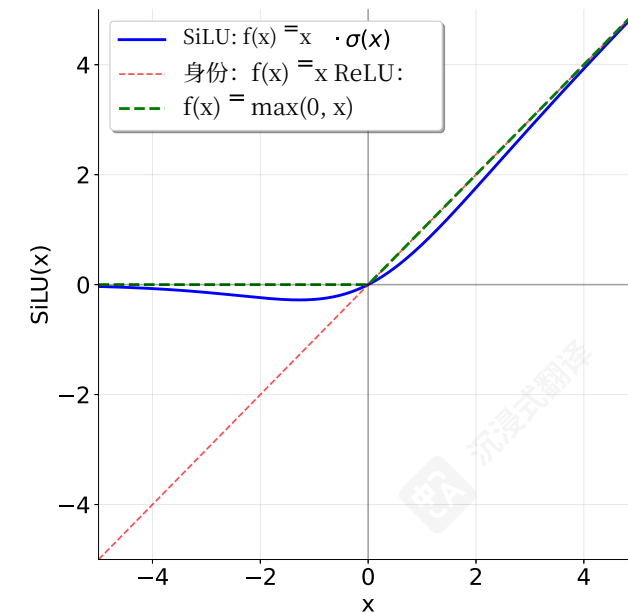


图 3：比较 SiLU（又名 Swish）和 ReLU 激活函数。

在原始 Transformer 论文 (Vaswani 等人 [2017] 的第 3.3 节) 中，Transformer 前馈网络由两个线性变换组成，它们之间有一个 ReLU 激活 ($\text{ReLU}(x) = \max(0, x)$)。内部前馈层的维度通常是输入维度大小的 4 倍。

然而，现代语言模型与原始设计相比，倾向于进行两个主要改变：它们使用另一种激活函数，并采用门控机制。具体来说，我们将实现 LLMs (如 Llama 3 [Grattafiori 等人, 2024] 和 Qwen 2.5 [Yang 等人, 2024]) 中采用的 “SwiGLU” 激活函数，该函数结合了 SiLU (通常称为 Swish) 激活和一种称为门控线性单元 (GLU) 的门控机制。我们还将省略线性层中有时使用的偏差项，遵循 PaLM [Chowdhery 等人, 2022] 和 LLaMA [Touvron 等人, 2023] 以来大多数现代 LLMs 的做法。

SiLU 或 Swish 激活函数 [Hendrycks and Gimpel, 2016, Elfwing et al., 2017] 定义如下：

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

如图 3 所示，SiLU 激活函数与 ReLU 激活函数相似，但在零点处平滑。

门控线性单元 (GLUs) 最初由 Dauphin 等人 [2017] 定义为通过一个 sigmoid 函数传递的线性变换与另一个线性变换的逐元素乘积：

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

其中 \odot 表示逐元素乘法。门控线性单元被建议 “通过为梯度提供线性路径来减少深度架构的梯度消失问题，同时保留非线性能力。”

将 SiLU/Swish 和 GLU 结合起来，我们得到 SwiGLU，我们将用于我们的前馈网络：

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

在 $x \in \mathbb{R}^{d_{\text{model}}}$ 、 $W_1, W_3 \in \mathbb{R}^{d_{\text{ff}} \times d_{\text{model}}}$ 、 $W_2 \in \mathbb{R}^{d_{\text{model}} \times d_{\text{ff}}}$ 以及规范上， $d_{\text{ff}} = \frac{8}{3}d_{\text{model}}$

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 (2 points)

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 2d vectors by the angle $\theta_{i,k} = \frac{i}{\Theta(2k-2)/d}$ 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 0s 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 2d 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.

Shazeer [2020] 首次提出了将SiLU/Swish激活函数与GLU结合，并进行了实验表明SwiGLU在语言建模任务上优于ReLU和SiLU（无门控）等基线。在作业的后面，你将比较SwiGLU和SiLU。尽管我们已经为这些组件提供了一些启发式论据（并且论文提供了更多支持证据），但保持经验性视角是好的：

Shazeer论文中有一句著名的话是

我们对这些架构为何似乎有效并不做解释；我们将它们的成功，如同其他一切，归因于神明的恩惠。

问题 (positionwise_feedforward): 实现位置感知的前馈网络 (2 分)

交付物：实现 SwiGLU 前馈网络，该网络由一个 SiLU 激活函数和一个 GLU 组成。

注意：在这种情况下，为了数值稳定性，你应该在实现中自由使用 `torch.sigmoid`。

你应该在实现中将 d_{ff} 设置为大约 $\frac{8}{3} \times d_{\text{model}}$ ，同时确保内部前馈层的维度是 64 的倍数，以充分利用你的硬件。为了使用我们提供的测试来测试你的实现，你需要实现测试适配器 `[adapters.run_swiglu]`。然后，运行 `uv run pytest -k test_swiglu` 来测试你的实现。

3.5.3 相对位置嵌入

为了将位置信息注入模型，我们将实现旋转位置嵌入 [Su et al.2021]，通常称为 RoPE。对于给定的查询标记 $q^{(i)} = W_q x^{(i)} \in \mathbb{R}^d$ 在标记位置 i ，我们将应用一个成对旋转矩阵 R^i ，得到 $q'^{(i)} = R^i q^{(i)} = R^i W_q x^{(i)}$ 。在这里， R^i 将作为 2d 向量通过角度 $\theta_{i,k} = \frac{i}{\Theta(2k-2)/d}$ 旋转嵌入元素对 $q_{2k-1:2k}^{(i)}$ ，用于 $k \in \{1, \dots, d/2\}$ 和一些常数 Θ 。因此，我们可以认为 R^i 是一个大小为 $d \times d$ 的块对角矩阵，其中块 R_k^i 用于 $k \in \{1, \dots, d/2\}$ ，

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

因此我们得到完整的旋转矩阵

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

其中 0 表示 2×2 零矩阵。虽然可以构造完整的 $d \times d$ 矩阵，但一个好的解决方案应该利用该矩阵的性质来更高效地实现转换。由于我们只关心给定序列中标记的相对旋转，我们可以重用我们为 $\cos(\theta_{i,k})$ 和 $\sin(\theta_{i,k})$ 计算的跨层和不同的批次。如果您想优化它，您可以使用一个被所有层引用的单一 RoPE 模块，它可以在 init 过程中用 `self.register_buffer(persistent=False)` 创建一个 2d 预计算的正弦和余弦值缓冲区，而不是一个 `nn.Parameter`（因为我们不想学习这些固定的余弦和正弦值）。然后，我们对 $k^{(j)}$ 执行与我们为 $q^{(i)}$ 所做的完全相同的旋转过程，旋转对应于 R^j 。请注意，该层没有可学习的参数。

Problem (rope): Implement RoPE (2 points)

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.

    theta: float Θ value for the RoPE
    d_k: int dimension of query and key vectors
    max_seq_len: int Maximum sequence length that will be inputted
    device: torch.device | None = None Device to store the buffer on

def forward(self, x: torch.Tensor, token_positions: torch.Tensor) -> torch.Tensor
Process an input tensor of shape (... , seq_len, d_k) and return a tensor of the same shape.
Note that you should tolerate x with an arbitrary number of batch dimensions. You should
assume that the token positions are a tensor of shape (... , seq_len) specifying the token
positions of x along the sequence dimension.

You should use the token positions to slice your (possibly precomputed) cos and 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:

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

(10)

Note that $\exp(v_i)$ can become `inf` for large values (then, `inf/inf = 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 (1 point)

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:

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

(11)

问题 (rope): 实现 RoPE (2 分)

交付物: 实现一个类 RotaryPositionalEmbedding, 它将 RoPE 应用于输入张量。建议的接口如下:

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

    theta: float Θ value for the RoPE
    d_k: int dimension of query and key vectors
    max_seq_len: int Maximum sequence length that will be inputted
    device: torch.device | None = None Device to store the buffer on

def forward(self, x: torch.Tensor, token_positions: torch.Tensor) -> torch.Tensor
处理一个形状为 (... , seq_len, d_k) 的输入张量并返回一个相同形状的张量。注意, 你应
该容忍 x, 它具有任意数量的批处理维度。你应该假设 token 位置是一个形状为
(... , seq_len) 的张量, 指定 x 沿序列维度的 token 位置。你应该使用 token 位置来沿序
列维度切片你的 (可能预先计算的) cos 和 sin 张量。

为了测试你的实现, 完成 [adapters.run_rope] 并确保它通过 uv runpytest -k test_rope。
```

3.5.4 缩放点积注意力

我们现在将实现 Vaswani 等人 [2017] (第 3.2.1 节) 中描述的缩放点积注意力。作为初步步骤, 注意力的定义将使用 softmax, 这是一个将未归一化的分数向量转换为归一化分布的操作:

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

(10)

请注意, $\exp(v_i)$ 对于大值可能会变成 `inf` (然后, `inf/inf = NaN`)。我们可以通过注意到 softmax 操作对于所有输入加上任何常数 c 都是不变的来避免这种情况。我们可以利用这个属性来提高数值稳定性——通常, 我们会从 o_i 的所有元素中减去最大值, 使新的最大值为 0。现在, 您将使用这个技巧来实现 softmax。

问题 (softmax): 实现 softmax (1 分)

交付物: 编写一个函数来应用 softmax 操作到一个张量上。您的函数应该接受两个参数: 一个张量和维度 i , 并将 softmax 应用到输入张量的第 i 维。输出张量应该与输入张量具有相同的形状, 但其第 i 维现在将是一个归一化的概率分布。使用从第 i 维的最大值中减去第 i 维的所有元素来避免数值稳定性问题的技巧。

要测试您的实现, 完成 [adapters.run_softmax] 并确保它通过 `uv run pytest -k test_softmax_matches_pytorch`。

现在我们可以将注意力操作数学定义为如下:

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

(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 $\left(\frac{Q^\top K}{\sqrt{d_k}}\right)$ and adding a $-\infty$ in any entry of the mask matrix that is False.

Problem (scaled_dot_product_attention): Implement scaled dot-product attention (5 points)

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:

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

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

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 multiplies.⁵

⁵As a stretch goal, try combining the key, query, and value projections into a single weight matrix so you only need a single matrix multiply.

其中 $Q \in \mathbb{R}^{n \times d_k}$, $K \in \mathbb{R}^{m \times d_k}$ 和 $V \in \mathbb{R}^{m \times d_v}$ 。在这里, Q , K 和 V 都是此操作的输入——请注意, 这些不是可学习参数。如果你想知道为什么这不是

QK^\top , 参见 3.3.1。
掩码: 有时为了方便, 需要掩码注意力操作的结果。掩码应该具有形状 $M \in \{\text{True}, \text{False}\}^{n \times m}$, 这个布尔矩阵的每一行 i 表示查询 i 应该关注哪些键。规范上 (并且有点令人困惑), 位置 (i, j) 处的值 **True** 表示查询 i 关注键 j , 而值 **False** 表示查询不关注键。换句话说, “信息流” 在 (i, j) 对中与值 **True** 匹配。例如, 考虑一个 1×3 掩码矩阵, 其条目为 `[[True, True, False]]`。单个查询向量只关注前两个键。

从计算角度看, 使用掩码比计算子序列的注意力效率高得多, 我们可以通过取预softmax值 $\left(\frac{Q^\top K}{\sqrt{d_k}}\right)$ 并在掩码矩阵中任何为 False 的条目中添加一个 $-\infty$ 来实现这一点。

Problem (scaled_dot_product_attention): Implement scaled dot-product attention (5 points)

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 因果多头自注意力

我们将实现 Vaswani 等人在 3.2.2 节中描述的多头自注意力 [2017]。回想一下, 从数学上讲, 应用多头注意力的操作定义如下:

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

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

其中 Q_i , K_i , V_i 分别是 Q , K 和 V 的嵌入维度大小为 d_k 或 d_v 的第 $i \in \{1, \dots, h\}$ 个切片。Attention 是 § 3.5.4 中定义的缩放点积注意力操作。由此我们可以形成多头自注意力操作:

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

这里, 可学习参数是 $W_Q \in \mathbb{R}^{hd_k \times d_{\text{model}}}$, $W_K \in \mathbb{R}^{hd_k \times d_{\text{model}}}$, $W_V \in \mathbb{R}^{hd_v \times d_{\text{model}}}$, 和 $W_O \in \mathbb{R}^{d_{\text{model}} \times hd_v}$ 。由于在多头注意力操作中, Q 、 K 和 V 被切片, 我们可以将 W_Q , W_K 和 W_V 沿输出维度视为每个头独立分离。当你实现这个功能时, 你应该在总共三个矩阵乘法中计算键、值和查询投影。⁵

⁵作为一个挑战目标, 尝试将键、查询和值投影合并到一个权重矩阵中, 这样你只需要一个矩阵乘法。

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 (5 points)

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.

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` to test your implementation.

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 ,

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

(15)

Problem (transformer_block): Implement the Transformer block (3 points)

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.

因果掩码。你的实现应该防止模型关注序列中的未来标记。换句话说，如果模型被给定一个标记序列 t_1, \dots, t_n , 并且我们想要计算前缀 t_1, \dots, t_i (其中 $i < n$) 的下一个词预测，模型不应该能够访问（关注）位置 $t_{i+1} \dots t_n$, 因为它在推理（生成文本）时无法访问这些标记（并且这些未来标记会泄露关于真实下一个词身份的信息，使语言建模预训练目标变得简单）。对于输入标记序列 $t_1 \dots t_n$, 我们可以通过运行 n 次多头自注意力（对于序列中的 n 个唯一前缀）来天真地防止访问未来标记。相反，我们将使用因果注意力掩码，它允许标记 i 关注序列中的所有位置 $j \leq i$ 。你可以使用 `torch.triu` 或广播索引比较来构建这个掩码，并且你应该利用你从 §3.5.4 中实现的缩放点积注意力已经支持注意力掩码这一事实。

应用 RoPE。RoPE 应该应用于查询和键向量，但不应用于值向量。此外，头维度应该作为批处理维度处理，因为在多头注意力中，注意力是独立地应用于每个头的。这意味着对于每个头，都应该对查询和键向量应用完全相同的 RoPE 旋转。

问题 (multihead_self_attention): 实现因果多头自注意力 (5 点)

交付物：将因果多头自注意力实现为 `torch.nn.Module`。您的实现应该接受（至少）以下参数：

`d_model`: `int` Transformer 块输入的维度。

`num_heads`: `int` 在多头自注意力中使用的头数。

遵循 Vaswani 等人 [2017], 设置 $d_k = d_v = d_{\text{model}}/h$ 。要测试您的实现与我们的提供的测试，实现测试适配器在 `[adapters.run_multihead_self_attention]`。然后，运行 `uv run pytest -k test_multihead_self_attention` 来测试您的实现。

3.6 The FullTransformer LM

让我们从组装 Transformer 模块开始（参考图 2 会很有帮助）。一个 Transformer 模块包含两个“子层”，一个用于多头自注意力，另一个用于前馈网络。在每个子层中，我们首先执行 RMSNorm，然后执行主要操作（MHA/FF），最后添加残差连接。

具体来说，Transformer 模块的前半部分（第一个“子层”）应该实现以下更新集，以从输入 x 产生输出 y ,

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

(15)

问题 (transformer_block): 实现 Transformer 模块 (3 分)

按照 §3.5 的描述和图 2 所示，实现预归一化 Transformer 模块。您的 Transformer 模块应接受（至少）以下参数。

`d_model`: `int` Transformer 模块输入的维度。

`num_heads`: `int` 用于多头自注意力的头数。

`d_ff`: `int` 位置前馈内层的维度。

To test your implementation, implement the adapter `[adapters.run_transformer_block]`. Then run `uv run pytest -k test_transformer_block` to test your implementation.
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 (3 points)

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` to test your implementation.

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

Resource accounting. It is useful to be able to understand how the various parts of the Transformer consume compute and memory. We will go through the steps to do some basic “FLOPs accounting.” The vast majority of FLOPS in a Transformer are matrix multiplies, so our core approach is simple:

- Write down all the matrix multiplies in a Transformer forward pass.
- Convert each matrix multiply into FLOPs required.

For this second step, the following fact will be useful:

Rule: Given $A \in \mathbb{R}^{m \times n}$ and $B \in \mathbb{R}^{n \times p}$, the matrix-matrix product AB requires $2mnp$ FLOPs.

To see this, note that $(AB)[i, j] = A[i, :] \cdot B[:, j]$, and that this dot product requires n additions and n multiplications ($2n$ FLOPs). Then, since the matrix-matrix product AB has $m \times p$ entries, the total number of FLOPS is $(2n)(mp) = 2mnp$.

Now, before you do the next problem, it can be helpful to go through each component of your Transformer block and Transformer LM, and list out all the matrix multiplies and their associated FLOPs costs.

Problem (transformer_accounting): Transformer LM resource accounting (5 points)

(a) Consider GPT-2 XL, which has the following configuration:

```
vocab_size : 50,257
context_length : 1,024
num_layers : 48
d_model : 1,600
```

要测试您的实现，请实现适配器 `[adapters.run_transformer_block]`。然后运行 `uv run pytest -k test_transformer_block` 来测试您的实现。交付物：通过提供的测试的 Transformer 模块代码。

现在我们将模块组合在一起，遵循图1中的高级设计图。按照我们在3.1.1节中对嵌入的描述，将其输入到 `num_layers` Transformer模块中，然后将其传递到三个输出层以获得词汇表上的分布。

问题 (transformer_lm): 实现Transformer LM (3分)

是时候将所有内容整合在一起了！根据 § 3.1中的描述和图1中的说明实现Transformer语言模型。至少，您的实现应接受上述Transformer模块的所有构造参数，以及这些附加参数：

- vocab_size: int** 词汇表的大小，用于确定标记的维度嵌入矩阵。
- context_length: int** 最大上下文长度，用于确定位置嵌入矩阵。
- num_layers: int** 要使用的 Transformer 模块数量。

要测试您的实现与我们的提供的测试，您首先需要在 `[adapters.run_transformer_lm]` 实现测试适配器。然后，运行 `uv run pytest -k test_transformer_lm` 来测试您的实现。

交付物：一个通过上述测试的 Transformer LM 模块。

资源核算。能够理解 Transformer 的各个部分如何消耗计算和内存是有用的。我们将通过一些基本的 “FLOPs 核算” 步骤来讲解。Transformer 中的绝大多数 FLOPS 都是矩阵乘法，因此我们的核心方法是简单的：

- 写下 Transformer 前向传递中的所有矩阵乘法。
- 将每个矩阵乘法转换为所需的 FLOPs。

对于这个第二步，以下事实将很有用：规则：给定 A 和 B ，矩阵

$$A \in \mathbb{R}^{m \times n} \quad B \in \mathbb{R}^{n \times p} \quad \text{-matrix product } AB \text{ requires } 2mnp \text{ FLOPs.}$$

要看到这一点，请注意 $(AB)[i, j] = A[i, :] \cdot B[:, j]$ ，并且这个点积需要 n 次加法和 n 次乘法 ($2n$ FLOPs)。然后，由于矩阵-矩阵乘法 AB 有 $m \times p$ 个元素，FLOPs 的总数是 $(2n)(mp) = 2mnp$ 。

现在，在你做下一个问题之前，回顾你的 Transformer 模块和 Transformer LM 的每个组件，并列出所有矩阵乘法及其相关的 FLOPs 成本可能会有所帮助。

问题 (transformer_accounting): Transformer LM 资源核算 (5 分)

(a) 考虑 GPT-2 XL，它具有以下配置：

```
vocab_size : 50,257
context_length : 1,024
num_layers : 48
d_model : 1,600
```


num_heads : 25
d_ff : 6,400

Suppose we constructed our model using this configuration. How many trainable parameters would our model have? Assuming each parameter is represented using single-precision floating point, how much memory is required to just load this model?

Deliverable: A one-to-two sentence response.

- (b) Identify the matrix multiplies required to complete a forward pass of our GPT-2 XL-shaped model. How many FLOPs do these matrix multiplies require in total? Assume that our input sequence has `context_length` tokens.

Deliverable: A list of matrix multiplies (with descriptions), and the total number of FLOPs required.

- (c) Based on your analysis above, which parts of the model require the most FLOPs?

Deliverable: A one-to-two sentence response.

- (d) Repeat your analysis with GPT-2 small (12 layers, 768 `d_model`, 12 heads), GPT-2 medium (24 layers, 1024 `d_model`, 16 heads), and GPT-2 large (36 layers, 1280 `d_model`, 20 heads). As the model size increases, which parts of the Transformer LM take up proportionally more or less of the total FLOPs?

Deliverable: For each model, provide a breakdown of model components and its associated FLOPs (as a proportion of the total FLOPs required for a forward pass). In addition, provide a one-to-two sentence description of how varying the model size changes the proportional FLOPs of each component.

- (e) Take GPT-2 XL and increase the context length to 16,384. How does the total FLOPs for one forward pass change? How do the relative contribution of FLOPs of the model components change?

Deliverable: A one-to-two sentence response.

num_heads : 25
d_ff : 6,400

假设我们使用此配置构建模型。我们的模型有多少可训练参数？假设每个参数使用单精度浮点数表示，仅加载此模型需要多少内存？

交付物：一段一到两句话的回应。

- (b) 确定完成我们 GPT-2 XL 形状模型正向传递所需的矩阵乘法。这些矩阵乘法总共需要多少 FLOPs？假设我们的输入序列有 `context_length` 个 token。

交付物：一个矩阵乘法列表（带描述），以及所需的 FLOPs 总数。

- (c) 根据你的分析，模型的哪些部分需要最多的 FLOPs？

可交付成果：一段至两段回应。

- (d) 使用 GPT-2 小型（12 层， 768 `d_model`， 12 个头）、GPT-2 中型（24 层， 1024 `d_model`， 16 个头）和 GPT-2 大型（36 层， 1280 `d_model`， 20 个头）重复你的分析。随着模型尺寸的增加，Transformer LM 的哪些部分占用了更多或更少的总 FLOPs？

可交付成果：对于每个模型，提供模型组件及其关联 FLOPs（作为正向传递所需总 FLOPs 比例）的分解。此外，提供一段至两段描述，说明如何通过改变模型尺寸来改变各组件的 FLOPs 比例的描述。

- (e) 使用 GPT-2 XL 并将上下文长度增加到 16,384。一个正向传递的总 FLOPs 如何变化？模型组件的 FLOPs 相对贡献如何变化？

可交付成果：一段一到两句话的回应。

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 $o_i \in \mathbb{R}^{\text{vocab_size}}$ for each position i , which results in:⁶

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

The cross entropy loss is generally defined with respect to the vector of logits $o_i \in \mathbb{R}^{\text{vocab_size}}$ and target x_{i+1} .⁷

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.

Implement `[adapters.run_cross_entropy]`, then run `uv run pytest -k test_cross_entropy` to test your implementation.

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 :

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

⁶Note that $o_i[k]$ refers to value at index k of the vector o_i .

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

4 训练 TransformerLM

我们现在有了预处理数据的步骤（通过 tokenizer）和模型的步骤（Transformer）。剩下的就是构建所有支持训练的代码。这包括以下内容：

- 损失：我们需要定义损失函数（交叉熵）。
- 优化器：我们需要定义优化器来最小化这个损失（AdamW）。
- 训练循环：我们需要所有支持的基础设施来加载数据、保存检查点，并管理训练。

4.1 交叉熵损失

回想一下，Transformer语言模型为每个长度为 $m+1$ 和 $i = 1 \dots m$ 的序列 x 定义了一个分布 $p_\theta(x_{i+1} \mid x_{1:i})$ ，，。给定一个由长度为 m 的序列组成的训练集 D ，我们定义了标准的交叉熵（负对数似然）损失函数：

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

(注意，Transformer的单次前向传递为所有 $i = 1, \dots, m$ 生成 $p_\theta(x_{i+1} \mid x_{1:i})$ 。)

特别是，Transformer为每个位置 i 计算logits $o_i \in \mathbb{R}^{\text{vocab_size}}$ ，这导致：⁶

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

交叉熵损失通常相对于logits向量 $o_i \in \mathbb{R}^{\text{vocab_size}}$ 和目标 x_{i+1} .⁷定义

实现交叉熵损失需要像softmax的情况一样，在数值问题上小心处理。

问题(cross_entropy): 实现交叉熵

可交付成果：编写一个函数来计算交叉熵损失，该函数接收预测logits (o_i) 和目标 (x_{i+1}) 并计算交叉熵 $\ell_i = -\log \text{softmax}(o_i)[x_{i+1}]$ 。您的函数应处理以下内容：

- 为数值稳定性减去最大元素。
- 尽可能取消log和exp。
- 处理任何额外的批次维度，并在批次上返回平均值。与第3.3节一样，我们假设批次类维度始终排在最前面，在词汇总量维度之前。

实现 `[adapters.run_cross_entropy]`，然后运行 `uv run pytest -k test_cross_entropy` 以测试您的实现。

困惑度交叉熵足以用于训练，但在我们评估模型时，我们也想报告困惑度。对于一个长度为 m 的序列，其中我们遭受交叉熵损失 $\ell_1 \dots \ell_m$ ：，

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

⁶请注意， $o_i[k]$ 指的是向量 o_i 的索引 k 处的值。⁷这对应于 x_{i+1} 上的狄拉克δ分布与预测的 $\text{softmax}(o_i)$ 分布之间的交叉熵。

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, ...)` should initialize your optimizer. Here, `params` will be a collection of parameters to be optimized (or parameter groups, in case the user wants to use different hyperparameters, such as learning rates, for different parts of the model). Make sure to pass `params` to the `__init__` method of the base class, which will store these parameters for use in `step`. You can take additional arguments depending on the optimizer (e.g., the learning rate is a common one), and pass them to the base class constructor as a dictionary, where keys are the names (strings) you choose for these parameters.

`def step(self)` should make one update of the parameters. During the training loop, this will be called after the backward pass, so you have access to the gradients on the last batch. This method should iterate through each parameter tensor `p` and modify them *in place*, i.e. setting `p.data`, which holds the tensor associated with that parameter based on the gradient `p.grad` (if it exists), the tensor representing the gradient of the loss with respect to that parameter.

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.
```

4.2 SGD 优化器

现在我们已经有了损失函数，我们将开始探索优化器。最简单的基于梯度的优化器是随机梯度下降（SGD）。我们从随机初始化的参数 θ_0 开始。然后对于每一步 $t = 0 \dots T - 1$ ，我们执行以下更新：

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

其中 B_t 是从数据集 D 中采样的随机批次数据，而学习率 α_t 和批次大小 $|B_t|$ 是超参数。

4.2.1 在 PyTorch 中实现 SGD

为了实现我们的优化器，我们将继承 PyTorch `torch.optim.Optimizer` 类。一个 `Optimizer` 子类必须实现两个方法：

`def __init__(self, params, ...)` 应初始化优化器。在这里，`params` 将是要优化的参数集合（或者参数组，如果用户希望对不同模型部分使用不同的超参数，例如学习率）。确保将 `params` 传递给基类的 `__init__` 方法，该方法将存储这些参数以供 `step` 使用。你可以根据优化器传入额外的参数（例如，学习率是一个常见的参数），并将它们作为字典传递给基类构造函数，其中键是你为这些参数选择的名称（字符串）。

`def step(self)` 应该进行一次参数更新。在训练循环中，这将在反向传播之后被调用，因此你可以获取最后一个批次的梯度。这个方法应该遍历每个参数张量 `p` 并就地修改它们，即设置 `p.data`，它根据梯度 `p.grad`（如果存在）持有与该参数相关的张量，即损失对该参数的梯度表示的张量。

PyTorch 优化器 API 有一些微妙之处，因此用示例来解释它更容易。为了使我们的示例更丰富，我们将实现 SGD 的一种轻微变体，其中学习率在训练过程中衰减，从初始学习率 α 开始，随着时间的推移逐步减小步长：

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

让我们看看这个版本的 SGD 如何作为 PyTorch `Optimizer` 实现：

```
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.

Problem (learning_rate_tuning): Tuning the learning rate (1 point)

As we will see, one of the hyperparameters that affects training the most is the learning rate. Let's see that in practice in our toy example. Run the SGD example above with three other values for the learning rate: `1e1`, `1e2`, and `1e3`, for just 10 training iterations. What happens with the loss for each of these learning rates? Does it decay faster, slower, or does it diverge (i.e., increase over the course of training)?

Deliverable: A one-two sentence response with the behaviors you observed.

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),

```

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

```

在 `__init__` 中，我们将参数传递给优化器，以及默认超参数传递给基类构造函数（参数可能会分组，每组具有不同的超参数）。如果参数只是一个 `torch.nn.Parameter` 对象的单一集合，基构造函数将创建一个单一组，并将其分配默认超参数。然后，在 `step` 中，我们遍历每个参数组，然后遍历该组中的每个参数，并应用公式 20。在这里，我们将迭代次数作为与每个参数关联的状态保留：我们首先读取此值，在梯度更新中使用它，然后更新它。API 指定用户可能会传递一个可调用的 `closure` 在优化器步骤之前重新计算损失。对于我们将要使用的优化器，我们不需要这个，但我们添加它以符合 API。

要看到这个效果，我们可以使用以下训练循环的最小示例：

```

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.

```

这是训练循环的典型结构：在每个迭代中，我们将计算损失并运行优化器的一步。在训练语言模型时，我们的可学习参数将来自模型（在 PyTorch 中，`m.parameters()` 给我们这个集合）。损失将在采样的数据批次上计算，但训练循环的基本结构将是相同的。

问题 (learning_rate_tuning): 调整学习率 (1 分)

正如我们将看到的，影响训练最多的超参数之一是学习率。让我们在我们的玩具示例中实际看看这一点。运行上面 SGD 示例的其他三个学习率值：`1e1`、`1e2` 和 `1e3`，仅进行 10 次训练迭代。对于这些学习率中的每一个，损失会发生什么？它是更快地衰减、更慢地衰减还是发散（即，在训练过程中增加）？

可交付成果：用一到两句话描述你观察到的行为。

4.3 AdamW

现代语言模型通常使用更复杂的优化器进行训练，而不是 SGD。最近使用的优化器大多是 Adam 优化器的衍生品 [Kingma 和 Ba, 2015]。我们将使用 AdamW [Loshchilov 和 Hutter, 2019]，它在最近的工作中得到了广泛应用。AdamW 对 Adam 提出了一种修改，通过添加权重衰减（在每次迭代中，我们将参数拉向 0）来改进正则化，

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

```
init( $\theta$ ) (Initialize learnable parameters)
 $m \leftarrow 0$  (Initial value of the first moment vector; same shape as  $\theta$ )
 $v \leftarrow 0$  (Initial value of the second moment vector; same shape as  $\theta$ )
for  $t = 1, \dots, T$  do
    Sample batch of data  $B_t$ 
     $g \leftarrow \nabla_{\theta} \ell(\theta; B_t)$  (Compute the gradient of the loss at the current time step)
     $m \leftarrow \beta_1 m + (1 - \beta_1)g$  (Update the first moment estimate)
     $v \leftarrow \beta_2 v + (1 - \beta_2)g^2$  (Update the second moment estimate)
     $\alpha_t \leftarrow \alpha \frac{\sqrt{1 - (\beta_2)^t}}{1 - (\beta_1)^t}$  (Compute adjusted  $\alpha$  for iteration  $t$ )
     $\theta \leftarrow \theta - \alpha_t \frac{m}{\sqrt{v + \epsilon}}$  (Update the parameters)
     $\theta \leftarrow \theta - \alpha \lambda \theta$  (Apply weight decay)
end for
```

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

Problem (adamw): Implement AdamW (2 points)

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 gives you 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). Implement `[adapters.get_adamw_cls]` and make sure it passes `uv run pytest -k test_adamw`.

Problem (adamwAccounting): Resource accounting for training with AdamW (2 points)

Let us compute how much memory and compute running AdamW requires. Assume we are using float32 for every tensor.

(a) How much peak memory does running AdamW require? Decompose your answer based on the memory usage of the parameters, activations, gradients, and optimizer state. Express your answer in terms of the `batch_size` and the model hyperparameters (`vocab_size`, `context_length`, `num_layers`, `d_model`, `num_heads`). Assume `d_ff = 4 × d_model`.

For simplicity, when calculating memory usage of activations, consider only the following components:

- Transformer block
 - RMSNorm(s)

这种方式与梯度更新是解耦的。我们将按照Loshchilov和Hutter在算法2中描述的方式实现AdamW [2019]。

AdamW 是有状态的：对于每个参数，它跟踪其第一和第二矩的运行估计。因此，AdamW 以额外的内存换取改进的稳定性和收敛性。除了学习率 α ，AdamW 还有一对控制矩估计更新的超参数 (β_1, β_2) 以及一个权重衰减率 λ 。典型应用将 (β_1, β_2) 设置为 $(0.9, 0.999)$ ，但像 LLaMA [Touvron 等人, 2023] 和 GPT-3 [Brown 等人, 2020] 通常使用 $(0.9, 0.95)$ 进行训练。算法可以写成如下形式，其中 ϵ 是一个小的值（例如， 10^{-8} ），用于在 v 中出现极小值时提高数值稳定性：

Algorithm 1 AdamW Optimizer

```
init( $\theta$ ) (Initialize learnable parameters)
 $m \leftarrow 0$  (Initial value of the first moment vector; same shape as  $\theta$ )
 $v \leftarrow 0$  (Initial value of the second moment vector; same shape as  $\theta$ )
for  $t = 1, \dots, T$  do
    Sample batch of data  $B_t$ 
     $g \leftarrow \nabla_{\theta} \ell(\theta; B_t)$  (Compute the gradient of the loss at the current time step)
     $m \leftarrow \beta_1 m + (1 - \beta_1)g$  (Update the first moment estimate)
     $v \leftarrow \beta_2 v + (1 - \beta_2)g^2$  (Update the second moment estimate)
     $\alpha_t \leftarrow \alpha \frac{\sqrt{1 - (\beta_2)^t}}{1 - (\beta_1)^t}$  (Compute adjusted  $\alpha$  for iteration  $t$ )
     $\theta \leftarrow \theta - \alpha_t \frac{m}{\sqrt{v + \epsilon}}$  (Update the parameters)
     $\theta \leftarrow \theta - \alpha \lambda \theta$  (Apply weight decay)
end for
```

请注意 t 从 1 开始。现在您将实现这个优化器。

问题 (adamw): 实现 AdamW (2 分)

可交付成果：将 AdamW 优化器实现为 `torch.optim.Optimizer` 的子类。您的类应接受学习率 α 在 `__init__` 中，以及 β 、 ϵ 和 λ 超参数。为了帮助您保持状态，基础 `Optimizer` 类为您提供了一个字典 `self.state`，它将 `nn.Parameter` 对象映射到一个字典，该字典存储您需要为该参数的任何信息（对于 AdamW，这将是对时刻的估计）。实现 `[adapters.get_adamw_cls]` 并确保它通过 `uv run pytest -k test_adamw`。

问题 (adamwAccounting): 使用 AdamW 进行训练的资源核算 (2 分)

让我们计算运行 AdamW 需要多少内存和计算资源。假设我们为每个张量使用 float32。

(a) 运行 AdamW 需要多少峰值内存？根据参数、激活值、梯度和优化器状态的内存使用情况分解您的答案。用 `batch_size` 和模型超参数 (`vocab_size`、`context_length`、`num_layers`、`d_model`、`num_heads`) 表示您的答案。假设 `d_ff = 4 × d_model`。

为简化起见，在计算激活内存使用时，仅考虑以下组件：

- Transformer block
 - RMSNorm(s)

- Multi-head self-attention sublayer: QKV projections, Q^TK matrix multiply, softmax, weighted sum of values, output projection.
- Position-wise feed-forward: W_1 matrix multiply, SiLU, W_2 matrix multiply
- final RMSNorm
- output embedding
- cross-entropy on logits

Deliverable: An algebraic expression for each of parameters, activations, gradients, and optimizer state, as well as the total.

- (b) Instantiate your answer for a GPT-2 XL-shaped model to get an expression that only depends on the `batch_size`. What is the maximum batch size you can use and still fit within 80GB memory?

Deliverable: An expression that looks like $a \cdot \text{batch_size} + b$ for numerical values a, b , and a number representing the maximum batch size.

- (c) How many FLOPs does running one step of AdamW take?

Deliverable: An algebraic expression, with a brief justification.

- (d) Model FLOPs utilization (MFU) is defined as the ratio of observed throughput (tokens per second) relative to the hardware's theoretical peak FLOP throughput [Chowdhery et al., 2022]. An NVIDIA A100 GPU has a theoretical peak of 19.5 teraFLOP/s for float32 operations. Assuming you are able to get 50% MFU, how long would it take to train a GPT-2 XL for 400K steps and a batch size of 1024 on a single A100? Following Kaplan et al. [2020] and Hoffmann et al. [2022], assume that the backward pass has twice the FLOPs of the forward pass.

Deliverable: The number of days training would take, with a brief justification.

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 trains⁸. 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 .

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 $\alpha_t = \alpha_{\min}$.

- Multi-head self-attention sublayer: QKV 投影, Q^TK 矩阵乘法, softmax, 值的加权和, 输出投影。
- Position-wise feed-forward: W_1 矩阵乘法, SiLU, W_2 矩阵乘法

- 最终 RMSNorm
- 输出嵌入
- logits 上的交叉熵

交付物: 每个参数、激活、梯度和优化器状态的代数表达式, 以及总计。

- (b) 为 GPT-2 XL 形状模型实例化你的答案以得到一个仅依赖于 `batch_size` 的表达式。你能使用的最大批处理大小是多少, 并且仍然适合在 80GB 内存中运行?

交付物: 一个看起来像 $a \cdot \text{batch_size} + b$ 的表达式, 用于数值 a, b , 以及一个表示最大批处理大小的数字。

- (c) 运行 AdamW 的一步需要多少 FLOPs?

交付物: 一个代数表达式, 附带简要的说明。

- (d) 模型 FLOPs 利用率 (MFU) 定义为观察到的吞吐量 (每秒令牌数) 相对于硬件的理论峰值 FLOP 吞吐量 [Chowdhery 等人, 2022] 的比率。一个 NVIDIA A100 GPU 对于 float32 操作的理论峰值是 19.5 太拉 FLOP/s。假设你能获得 50% 的 MFU, 在单个 A100 上训练 GPT-2 XL 400K 步, 批处理大小为 1024, 需要多长时间? 遵循 Kaplan 等人 [2020] 和 Hoffmann 等人 [2022], 的假设, 反向传播的 FLOPs 是正向传播的两倍。

可交付成果: 训练所需的天数, 并附简要说明。

4.4 学习率调度

导致损失快速下降的学习率值在训练过程中往往会变化。在训练 Transformers 时, 通常使用学习率调度, 我们从一个较大的学习率开始, 在训练初期进行快速更新, 并随着模型训练逐渐衰减到较小的值⁸在本作业中, 我们将实现用于训练 LLaMA [Touvron 等人, 2023]所使用的余弦退火调度。

调度器本质上是一个函数, 它接受当前步骤 t 以及其他相关参数 (例如初始和最终学习率), 并返回在步骤 t 用于梯度更新的学习率。最简单的调度是常数函数, 它将始终返回相同的学习率, 无论给定任何 t 。

余弦退火学习率调度采用 (i) 当前迭代次数 t , (ii) 最大学习率 α_{\max} , (iii) 最小 (最终) 学习率 α_{\min} , (iv) 预热迭代次数 T_w , 以及 (v) 余弦退火迭代次数 T_c 。迭代次数 t 的学习率定义为:

(预热) 如果 $t < T_w$, 则 $\alpha_t = \frac{t}{T_w} \alpha_{\max}$ 。

(余弦退火) 如果 $T_w \leq t \leq T_c$, 则 $\alpha_t = \alpha_{\min} + \frac{1}{2} \left(1 + \cos \left(\frac{t - T_w}{T_c - T_w} \pi \right) \right) (\alpha_{\max} - \alpha_{\min})$ 。

(退火后) 如果 $t > T_c$, 那么 $\alpha_t = \alpha_{\min}$ 。

⁸It's sometimes common to use a schedule where the learning rate rises back up (restarts) to help get past local minima.

⁸有时常见使用一个学习率回升 (重启) 的调度方案, 以帮助越过局部最小值。

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

Write a function that takes t , α_{\max} , α_{\min} , T_w and 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`.

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 (1 point)

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`.

问题 (learning_rate_schedule): 实现带预热的余弦学习率调度

编写一个函数，该函数接受 t 、 α_{\max} 、 α_{\min} 、 T_w 和 T_c ，并根据上述定义的调度程序返回学习率 α_t 。然后实现 `[adapters.get_lr_cosine_schedule]` 并确保它通过 `uv run pytest -k test_get_lr_cosine_schedule`。

4.5 梯度裁剪

在训练过程中，我们有时会遇到产生较大梯度的训练样本，这可能会使训练不稳定。为了缓解这种情况，实践中常用的一种技术是梯度裁剪。其思想是在每次反向传播后、在执行优化器步骤之前，对梯度的范数进行限制。

给定梯度（所有参数的梯度） g ，我们计算其 ℓ_2 -范数 $\|g\|_2$ 。如果该范数小于一个最大值 M ，则将 g 保持不变；否则，我们将 g 按因子 $\frac{M}{\|g\|_2 + \epsilon}$ 进行缩放（其中添加一个小的 ϵ ，如 10^{-6} ，以保持数值稳定性）。请注意，最终得到的范数将略低于 M 。

问题 (gradient_clipping): 实现梯度裁剪 (1 分)

编写一个实现梯度裁剪的函数。你的函数应该接受一个参数列表和一个最大 ℓ_2 -范数。它应该原地修改每个参数梯度。使用 $\epsilon = 10^{-6}$ (PyTorch 默认)。然后，实现适配器 `[adapters.run_gradient_clipping]` 并确保它通过 `uv run pytest -k test_gradient_clipping`。

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 (e.g., that you prepared in `tokenizer_experiments`) 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 $1 \leq i < n - m$ 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 (2 points)

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 $(\text{batch_size}, \text{context_length})$ 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` to test your implementation.

Low-Resource/Downscaling Tip: Data loading on CPU or Apple Silicon

If you are planning to train your LM on CPU or Apple Silicon, you need to move your data to the correct device (and similarly, you should use the same device for your model later on). If you are on CPU, you can use the `'cpu'` device string, and on Apple Silicon (M* chips), you can use the `'mps'` device string. For more on MPS, checkout these resources:

- <https://developer.apple.com/metal/pytorch/>
- <https://pytorch.org/docs/main/notes/mps.html>

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 memory-mapped 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 训练循环

我们现在终于要把到目前为止构建的主要组件组合起来：标记化的数据、模型和优化器。

5.1 数据加载器

分词后的数据（例如，你在 `tokenizer_experiments` 中准备的数据）是一个由 token 组成的单个序列 $x = (x_1, \dots, x_n)$ 。尽管源数据可能由单独的文档组成（例如，不同的网页或源代码文件），但一种常见的做法是将它们全部连接成一个由 token 组成的单个序列，并在它们之间添加一个分隔符（例如，`<|endoftext|>` token）。

数据加载器将这转换为一批次的流，其中每一批包含 B 长度为 m 的序列，并配对相应的下一个 token，其长度也为 m 。例如，对于 $B = 1, m = 3$, $([x_2, x_3, x_4], [x_3, x_4, x_5])$ 将是一个潜在批次。

以这种方式加载数据简化了训练的原因有多个。首先，任何 $1 \leq i < n - m$ 都给出一个有效的训练序列，因此采样序列非常简单。由于所有训练序列的长度相同，因此无需填充输入序列，这提高了硬件利用率（同时通过增加批次大小 B ）。最后，我们也不需要完全加载整个数据集来采样训练数据，这使得处理可能无法全部装入内存的大型数据集变得容易。

问题 (data_loading): 实现数据加载 (2分)

交付物：编写一个函数，该函数接受一个 numpy 数组 x （包含 token ID 的整数数组）、`a batch_size`、一个 `context_length` 以及一个 PyTorch 设备字符串（例如，`'cpu'` 或 `'cuda:0'`），并返回一对张量：采样的输入序列和相应的下一个 token 目标。这两个张量应该具有形状 $(\text{batch_size}, \text{context_length})$ ，包含 token ID，并且都应该放置在请求的设备上。要针对我们提供的测试测试你的实现，你首先需要实现 `[adapters.run_get_batch]` 中的测试适配器。然后，运行 `uv run pytest -k test_get_batch` 来测试你的实现。

低资源/降采样提示：在 CPU 或 Apple Silicon 上加载数据

如果你计划在 CPU 或 Apple Silicon 上训练你的 LM，你需要将你的数据移动到正确的设备上（同样，你以后应该使用相同的设备来训练你的模型）。如果你在 CPU 上，你可以使用 `'cpu'` 设备字符串，而在 Apple Silicon (M* 芯片) 上，你可以使用 `'mps'` 设备字符串。For more on MPS, checkout these resources:

- <https://developer.apple.com/metal/pytorch/>
- <https://pytorch.org/docs/main/notes/mps.html>

如果数据集太大而无法加载到内存中怎么办？我们可以使用一个名为 `mmap` 的 Unix 系统调用，它将磁盘上的文件映射到虚拟内存，并在访问该内存位置时才懒加载文件内容。因此，你可以“假装”你拥有整个数据集在内存中。Numpy 通过 `np.memmap`（或者如果你最初用 `np.save` 保存数组，则通过标志 `mmap_mode='r'` 到 `np.load`）实现这一点，它将返回一个类似 numpy 数组的对象，该对象在访问时按需加载条目。在训练期间从你的数据集（即 numpy 数组）中采样时，请确保以内存映射模式加载数据集（通过 `np.memmap` 或标志 `mmap_mode='r'` 到 `np.load`，具体取决于你如何保存数组）。确保你还指定了一个与你要加载的数组匹配的 `dtype`。最好明确验证内存映射数据是否正确（例如，不包含超出预期词汇大小的值）。

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 (1 point)

Implement the following two functions to load and save checkpoints:

`def save_checkpoint(model, optimizer, iteration, out)` should dump all the state from the first three parameters into the file-like object `out`. You can use the `state_dict` method of both the model and the optimizer to get their relevant states and use `torch.save(obj, out)` to dump `obj` into `out` (PyTorch supports either a path or a file-like object here). A typical choice is to have `obj` be a dictionary, but you can use whatever format you want as long as you can load your checkpoint later.

This function expects the following 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)` should load a checkpoint from `src` (path or file-like object), and then recover the model and optimizer states from that checkpoint. Your function should return the iteration number that was saved to the checkpoint. You can use `torch.load(src)` to recover what you saved in your `save_checkpoint` implementation, and the `load_state_dict` method in both the model and optimizers to return them to their previous states.

This function expects the following parameters:

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

Implement the `[adapters.run_save_checkpoint]` and `[adapters.run_load_checkpoint]` adapters, and make sure they pass `uv run pytest -k test_checkpointing`.

5.2 检查点

除了加载数据外，我们还需要在训练过程中保存模型。在运行作业时，我们通常希望能够恢复中途停止的训练运行（例如，由于作业超时、机器故障等原因）。即使一切顺利，我们可能还想稍后访问中间模型（例如，为了事后研究训练动态、从训练不同阶段的模型中采样等）。

一个检查点应该包含我们恢复训练所需的所有状态。我们当然希望至少能够恢复模型权重。如果使用有状态优化器（例如 AdamW），我们还需要保存优化器的状态（例如，在 AdamW 的情况下，是动量估计）。最后，为了恢复学习率计划，我们需要知道我们停止时的迭代次数。PyTorch 使保存所有这些变得容易：每个 `nn.Module` 都有一个 `state_dict()` 方法，该方法返回一个包含所有可学习权重的字典；我们可以稍后使用姐妹方法 `load_state_dict()` 恢复这些权重。任何 `nn.optim.Optimizer` 都是同样的情况。最后，`torch.save(obj, dest)` 可以将一个对象（例如，一个包含张量作为某些值的字典，但也包含像整数这样的常规 Python 对象）转储到文件（路径）或文件类对象中，然后可以使用 `torch.load(src)` 将其重新加载到内存中。

问题 (checkpointing)：实现模型检查点 (1分)

实现以下两个函数以加载和保存检查点：

`def save_checkpoint(model, optimizer, iteration, out)` 应该将前三个参数的所有状态序列化到文件对象 `out` 中。你可以使用模型和优化器的 `state_dict` 方法来获取它们的相关状态，并使用 `torch.save(obj, out)` 将 `obj` 序列化到 `out` 中（PyTorch 支持路径或文件对象）。一个典型选择是让 `obj` 是一个字典，但你可以使用任何你想要的格式，只要以后可以加载你的检查点。

此函数期望以下参数：

`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)` 应该从 `src`（路径或文件对象）加载检查点，然后从该检查点恢复模型和优化器状态。你的函数应该返回保存到检查点的迭代次数。你可以使用 `torch.load(src)` 来恢复你在 `save_checkpoint` 实现中保存的内容，并使用模型和优化器中的 `load_state_dict` 方法将它们恢复到之前的状态。

此函数期望以下参数：

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

Implement the `[adapters.run_save_checkpoint]` and `[adapters.run_load_checkpoint]` adapters, and make sure they pass `uv run pytest -k test_checkpointing`.

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 (4 points)

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 training and validation large 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 and Biases).^a

^awandb.ai

5.3 训练循环

现在，终于到了将您实现的所有组件整合到主训练脚本中的时候了。由于您之后会多次进行这些操作以研究不同选择如何影响训练，因此让使用不同的超参数（例如，通过将它们作为命令行参数）轻松启动训练运行将是有益的。

问题 (training_together): 组合起来 (4分)

交付物: 编写一个脚本，运行训练循环来训练你的模型，使用用户提供的输入。特别是，我们建议你的训练脚本允许（至少）以下功能:

- 配置和控制各种模型和优化器超参数的能力。
- 使用 `np.memmap` 高效加载训练和验证大型数据集。
- 将检查点序列化到用户提供的路径。
- 定期记录训练和验证性能（例如，输出到控制台和/或外部类似 Weights and Biases 的服务）。^a

^awandb.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 × 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 .

6 生成文本

现在我们可以训练模型了，我们还需要的是从模型中生成文本的能力。回想一下，语言模型接收一个长度为 (`sequence_length`) 的（可能批量的）整数序列，并生成一个大小为 (`sequence_length × vocab_size`) 的矩阵，其中序列中的每个元素都是一个概率分布，预测该位置之后的下一个词。我们将编写几个函数，将这转换为生成新序列的采样方案。

Softmax 标准惯例中，语言模型的输出是最终线性层的输出（“logits”），因此我们需要通过 softmax 操作将其转换为归一化概率，正如我们之前在公式10中看到的。

解码 要从我们的模型中生成文本（解码），我们将向模型提供一系列前缀标记（“提示”），并要求它对词汇表生成一个概率分布，以预测序列中的下一个词。然后，我们将从这个词汇项的概率分布中进行采样，以确定下一个输出标记。

具体来说，解码过程的一步应该接收一个序列 $x_{1...t}$ 并通过以下公式返回一个标记 x_{t+1} ，

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

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

其中 TransformerLM 是我们的模型，它将 `sequence_length` 序列作为输入，并产生一个大小为 (`sequence_length × vocab_size`) 的矩阵，我们取该矩阵的最后一个元素，因为我们正在寻找在 t -th 位置的下个词预测。

这给我们提供了一个基本的解码器，通过反复从这些单步条件中采样（将我们之前生成的输出 token 添加到下一个解码时间步的输入中），直到我们生成序列结束 token `<|endoftext|>`（或用户指定的最大 token 数量）。

解码器技巧 我们将使用小型模型进行实验，小型模型有时会生成非常低质量的文本。两种简单的解码器技巧可以帮助解决这些问题。首先，在温度缩放中，我们使用温度参数 τ 修改我们的 softmax，其中新的 softmax 是

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

注意设置 $\tau \rightarrow 0$ 使得 v 的最大元素占主导地位，softmax 的输出变成一个集中在该最大元素上的单热向量。

其次，另一种技巧是核或 top-p 采样，其中我们通过截断低概率词来修改采样分布。设 q 是我们从（温度缩放的）softmax 得到的概率分布，大小为 (`vocab_size`)。核采样使用超参数 p 产生下一个 token，

到方程式

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

where $V(p)$ 是最小的索引集合，使得 $\sum_{j \in V(p)} q_j \geq p$ 。你可以通过首先按大小对概率分布 q 进行排序，并选择最大的词汇元素，直到达到目标级别 α ，轻松计算这个量。

Problem (decoding): Decoding (3 points)

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., 2020; also referred to as nucleus sampling), given a user-specified threshold value.

问题(decoding): 解码 (3分)

交付物：实现一个从你的语言模型解码的函数。我们建议你支持以下功能：

- 为用户提供提示生成补全 (即，接收一些 $x_{1...t}$ 并采样补全，直到命中一个 `<|endoftext|>` 标记)。
- 允许用户控制生成的最大标记数。
- 给定期望的温度值，在采样前对预测的下一个词分布应用softmax温度缩放。
- Top-p采样 (Holtzman等人，2020年；也称为核采样)，给定用户指定的阈值。

7 Experiments

Now it is time to put everything together and train (small) language models on a pretaining 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 (17M 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 (3 points)

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 10000. Typical vocabulary sizes are in the tens to hundreds of thousands. You should vary this and see how the vocabulary and model behavior changes.

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.

7 实验

现在，是时候将所有内容整合起来，并在预处理的数据集上训练（小型）语言模型了。

7.1 如何运行实验和交付成果

要理解 Transformer 架构组件背后的原理，最好的方法就是实际修改它并亲自运行。没有比动手实践更好的替代方案了。

为此，能够快速、一致地进行实验并记录你所做的事情非常重要。为了快速实验，我们将在一个小型模型（17M 参数）和简单数据集（TinyStories）上运行许多实验。为了保持一致性，你将系统地消融组件并调整超参数，为了记录结果，我们将要求你提交每个实验的日志和学习曲线。

为了能够提交损失曲线，请确保定期评估验证损失，并记录步骤数和实际运行时间。你可能会发现 Weights and Biases 等日志基础设施很有帮助。

问题 (experiment_log): 实验记录 (3分)

为您的训练和评估代码，创建实验跟踪基础设施，以便您能够根据梯度步数和墙上时间跟踪实验和损失曲线。

交付物：为您的实验准备的日志基础设施代码，以及本节下方作业问题中所有您尝试过的内容的实验日志（一个文档）。

7.2 TinyStories

我们将从一个非常简单的数据集（TinyStories；Eldan和Li，2023）开始，其中模型可以快速训练，并且我们可以观察到一些有趣的行为。获取这个数据集的说明在1节。这个数据集的一个示例如下。

一个来自 TinyStories 的例子

从前有个名叫Ben的小男孩。Ben喜欢探索他周围的世界。他看到了很多奇妙的东西，比如商店里陈列的漂亮花瓶。有一天，Ben在商店里走着时，偶然发现了一个非常特别的花瓶。当Ben看到它时，他感到非常惊讶！他说：“哇，那真是一个非常奇妙的花瓶！我能买它吗？”店主微笑着说：“当然可以。你可以把它带回家，向所有朋友展示它有多奇妙！”所以Ben把花瓶带回家，他对此感到非常自豪！他叫来他的朋友们，向他们展示了那个奇妙的花瓶。他所有的朋友都认为花瓶很漂亮，简直不敢相信Ben有多幸运。这就是Ben在商店里找到那个奇妙花瓶的故事！

超参数调整 我们会告诉你一些非常基本的超参数来开始，并请你找到其他一些效果好的设置。

vocab_size 10000。典型的词汇量大小通常在几十万到上百万。你应该改变这个值，看看词汇量和模型行为如何变化。

context_length 256。简单数据集如 TinyStories 可能不需要很长的序列长度，但对于后面的 OpenWebText 数据，你可能需要调整这个值。尝试调整这个值并观察对迭代运行时和最终困惑度。

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

RoPE theta parameter Θ 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 327,680,000 (your batch size \times total step count \times context length should equal roughly this value).

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. **Important note:** If your implementation is correct and efficient, the above hyperparameters should result in a roughly 30-40 minute runtime on 1 H100 GPU. If you have runtimes that are much longer, please check and make sure your dataloading, checkpointing, or validation loss code is not bottlenecking your runtimes and that your implementation is properly batched.

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 (3 points) (4 H100 hrs)

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

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

d_model 512. 这比许多小型 Transformer 论文中使用的 768 维度稍微小一些，但这会使事情更快。

d_ff 1344。这大约是 $\frac{8}{3}d_{\text{model}}$ ，同时是 64 的倍数，这对 GPU 性能很好。

RoPE theta parameter Θ 10000。

number of layers and heads 4 层，16 个头。总而言之，这将提供大约 17M 个非嵌入参数，这是一个相当小的 Transformer。

total tokens processed 327, 680,000（你的批大小 \times 总步数 \times 上下文长度应该等于大约这个值）。

你应该通过一些尝试和错误来为以下其他超参数找到好的默认值：**learning rate**, **learning rate warmup**, **other AdamW hyperparameters** ($\beta_1, \beta_2, \epsilon$), 和 **weight decay**。你可以在 Kingma 和 Ba [2015] 中找到这些超参数的一些典型选择。

将它们结合起来 现在，你可以通过获取一个训练好的 BPE 分词器，对训练数据集进行分词，并在你编写的训练循环中运行这些操作来将所有内容结合起来。重要提示：如果你的实现是正确且高效的，上述超参数应该在 1 个 H100 GPU 上产生大约 30-40 分钟的运行时间。如果你有远超这个时间的运行时间，请检查并确保你的数据加载、检查点或验证损失代码没有成为运行时间的瓶颈，并且你的实现是正确批处理的。

调试模型架构的技巧和窍门 我们强烈建议你熟悉你的 IDE 的内置调试器（例如 VSCode/PyCharm），这比使用打印语句调试要节省时间。如果你使用文本编辑器，你可以使用类似 `pdb` 的东西。在调试模型架构时，还有一些其他的好习惯：

- 在开发任何神经网络架构时，一个常见的第一步是过拟合到一个单独的小批量。如果你的实现是正确的，你应该能够快速将训练损失驱动到接近零。
- 在各种模型组件中设置调试断点，并检查中间张量的形状以确保它们符合你的预期。
- 监控激活值、模型权重和梯度的范数，以确保它们没有爆炸或消失。

问题 (learning_rate): 调整学习率 (3分) (4 H100 小时)

学习率是调整最重要的超参数之一。针对你已训练的基础模型，回答以下问题：

- (a) 对学习率进行超参数扫描，并报告最终损失（如果优化器发散，请注明）。

交付物：与多个学习率相关的学习曲线。解释你的超参数搜索策略。

可交付成果：在 TinyStories 上的验证损失（按标记）最多为 1.45

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

If you are running on `cpu` or `mps`, 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 $\text{batch size} \times \text{total step count} \times \text{context length} = 32 \times 5000 \times 256 = 40,960,000$ 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, we suggest adjusting the cosine learning rate decay schedule to terminate its decay (i.e., reach the minimum learning rate) at precisely step X .

- When using `mps`, do **not** use TF32 kernels, i.e., do **not** set

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

as you might with `cuda` devices. We tried enabling TF32 kernels with `mps` (`torch` version 2.6.0) and found the backend will use silently broken kernels that cause unstable training.

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

- On `cpu`, compile your model with

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

- On `mps`, you can somewhat optimize the backward pass using

```
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.

Now let’s vary the batch size and see what happens to training. Batch sizes are important – they let us get higher efficiency from our GPUs by doing larger matrix multiplies, but is it true that we always want batch sizes to be large? Let’s run some experiments to find out.

Problem (batch_size_experiment): Batch size variations (1 point) (2 H100 hrs)

Vary your batch size all the way from 1 to the GPU memory limit. Try at least a few batch sizes in between, including typical sizes like 64 and 128.

Deliverable: Learning curves for runs with different batch sizes. The learning rates should be optimized again if necessary.

Deliverable: A few sentences discussing of your findings on batch sizes and their impacts on training.

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.

低资源/降尺度提示：在 CPU 或 Apple Silicon 上进行少量步骤的训练

如果你在 `cpu` 或 `mps` 上运行，你应该将处理的标记总数减少到 40,000,000，这将足以生成相对流畅的文本。你也可以将目标验证损失从 1.45 增加到 2.00。

在 M3 Max 芯片和 36 GB 内存上使用调整后的学习率运行我们的解决方案代码，我们使用批量大小 \times 总步数 \times 上下文长度 $= 32 \times 5000 \times 256 = 40,960,000$ 个 token，在 `cpu` 上耗时 1 小时 22 分钟，在 `mps` 上耗时 36 分钟。在步骤 5000 时，我们实现了 1.80 的验证损失。

一些额外的提示：

- 在使用 X 训练步骤时，我们建议调整余弦学习率衰减计划，使其在精确的步骤 X 结束衰减（即达到最小学习率）。

- 在使用 `mps` 时，不要使用 TF32 内核，即不要设置

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

就像你可能在使用 `cuda` 设备时那样。我们尝试使用 `mps` (`torch` 版本 2.6.0) 启用 TF32 内核，并发现后端会静默地使用损坏的内核，导致训练不稳定。

- 您可以通过使用 `torch.compile` 对模型进行 JIT 编译来加快训练速度。具体来说：

- 在 `cpu` 上，使用

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

- 在 `mps` 上，您可以使用

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

在 `mps` 上，使用 Inductor 进行编译自 `torch` 版本 2.6.0 起不受支持。

- (b) 民间智慧认为最佳学习率是“处于稳定边缘”。调查学习率发散点与您的最佳学习率之间的关系。学习率发散点与您的最佳学习率之间的关系。

可交付成果：学习率递增的学习曲线，其中至少包含一个发散的运行，以及分析其与收敛速度的关系。

现在让我们改变批大小，看看这对训练有什么影响。批大小很重要——它们通过执行更大的矩阵乘法，让我们能够从我们的 GPU 中获得更高的效率，但是批大小是否总是要大是真的吗？让我们做一些实验来找出答案。

问题 (batch_size_experiment)：批大小变化 (1分) (2 H100 小时)

将你的批大小从 1 一直变化到 GPU 内存限制。尝试在中间至少有几个批大小，包括典型的 64 和 128 大小。

可交付成果：不同批大小的运行的学习曲线。如果需要，应该再次优化学习率。

可交付成果：几句话讨论您关于批处理大小及其对训练影响的发现。

手握你的解码器，我们现在可以生成文本了！我们将从模型中生成文本，看看它有多好。作为参考，你应该得到至少与下面示例一样好的输出。

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 or Apple Silicon

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 (1 point)

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

7.3 Ablations and architecture modification

The best way to understand the Transformer is to actually modify it and see how it behaves. We will now do a few simple ablations and modifications.

Ablation 1: layer normalization It is often said that layer normalization is important for the stability of Transformer training. But perhaps we want to live dangerously. Let’s remove RMSNorm from each of our Transformer blocks and see what happens.

Problem (layer_norm_ablation): Remove RMSNorm and train (1 point) (1 H100 hr)

Remove all of the RMSNorms from your Transformer and train. What happens at the previous optimal learning rate? Can you get stability by using a lower learning rate?
Deliverable: A learning curve for when you remove RMSNorms and train, as well as a learning curve for the best learning rate.
Deliverable: A few sentence commentary on the impact of RMSNorm.

示例 (ts_generate_example): 来自 TinyStories 语言模型的样本输出

从前，有一个名叫莉莉的漂亮女孩。她喜欢嚼口香糖，尤其是那个大黑糖。有一天，莉莉的妈妈让她帮忙做饭。莉莉非常兴奋！她喜欢帮助妈妈。莉莉的妈妈做了一大锅汤当晚餐。莉莉非常高兴，说：“谢谢你，妈妈！我爱你。”她帮妈妈把汤倒进一个大碗里。晚餐后，莉莉的妈妈又做了一些美味的汤。莉莉非常喜欢！她说：“谢谢你，妈妈！这汤太好吃了！”她妈妈微笑着说：“很高兴你喜欢，莉莉。”他们做完饭，继续一起做饭。结束。

低资源/降采样提示：在 CPU 或 Apple Silicon 上生成文本

如果你使用了低资源配置并处理了 40M 个 token，你应该看到仍然类似英语但不如上述流畅的生成。例如，我们在 40M 个 token 上训练的 TinyStories 语言模型的样本输出如下：

从前，有一个名叫 Sue 的小女孩。Sue 有一颗她非常喜欢的牙齿。那是她最好的头。有一天，Sue 去散步时遇到了一只瓢虫！他们成了好朋友，一起在小路上玩耍。“嘿，波莉！我们去外面玩吧！”Tim 说。Sue 看着天空，发现很难找到跳舞的明亮方式。她微笑着同意帮助说话！”当 Sue 看着天空移动时，她在想什么。她

以下是精确的问题陈述以及我们要求的内容：

问题(generate)：生成文本 (1分)

使用你的解码器和你的训练好的检查点，报告你的模型生成的文本。你可能需要调整解码器参数（温度、top-p 等）来获得流畅的输出。
交付物：至少 256 个 token 的文本（或直到第一个 <|endoftext|> token），以及对该输出流畅性的简要评论，以及至少两个影响该输出好坏的因素。

7.3 消融和架构修改

理解 Transformer 的最佳方法实际上就是修改它并观察它的行为。我们现在将进行一些简单的消融和修改。

消融 1: 层归一化 人们常说层归一化对 Transformer 训练的稳定性很重要。但也许我们想冒点风险。让我们从每个 Transformer 模块中移除 RMSNorm，看看会发生什么。

问题 (layer_norm_ablation): 移除 RMSNorm 并训练 (1 分) (1 H100 小时)

从你的 Transformer 中移除所有的 RMSNorm 并训练。在之前的最优学习率下会发生什么？你能通过使用更低的学习率来获得稳定性吗？
交付物：移除 RMSNorm 并训练时的学习曲线，以及最佳学习率的学习曲线。

可交付成果：关于 RMSNorm 影响的简短评论。

Let’s now investigate another layer normalization choice that seems arbitrary at first glance. *Pre-norm* Transformer blocks are defined as

$$\begin{aligned} z &= x + \text{MultiHeadedSelfAttention}(\text{RMSNorm}(x)) \\ y &= z + \text{FFN}(\text{RMSNorm}(z)). \end{aligned}$$

This is one of the few ‘consensus’ modifications to the original Transformer architecture, which used a *post-norm* approach as

$$\begin{aligned} z &= \text{RMSNorm}(x + \text{MultiHeadedSelfAttention}(x)) \\ y &= \text{RMSNorm}(z + \text{FFN}(z)). \end{aligned}$$

Let’s revert back to the *post-norm* approach and see what happens.

Problem (pre_norm_ablation): Implement post-norm and train (1 point) (1 H100 hr)

Modify your pre-norm Transformer implementation into a post-norm one. Train with the post-norm model and see what happens.

Deliverable: A learning curve for a post-norm transformer, compared to the pre-norm one.

We see that layer normalization has a major impact on the behavior of the transformer, and that even the position of the layer normalization is important.

Ablation 2: position embeddings We will next investigate the impact of the position embeddings on the performance of the model. Specifically, we will compare our base model (with RoPE) with not including position embeddings at all (NoPE). It turns out that decoder-only transformers, i.e., those with a causal mask as we have implemented, can in theory infer relative or absolute position information without being provided with position embeddings explicitly [Tsai et al., 2019, Kazemnejad et al., 2023]. We will now test empirically how NoPE performs compare to RoPE.

Problem (no_pos_emb): Implement NoPE (1 point) (1 H100 hr)

Modify your Transformer implementation with RoPE to remove the position embedding information entirely, and see what happens.

Deliverable: A learning curve comparing the performance of RoPE and NoPE.

Ablation 3: SwiGLU vs. SiLU Next, we will follow Shazeer [2020] and test the importance of gating in the feed-forward network, by comparing the performance of SwiGLU feed-forward networks versus feed-forward networks using SiLU activations but no gated linear unit (GLU):

$$\text{FFN}_{\text{SiLU}}(x) = W_2 \text{SiLU}(W_1 x). \tag{25}$$

Recall that in our SwiGLU implementation, we set the dimensionality of the inner feed-forward layer to be roughly $d_{\text{ff}} = \frac{8}{3}d_{\text{model}}$ (while ensuring that $d_{\text{ff}} \bmod 64 = 0$, to make use of GPU tensor cores). In your FFN_{SiLU} implementation you should set $d_{\text{ff}} = 4 \times d_{\text{model}}$, to approximately match the parameter count of the SwiGLU feed-forward network (which has three instead of two weight matrices).

Problem (swiglu_ablation): SwiGLU vs. SiLU (1 point) (1 H100 hr)

Deliverable: A learning curve comparing the performance of SwiGLU and SiLU feed-forward networks, with approximately matched parameter counts.

现在让我们研究另一种看起来任意的层归一化选择。预归一化 Transformer 模块定义为

$$\begin{aligned} z &= x + \text{MultiHeadedSelfAttention}(\text{RMSNorm}(x)) \\ y &= z + \text{FFN}(\text{RMSNorm}(z)). \end{aligned}$$

这是原始 Transformer 架构中少数几个 ‘共识’ 修改之一，该架构使用了后规范方法as

$$\begin{aligned} z &= \text{RMSNorm}(x + \text{MultiHeadedSelfAttention}(x)) \\ y &= \text{RMSNorm}(z + \text{FFN}(z)). \end{aligned}$$

让我们回到后规范方法，看看会发生什么。

问题 (pre_norm_ablation): 实现后规范并训练 (1 分) (1 H100 小时)

将你的预规范 Transformer 实现修改为后规范的一个。使用后规范模型并看看会发生什么。

交付物：后规范 Transformer 的学习曲线，与预规范的一个进行比较。

我们看到层归一化对 Transformer 的行为有重大影响，而且层归一化的位置也很重要。

消融实验 2：位置嵌入 我们接下来将研究位置嵌入对模型性能的影响。具体来说，我们将比较我们的基础模型（带有 RoPE）与完全不包括位置嵌入（NoPE）。结果表明，仅解码器 Transformer，即我们实现的带有因果掩码的模型，理论上可以在不显式提供位置嵌入的情况下推断相对或绝对位置信息 [Tsai 等人, 2019, Kazemnejad 等人, 2023]。我们现在将通过实验测试 NoPE 与 RoPE 的性能比较。

问题 (no_pos_emb): 实现 NoPE (1 分) (1 H100 小时)

修改你的带有 RoPE 的 Transformer 实现，以完全移除位置嵌入信息看看会发生什么。

交付物：比较 RoPE 和 NoPE 性能的学习曲线。

消融实验 3：SwiGLU 与 SiLU 下一步，我们将遵循 Shazeer [2020] 并测试门控在前馈网络中的重要性，通过比较 SwiGLU 前馈网络的性能与使用 SiLU 激活但没有门控线性单元（GLU）的前馈网络：

$$\text{FFN}_{\text{SiLU}}(x) = W_2 \text{SiLU}(W_1 x). \tag{25}$$

回想一下，在我们的 SwiGLU 实现中，我们将内部前馈层的维度设置为大约 $d_{\text{ff}} = \frac{8}{3}d_{\text{model}}$ （同时确保 $d_{\text{ff}} \bmod 64 = 0$ ，以利用 GPU 张量核）。在您的 FFN_{SiLU} 实现中，您应该设置 $d_{\text{ff}} = 4 \times d_{\text{model}}$ ，以大约匹配 SwiGLU 前馈网络的参数数量（该网络有三个而不是两个权重矩阵）。

问题 (swiglu_ablation): SwiGLU 与 SiLU (1 分) (1 H100 小时)

可交付成果：一个比较 SwiGLU 和 SiLU 前馈网络性能的学习曲线，参数数量大致匹配。

Deliverable: A few sentences discussing your findings.

Low-Resource/Downscaling Tip: Online students with limited GPU resources should test modifications on TinyStories

In the remainder of the assignment, we will move to a larger-scale, noisier web dataset (OpenWebText), experimenting with architecture modifications and (optionally) making a submission to the course leaderboard.

It takes a long time to train an LM to fluency on OpenWebText, so we suggest that online students with limited GPU access continue testing modifications on TinyStories (using validation loss as a metric to evaluate performance).

7.4 Running on OpenWebText

We will now move to a more standard pretraining dataset created from a webcrawl. A small sample of OpenWebText [Gokaslan et al., 2019] is also provided as a single text file: see section 1 for how to access this file.

Here is an example from OpenWebText. Note how the text is much more realistic, complex, and varied. You may want to look through the training dataset to get a sense of what training data looks like for a webscraped corpus.

Example (owt_example): One example from OWT

Baseball Prospectus director of technology Harry Pavlidis took a risk when he hired Jonathan Judge. Pavlidis knew that, as Alan Schwarz wrote in The Numbers Game, “no corner of American culture is more precisely counted, more passionately quantified, than performances of baseball players.” With a few clicks here and there, you can findout that Noah Syndergaard’s fastball revolves more than 2,100 times per minute on its way to the plate, that Nelson Cruz had the game’s highest average exit velocity among qualified hitters in 2016 and myriad other tidbits that seem ripped from a video game or science fiction novel. The rising ocean of data has empowered an increasingly important actor in baseball’s culture: the analytical hobbyist.

That empowerment comes with added scrutiny – on the measurements, but also on the people and publications behind them. With Baseball Prospectus, Pavlidis knew all about the backlash that accompanies quantitative imperfection. He also knew the site’s catching metrics needed to be reworked, and that it would take a learned mind – someone who could tackle complex statistical modeling problems – to complete the job.

“He freaks us out.” Harry Pavlidis

Pavlidis had a hunch that Judge “got it” based on the latter’s writing and their interaction at a site-sponsored ballpark event. Soon thereafter, the two talked over drinks. Pavlidis’ intuition was validated. Judge was a fit for the position – better yet, he was a willing fit. “I spoke to a lot of people,” Pavlidis said, “he was the only one brave enough to take it on.” [...]

Note: You may have to re-tune your hyperparameters such as learning rate or batch size for this experiment.

Problem (main_experiment): Experiment on OWT (2 points) (3 H100 hrs)

Train your language model on OpenWebText with the same model architecture and total training iterations as TinyStories. How well does this model do?

Deliverable: A learning curve of your language model on OpenWebText. Describe the difference in losses from TinyStories – how should we interpret these losses?

可交付成果：用几句话讨论你的发现。

低资源/降采样提示：GPU 资源有限的在线学生应在 TinyStories 上测试修改

在作业的其余部分，我们将转向一个更大规模、更嘈杂的网页数据集（Open-WebText），尝试进行架构修改，并（可选地）向课程排行榜提交成绩。

在OpenWebText上训练一个语言模型达到熟练状态需要很长时间，因此我们建议GPU访问受限的在线学生继续在TinyStories上测试修改（使用验证损失作为评估性能的指标）。

7.4 在OpenWebText上运行

我们现在将转向一个从网页爬取中创建的更标准的预训练数据集。OpenWebText的一个小样本 [Gokaslan et al. 2019] 也作为单个文本文件提供：有关如何访问此文件的说明，请参见第1节。

这里有一个来自OpenWebText的示例。注意文本是多么更真实、复杂和多样化。你可能想浏览训练数据集，以了解网页爬取语料库的训练数据是什么样的。

示例 (owt_example): 来自OWT的一个示例

Baseball Prospectus技术总监Harry Pavlidis在聘用Jonathan Judge时冒了险。Pavlidis知道，正如Alan Schwarz在《The Numbers Game》中所写，“美国文化中没有哪个角落被如此精确地统计，如此热情地量化，就像棒球运动员的表现一样。”在这里和那里点击几下，你就能发现诺亚·辛德伯格的快速球在飞向本垒时每分钟旋转超过2,100次，内尔森·克鲁斯在2016年合格击球手中的平均出球速度最高，以及其他无数似乎来自电子游戏或科幻小说的细节。不断上涨的数据海洋赋予了棒球文化中一个越来越重要的角色：分析爱好者。

这种赋权伴随着额外的审查——不仅对测量方法，也对背后的个人和出版物。对于Baseball Prospectus, Pavlidis深知定量不完美所带来的反噬。他也知道该网站的本垒打指标需要重新设计，并且需要一位博学的人——一个能够处理复杂统计建模问题的人——来完成这项工作。“他让我们害怕。” Harry Pavlidis

Pavlidis 有一个预感，认为法官“理解了”这是基于后者写的文字以及他们在场地赞助的棒球场活动中的互动。不久之后，两人边喝边谈。Pavlidis 的直觉得到了验证。法官适合这个职位——更妙的是，他愿意接受这个职位。“我谈了很多人的意见，” Pavlidis 说，“他是唯一敢接手的人。” [...]

注意：您可能需要重新调整此实验的超参数，例如学习率或批处理大小。

问题 (main_experiment): OWT 实验 (2 分) (3 H100 小时)

在 OpenWebText 上用与 TinyStories 相同的模型架构和总训练迭代次数训练您的语言模型。这个模型表现如何？

可交付成果：您的语言模型在 OpenWebText 上的学习曲线。描述与 TinyStories 的损失差异——我们应该如何解释这些损失？

Deliverable: Generated text from OpenWebText LM, in the same format as the TinyStories outputs. How is the fluency of this text? Why is the output quality worse even though we have the same model and compute budget as TinyStories?

7.5 Your own modification + leaderboard

Congratulations on getting to this point. You’re almost done! You will now try to improve upon the Transformer architecture, and see how your hyperparameters and architecture stack up against other students in the class.

Rules for the leaderboard There are no restrictions other than the following:

Runtime Your submission can run for at most 1.5 hours on an H100. You can enforce this by setting `--time=01:30:00` in your slurm submission script.

Data You may only use the OpenWebText training dataset that we provide.

Otherwise, you are free to do whatever your heart desires.
If you are looking for some ideas on what to implement, you can checkout some of these resources:

- State-of-the-art open-source LLM families, such as Llama 3 [Grattafiori et al., 2024] or Qwen 2.5 [Yang et al., 2024].
- The NanoGPT speedrun repository (<https://github.com/KellerJordan/modded-nanogpt>), where community members post many interesting modifications for “speedrunning” small-scale language model pretraining. For example, a common modification that dates back to the original Transformer paper is to tie the weights of the input and output embeddings together (see Vaswani et al. [2017] (Section 3.4) and Chowdhery et al. [2022] (Section 2)). If you do try weight tying, you may have to decrease the standard deviation of the embedding/LM head init.

You will want to test these on either a small subset of OpenWebText or on TinyStories before trying the full 1.5-hour run.

As a caveat, we do note that some of the modifications you may find working well in this leaderboard may not generalize to larger-scale pretraining. We will explore this idea further in the scaling laws unit of the course.

Problem (leaderboard): Leaderboard (6 points) (10 H100 hrs)

You will train a model under the leaderboard rules above with the goal of minimizing the validation loss of your language model within 1.5 H100-hour.

Deliverable: The final validation loss that was recorded, an associated learning curve that clearly shows a wallclock-time x-axis that is less than 1.5 hours and a description of what you did. We expect a leaderboard submission to beat at least the naive baseline of a 5.0 loss. Submit to the leaderboard here: <https://github.com/stanford-cs336/assignment1-basics-leaderboard>.

可交付成果：从 OpenWebText LM 生成的文本，与 TinyStories 输出格式相同。这段文本的流畅性如何？为什么即使我们使用相同的模型和计算预算，输出质量也更差？

7.5 您的修改 + 排行榜

恭喜您到达这一步。您几乎完成了！您现在将尝试改进 Transformer 架构，并看看您的超参数和架构与其他班级学生的表现如何。

排行榜规则 除了以下限制外，没有其他限制：

Runtime您的提交在 H100 上最多运行 1.5 小时。您可以通过在您的 slurm 提交脚本中设置`--time=01:30:00`来强制执行此限制。

Data 您只能使用我们提供的 OpenWebText 训练数据集。

否则，你可以随心所欲地做任何事。
如果你在寻找一些关于要实现的内容的想法，你可以查看这些资源：

- 最先进的开源LLM系列，例如Llama 3 [Grattafiori等人， 2024] 或Qwen 2.5 [Yang 等人， 2024]。
- NanoGPT速成仓库 (<https://github.com/KellerJordan/modded-nanogpt>)，社区成员在那里发布了许多有趣的“速成”小规模语言模型预训练修改。例如，一个追溯到原始Transformer论文的常见修改是将输入和输出嵌入的权重绑定在一起（参见Vaswani等人 [2017]（第3.4节）和Chowdhery等人 [2022]（第2节））。如果你尝试权重绑定，你可能需要降低嵌入/LM头的初始化标准差。

你需要在尝试完整的 1.5 小时运行之前，在这些数据集的小子集（OpenWebText 或 TinyStories）上测试这些模型。

需要注意的是，你在本次排行榜上发现的一些有效修改可能无法泛化到更大规模的预训练。我们将在课程的扩展规律单元中进一步探讨这一想法。

问题 (leaderboard): 排行榜 (6 分) (10 H100 小时)

你需要在排行榜规则下训练一个模型，目标是在 1.5 H100 小时内最小化你的语言模型的验证损失。

提交内容：最终记录的验证损失、一个清晰显示时钟时间 x 轴小于 1.5 小时的相关学习曲线，以及你所做的工作的描述。我们期望排行榜提交的成绩至少要优于 5.0 损失的朴素基线。请在此处提交排行榜: <https://github.com/stanford-cs336/assignment1-basics-leaderboard>。

References

Ronen Eldan and Yuanzhi Li. TinyStories: How small can language models be and still speak coherent English?, 2023. arXiv:2305.07759.

Aaron Gokaslan, Vanya Cohen, Ellie Pavlick, and Stefanie Tellex. OpenWebText corpus. <http://Skylion007.github.io/OpenWebTextCorpus>, 2019.

Rico Sennrich, Barry Haddow, and Alexandra Birch. Neural machine translation of rare words with subword units. In *Proc. of ACL*, 2016.

Changhan Wang, Kyunghyun Cho, and Jiatao Gu. Neural machine translation with byte-level subwords, 2019. arXiv:1909.03341.

Philip Gage. A new algorithm for data compression. *C Users Journal*, 12(2):23–38, February 1994. ISSN 0898-9788.

Alec Radford, Jeff Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. Language models are unsupervised multitask learners, 2019.

Alec Radford, Karthik Narasimhan, Tim Salimans, and Ilya Sutskever. Improving language understanding by generative pre-training, 2018.

Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Ł ukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *Proc. of NeurIPS*, 2017.

Toan Q. Nguyen and Julian Salazar. Transformers without tears: Improving the normalization of self-attention. In *Proc. of IWSWLT*, 2019.

Ruibin Xiong, Yunchang Yang, Di He, Kai Zheng, Shuxin Zheng, Chen Xing, Huishuai Zhang, Yanyan Lan, Liwei Wang, and Tie-Yan Liu. On layer normalization in the Transformer architecture. In *Proc. of ICML*, 2020.

Jimmy Lei Ba, Jamie Ryan Kiros, and Geoffrey E. Hinton. Layer normalization, 2016. arXiv:1607.06450.

Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, Aurelien Rodriguez, Armand Joulin, Edouard Grave, and Guillaume Lample. Llama: Open and efficient foundation language models, 2023. arXiv:2302.13971.

Biao Zhang and Rico Sennrich. Root mean square layer normalization. In *Proc. of NeurIPS*, 2019.

Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Alex Vaughan, Amy Yang, Angela Fan, Anirudh Goyal, Anthony Hartshorn, Aobo Yang, Archi Mitra, Archie Sravankumar, Artem Korenev, Arthur Hinsvark, Arun Rao, Aston Zhang, Aurelien Rodriguez, Austen Gregerson, Ava Spataru, Baptiste Roziere, Bethany Biron, Binh Tang, Bobbie Chern, Charlotte Caucheteux, Chaya Nayak, Chloe Bi, Chris Marra, Chris McConnell, Christian Keller, Christophe Touret, Chunyang Wu, Corinne Wong, Cristian Canton Ferrer, Cyrus Nikolaidis, Damien Allonsius, Daniel Song, Danielle Pintz, Danny Livshits, Danny Wyatt, David Esiobu, Dhruv Choudhary, Dhruv Mahajan, Diego Garcia-Olano, Diego Perino, Dieuwke Hupkes, Egor Lakomkin, Ehab AlBadawy, Elina Lobanova, Emily Dinan, Eric Michael Smith, Filip Radenovic, Francisco Guzmán, Frank Zhang, Gabriel Synnaeve, Gabrielle Lee, Georgia Lewis Anderson, Govind Thattai, Graeme Nail, Gregoire Mialon, Guan Pang, Guillem Cucurell, Hailey Nguyen, Hannah Korevaar, Hu Xu, Hugo Touvron, Iliyan Zarov, Imanol Arrieta Ibarra, Isabel Kloumann, Ishan Misra, Ivan Evtimov, Jack Zhang, Jade Copet, Jaewon Lee, Jan Geffert, Jana Vranes, Jason Park, Jay Mahadeokar, Jeet Shah, Jelmer van der Linde, Jennifer Billock, Jenny Hong, Jenya Lee, Jeremy Fu, Jianfeng Chi, Jianyu

参考文献

Ronen Eldan and Yuanzhi Li. TinyStories: How small can language models be and still speak coherent English?, 2023. arXiv:2305.07759.

Aaron Gokaslan, Vanya Cohen, Ellie Pavlick, and Stefanie Tellex. OpenWebText corpus. <http://Skylion007.github.io/OpenWebTextCorpus>, 2019.

Rico Sennrich, Barry Haddow, and Alexandra Birch. Neural machine translation of rare words with subword units. In *Proc. of ACL*, 2016.

Changhan Wang, Kyunghyun Cho, and Jiatao Gu. Neural machine translation with byte-level subwords, 2019. arXiv:1909.03341.

Philip Gage. A new algorithm for data compression. *C Users Journal*, 12(2):23–38, February 1994. ISSN 0898-9788.

Alec Radford, Jeff Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. Language models are unsupervised multitask learners, 2019.

Alec Radford, Karthik Narasimhan, Tim Salimans, and Ilya Sutskever. Improving language understanding by generative pre-training, 2018.

Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Ł ukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *Proc. of NeurIPS*, 2017.

Toan Q. Nguyen and Julian Salazar. Transformers without tears: Improving the normalization of self-attention. In *Proc. of IWSWLT*, 2019.

Ruibin Xiong, Yunchang Yang, Di He, Kai Zheng, Shuxin Zheng, Chen Xing, Huishuai Zhang, Yanyan Lan, Liwei Wang, and Tie-Yan Liu. On layer normalization in the Transformer architecture. In *Proc. of ICML*, 2020.

Jimmy Lei Ba, Jamie Ryan Kiros, and Geoffrey E. Hinton. Layer normalization, 2016. arXiv:1607.06450.

Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, Aurelien Rodriguez, Armand Joulin, Edouard Grave, and Guillaume Lample. Llama: Open and efficient foundation language models, 2023. arXiv:2302.13971.

Biao Zhang and Rico Sennrich. Root mean square layer normalization. In *Proc. of NeurIPS*, 2019.

Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Alex Vaughan, Amy Yang, Angela Fan, Anirudh Goyal, Anthony Hartshorn, Aobo Yang, Archi Mitra, Archie Sravankumar, Artem Korenev, Arthur Hinsvark, Arun Rao, Aston Zhang, Aurelien Rodriguez, Austen Gregerson, Ava Spataru, Baptiste Roziere, Bethany Biron, Binh Tang, Bobbie Chern, Charlotte Caucheteux, Chaya Nayak, Chloe Bi, Chris Marra, Chris McConnell, Christian Keller, Christophe Touret, Chunyang Wu, Corinne Wong, Cristian Canton Ferrer, Cyrus Nikolaidis, Damien Allonsius, Daniel Song, Danielle Pintz, Danny Livshits, Danny Wyatt, David Esiobu, Dhruv Choudhary, Dhruv Mahajan, Diego Garcia-Olano, Diego Perino, Dieuwke Hupkes, Egor Lakomkin, Ehab AlBadawy, Elina Lobanova, Emily Dinan, Eric Michael Smith, Filip Radenovic, Francisco Guzmán, Frank Zhang, Gabriel Synnaeve, Gabrielle Lee, Georgia Lewis Anderson, Govind Thattai, Graeme Nail, Gregoire Mialon, Guan Pang, Guillem Cucurell, Hailey Nguyen, Hannah Korevaar, Hu Xu, Hugo Touvron, Iliyan Zarov, Imanol Arrieta Ibarra, Isabel Kloumann, Ishan Misra, Ivan Evtimov, Jack Zhang, Jade Copet, Jaewon Lee, Jan Geffert, Jana Vranes, Jason Park, Jay Mahadeokar, Jeet Shah, Jelmer van der Linde, Jennifer Billock, Jenny Hong, Jenya Lee, Jeremy Fu, Jianfeng Chi, Jianyu

Huang, Jiawen Liu, Jie Wang, Jiecao Yu, Joanna Bitton, Joe Spisak, Jongsoo Park, Joseph Rocca, Joshua Johnstun, Joshua Saxe, Junteng Jia, Kalyan Vasuden Alwala, Karthik Prasad, Kartikeya Upasani, Kate Plawiak, Ke Li, Kenneth Heafield, Kevin Stone, Khalid El-Arini, Krithika Iyer, Kshitiz Malik, Kuenley Chiu, Kunal Bhalla, Kushal Lakhotia, Lauren Rantala-Yearly, Laurens van der Maaten, Lawrence Chen, Liang Tan, Liz Jenkins, Louis Martin, Lovish Madaan, Lubo Malo, Lukas Blecher, Lukas Landzaat, Luke de Oliveira, Madeline Muzzi, Mahesh Pasupuleti, Mannat Singh, Manohar Paluri, Marcin Kardas, Maria Tsimpoukelli, Mathew Oldham, Mathieu Rita, Maya Pavlova, Melanie Kambadur, Mike Lewis, Min Si, Mitesh Kumar Singh, Mona Hassan, Naman Goyal, Narjes Torabi, Nikolay Bashlykov, Nikolay Bogoychev, Niladri Chatterji, Ning Zhang, Olivier Duchenne, Onur Çelebi, Patrick Alrassy, Pengchuan Zhang, Pengwei Li, Petar Vasic, Peter Weng, Prajwal Bhargava, Pratik Dubal, Praveen Krishnan, Punit Singh Koura, Puxin Xu, Qing He, Qingxiao Dong, Ragavan Srinivasan, Raj Ganapathy, Ramon Calderer, Ricardo Silveira Cabral, Robert Stojnic, Roberta Raileanu, Rohan Maheswari, Rohit Girdhar, Rohit Patel, Romain Sauvestre, Ronnie Polidoro, Roshan Sumbaly, Ross Taylor, Ruan Silva, Rui Hou, Rui Wang, Saghar Hosseini, Sahana Chennabasappa, Sanjay Singh, Sean Bell, Seohyun Sonia Kim, Sergey Edunov, Shaoliang Nie, Sharan Narang, Sharath Raparthi, Sheng Shen, Shengye Wan, Shruti Bhosale, Shun Zhang, Simon Vandenhende, Soumya Batra, Spencer Whitman, Sten Sootla, Stephane Collet, Suchin Gururangan, Sydney Borodinsky, Tamar Herman, Tara Fowler, Tarek Sheasha, Thomas Georgiou, Thomas Scialom, Tobias Speckbacher, Todor Mihaylov, Tong Xiao, Ujjwal Karn, Vedanuj Goswami, Vibhor Gupta, Vignesh Ramanathan, Viktor Kerkez, Vincent Conguet, Virginie Do, Vish Vogeti, Vitor Albiero, Vladan Petrovic, Weiwei Chu, Wenhan Xiong, Wenxin Fu, Whitney Meers, Xavier Martinet, Xiaodong Wang, Xiaofang Wang, Xiaoqing Ellen Tan, Xide Xia, Xinfeng Xie, Xuchao Jia, Xuwei Wang, Yaelle Goldschlag, Yashesh Gaur, Yasmine Babaei, Yi Wen, Yiwen Song, Yuchen Zhang, Yue Li, Yuning Mao, Zacharie Delpierre Coudert, Zheng Yan, Zhengxing Chen, Zoe Papakipos, Aaditya Singh, Aayushi Srivastava, Abha Jain, Adam Kelsey, Adam Shajnfeld, Adithya Gangidi, Adolfo Victoria, Ahuva Goldstand, Ajay Menon, Ajay Sharma, Alex Boesenberg, Alexei Baevski, Allie Feinstein, Amanda Kallet, Amit Sangani, Amos Teo, Anam Yunus, Andrei Lupu, Andres Alvarado, Andrew Caples, Andrew Gu, Andrew Ho, Andrew Poulton, Andrew Ryan, Ankit Ramchandani, Annie Dong, Annie Franco, Anuj Goyal, Aparajita Saraf, Arka-bandhu Chowdhury, Ashley Gabriel, Ashwin Bharambe, Assaf Eisenman, Azadeh Yazdan, Beau James, Ben Maurer, Benjamin Leonhardi, Bernie Huang, Beth Loyd, Beto De Paola, Bhargavi Paranjape, Bing Liu, Bo Wu, Boyu Ni, Braden Hancock, Bram Wasti, Brandon Spence, Brani Stojkovic, Brian Gamido, Britt Montalvo, Carl Parker, Carly Burton, Catalina Mejia, Ce Liu, Changhan Wang, Changkyu Kim, Chao Zhou, Chester Hu, Ching-Hsiang Chu, Chris Cai, Chris Tindal, Christoph Feichtenhofer, Cynthia Gao, Damon Civan, Dana Beaty, Daniel Kreymer, Daniel Li, David Adkins, David Xu, Davide Testuggine, Delia David, Devi Parikh, Diana Liskovich, Didem Foss, Dingkan Wang, Duc Le, Dustin Holland, Edward Dowling, Eissa Jamil, Elaine Montgomery, Eleonora Presani, Emily Hahn, Emily Wood, Eric-Tuan Le, Erik Brinkman, Esteban Arcaute, Evan Dunbar, Evan Smothers, Fei Sun, Felix Kreuk, Feng Tian, Filippas Kokkinos, Firat Ozgenel, Francesco Caggioni, Frank Kanayet, Frank Seide, Gabriela Medina Florez, Gabriella Schwarz, Gada Badeer, Georgia Swee, Gil Halpern, Grant Herman, Grigory Sizov, Guangyi Zhang, Guna Lakshminarayanan, Hakan Inan, Hamid Shojanazeri, Han Zou, Hannah Wang, Hanwen Zha, Haroun Habeeb, Harrison Rudolph, Helen Suk, Henry Aspegren, Hunter Goldman, Hongyuan Zhan, Ibrahim Damlaj, Igor Molybog, Igor Tufanov, Ilias Leontiadis, Irina-Elena Veliche, Itai Gat, Jake Weissman, James Geboski, James Kohli, Janice Lam, Japhet Asher, Jean-Baptiste Gaya, Jeff Marcus, Jeff Tang, Jennifer Chan, Jenny Zhen, Jeremy Reizenstein, Jeremy Teboul, Jessica Zhong, Jian Jin, Jingyi Yang, Joe Cummings, Jon Carvill, Jon Shepard, Jonathan McPhie, Jonathan Torres, Josh Ginsburg, Junjie Wang, Kai Wu, Kam Hou U, Karan Saxena, Kartikay Khandelwal, Katayoun Zand, Kathy Matosich, Kaushik Veeraraghavan, Kelly Michelena, Keqian Li, Kiran Jagadeesh, Kun Huang, Kunal Chawla, Kyle Huang, Lailin Chen, Lakshya Garg, Lavender A, Leandro Silva, Lee Bell, Lei Zhang, Liangpeng Guo, Licheng Yu, Liron Moshkovich, Luca Wehrstedt, Madian Khabsa, Manav Avalani, Manish Bhatt, Martynas Mankus, Matan Hasson, Matthew Lennie, Matthias Reso, Maxim Groshev, Maxim Naumov, Maya Lathi, Meghan Keneally, Miao Liu, Michael L. Seltzer, Michal Valko, Michelle Restrepo, Mihir Patel, Mik Vyatskov, Mikayel Samvelyan, Mike Clark, Mike Macey, Mike Wang, Miquel Jubert Hermoso, Mo Metanat, Moham-

Huang, Jiawen Liu, Jie Wang, Jiecao Yu, Joanna Bitton, Joe Spisak, Jongsoo Park, Joseph Rocca, Joshua Johnstun, Joshua Saxe, Junteng Jia, Kalyan Vasuden Alwala, Karthik Prasad, Kartikeya Upasani, Kate Plawiak, Ke Li, Kenneth Heafield, Kevin Stone, Khalid El-Arini, Krithika Iyer, Kshitiz Malik, Kuenley Chiu, Kunal Bhalla, Kushal Lakhotia, Lauren Rantala-Yearly, Laurens van der Maaten, Lawrence Chen, Liang Tan, Liz Jenkins, Louis Martin, Lovish Madaan, Lubo Malo, Lukas Blecher, Lukas Landzaat, Luke de Oliveira, Madeline Muzzi, Mahesh Pasupuleti, Mannat Singh, Manohar Paluri, Marcin Kardas, Maria Tsimpoukelli, Mathew Oldham, Mathieu Rita, Maya Pavlova, Melanie Kambadur, Mike Lewis, Min Si, Mitesh Kumar Singh, Mona Hassan, Naman Goyal, Narjes Torabi, NikolayBashlykov, NikolayBogoychev, Niladri Chatterji, Ning Zhang, Olivier Duchenne, Onur Çelebi, Patrick Alrassy, Pengchuan Zhang, Pengwei Li, Petar Vasic, Peter Weng, Prajwal Bhargava, Pratik Dubal, Praveen Krishnan, Punit Singh Koura, Puxin Xu, Qing He, Qingxiao Dong, Ragavan Srinivasan, Raj Ganapathy, Ramon Calderer, Ricardo Silveira Cabral, Robert Stojnic, Roberta Raileanu, Rohan Maheswari, Rohit Girdhar, Rohit Patel, Romain Sauvestre, Ronnie Polidoro, Roshan Sumbaly, Ross Taylor, Ruan Silva, Rui Hou, Rui Wang, Saghar Hosseini, Sahana Chennabasappa, Sanjay Singh, Sean Bell, Seohyun Sonia Kim, Sergey Edunov, Shaoliang Nie, Sharan Narang, Sharath Raparthi, Sheng Shen, Shengye Wan, Shruti Bhosale, Shun Zhang, Simon Vandenhende, Soumya Batra, Spencer Whitman, Sten Sootla, Stephane Collet, Suchin Gururangan, Sydney Borodinsky, Tamar Herman, Tara Fowler, Tarek Sheasha, Thomas Georgiou, Thomas Scialom, Tobias Speckbacher, Todor Mihaylov, Tong Xiao, Ujjwal Karn, Vedanuj Goswami, Vibhor Gupta, Vignesh Ramanathan, Viktor Kerkez, Vincent Conguet, Virginie Do, Vish Vogeti, Vitor Albiero, Vladan Petrovic, Weiwei Chu, Wenhan Xiong, Wenxin Fu, Whitney Meers, Xavier Martinet, Xiaodong Wang, Xiaofang Wang, Xiaoqing Ellen Tan, Xide Xia, Xinfeng Xie, Xuchao Jia, Xuwei Wang, Yaelle Goldschlag, Yashesh Gaur, Yasmine Babaei, Yi Wen, Yiwen Song, Yuchen Zhang, Yue Li, Yuning Mao, Zacharie Delpierre Coudert, Zheng Yan, Zhengxing Chen, Zoe Papakipos, Aaditya Singh, Aayushi Srivastava, Abha Jain,

亚当·凯尔斯, 亚当·沙申菲尔德, 阿迪提亚·冈迪, 阿道夫·维多利亚, 阿胡瓦·戈尔斯坦, 阿杰伊·门农, 阿杰伊·夏尔马, 亚历克斯·贝森伯格, 阿列克谢·巴夫斯基, 艾莉·费因斯坦, 阿曼达·卡莱卡特, 阿米特·桑加尼, 阿莫斯·李, Anam Yunus, Andrei Lupu, Andres Alvarado, Andrew Caples, Andrew Gu, Andrew Ho, Andrew Poulton, Andrew Ryan, Ankit Ramchandani, Annie Dong, Annie Franco, Anuj Goyal, Aparajita Saraf, Arka-bandhu Chowdhury, Ashley Gabriel, Ashwin Bharambe, Assaf Eisenman, Azadeh Yazdan, Beau James, Ben Maurer, Benjamin Leonhardi, Bernie Huang, Beth Loyd, Beto De Paola, Bhargavi Paranjape, Bing Liu, Bo Wu, Boyu Ni, Braden Hancock, Bram Wasti, Brandon Spence, Brani Stojkovic, Brian Gamido, Britt Montalvo, Carl Parker, Carly Burton, Catalina Mejia, Ce Liu, Changhan Wang, Changkyu Kim, Chao Zhou, Chester Hu, Ching-Hsiang Chu, Chris Cai, Chris Tindal, Christoph Feichtenhofer, Cynthia Gao, Damon Civan, Dana Beaty, Daniel Kreymer, Daniel Li, David Adkins, David Xu, Davide Testuggine, Delia David, Devi Parikh, Diana Liskovich, Didem Foss, Dingkan Wang, Duc Le, Dustin Holland, Edward Dowling, Eissa Jamil, Elaine Montgomery, Eleonora Presani, Emily Hahn, Emily Wood, Eric-Tuan Le, Erik Brinkman, Esteban Arcaute, Evan Dunbar, Evan Smothers, Fei Sun, Felix Kreuk, Feng Tian, Filippas Kokkinos, Firat Ozgenel, Francesco Caggioni, Frank Kanayet, Frank Seide, Gabriela Medina Flores, Gabriella Schwarz, Gada Badeer, Georgia Swee, Gil Halpern, Grant Herman, Grigory Sizov, Guangyi Zhang, Guna Lakshminarayanan, Hakan Inan, Hamid Shojanazeri, Han Zou, Hannah Wang, Hanwen Zha, Haroun Habeeb, Harrison Rudolph, Helen Suk, Henry Aspegren, Hunter Goldman, Hongyuan Zhan, Ibrahim Damlaj, Igor Molybog, Igor Tufanov, Ilias Leontiadis, Irina-Elena Veliche, Itai Gat, Jake Weissman, JamesGeboski, James Kohli, Janice Lam, Japhet Asher, Jean-Baptiste Gaya, Jeff Marcus, Jeff Tang, Jennifer Chan, Jenny Zhen, Jeremy Reizenstein, Jeremy Teboul, Jessica Zhong, Jian Jin, Jingyi Yang, Joe Cummings, Jon Carvill, Jon Shepard, Jonathan McPhie, Jonathan Torres, Josh Ginsburg, Junjie Wang, Kai Wu, Kam Hou U, Karan Saxena, Kartikay Khandelwal, Katayoun Zand, Kathy Matosich, Kaushik Veeraraghavan, Kelly Michelena, Keqian Li, Kiran Jagadeesh, Kun Huang, Kunal Chawla, Kyle Huang, Lailin Chen, Lakshya Garg, Lavender A, Leandro Silva, Lee Bell, Lei Zhang, Liangpeng Guo, Licheng Yu, Liron Moshkovich, Luca Wehrstedt, Madian Khabsa, Manav Avalani, Manish Bhatt, Martynas Mankus, Matan Hasson, Matthew Lennie, Matthias Reso, Maxim Groshev, Maxim Naumov, Maya Lathi, Meghan Keneally, Miao Liu, Michael L. Seltzer, Michal Valko, Michelle Restrepo, Mihir Patel, Mik Vyatskov, Mikayel Samvelyan, Mike Clark, Mike Macey, Mike Wang, Miquel Jubert Hermoso, Mo Metanat, Moham

mad Rastegari, Munish Bansal, Nandhini Santhanam, Natascha Parks, Natasha White, Navyata Bawa, Nayan Singhal, Nick Egebo, Nicolas Usunier, Nikhil Mehta, Nikolay Pavlovich Laptev, Ning Dong, Norman Cheng, Oleg Chernoguz, Olivia Hart, Omkar Salpekar, Ozlem Kalinli, Parkin Kent, Parth Parekh, Paul Saab, Pavan Balaji, Pedro Rittner, Philip Bontrager, Pierre Roux, Piotr Dollar, Polina Zvyagina, Prashant Ratanchandani, Pritish Yuvraj, Qian Liang, Rachad Alao, Rachel Rodriguez, Rafi Ayub, Raghotham Murthy, Raghu Nayani, Rahul Mitra, Rangaprabhu Parthasarathy, Raymond Li, Rebekkah Hogan, Robin Battey, Rocky Wang, Russ Howes, Ruty Rinott, Sachin Mehta, Sachin Siby, Sai Jayesh Bondu, Samyak Datta, Sara Chugh, Sara Hunt, Sargun Dhillon, Sasha Sidorov, Satadru Pan, Saurabh Mahajan, Saurabh Verma, Seiji Yamamoto, Sharadh Ramaswamy, Shaun Lindsay, Shaun Lindsay, Sheng Feng, Shenghao Lin, Shengxin Cindy Zha, Shishir Patil, Shiva Shankar, Shuqiang Zhang, Shuqiang Zhang, Sinong Wang, Sneha Agarwal, Soji Sajuyigbe, Soumith Chintala, Stephanie Max, Stephen Chen, Steve Kehoe, Steve Satterfield, Sudarshan Govindaprasad, Sumit Gupta, Summer Deng, Sungmin Cho, Sunny Virk, Suraj Subramanian, Sy Choudhury, Sydney Goldman, Tal Remez, Tamar Glaser, Tamara Best, Thilo Koehler, Thomas Robinson, Tianhe Li, Tianjun Zhang, Tim Matthews, Timothy Chou, Tzook Shaked, Varun Vontimitta, Victoria Ajayi, Victoria Montanez, Vijai Mohan, Vinay Satish Kumar, Vishal Mangla, Vlad Ionescu, Vlad Poenaru, Vlad Tiberiu Mihailescu, Vladimir Ivanov, Wei Li, Wenchen Wang, Wenwen Jiang, Wes Bouaziz, Will Constable, Xiaocheng Tang, Xiaojuan Wu, Xiaolan Wang, Xilun Wu, Xinbo Gao, Yaniv Kleinman, Yanjun Chen, Ye Hu, Ye Jia, Ye Qi, Yenda Li, Yilin Zhang, Ying Zhang, Yossi Adi, Youngjin Nam, Yu, Wang, Yu Zhao, Yuchen Hao, Yundi Qian, Yunlu Li, Yuzi He, Zach Rait, Zachary DeVito, Zef Rosnbrick, Zhaoduo Wen, Zhenyu Yang, Zhiwei Zhao, and Zhiyu Ma. The llama 3 herd of models, 2024. URL <https://arxiv.org/abs/2407.21783>.

An Yang, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chengyuan Li, Dayiheng Liu, Fei Huang, Haoran Wei, Huan Lin, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin Yang, Jiayi Yang, Jingren Zhou, Junyang Lin, Kai Dang, Keming Lu, Keqin Bao, Kexin Yang, Le Yu, Mei Li, Mingfeng Xue, Pei Zhang, Qin Zhu, Rui Men, Runji Lin, Tianhao Li, Tingyu Xia, Xingzhang Ren, Xuancheng Ren, Yang Fan, Yang Su, Yichang Zhang, Yu Wan, Yuqiong Liu, Zeyu Cui, Zhenru Zhang, and Zihan Qiu. Qwen2.5 technical report. *arXiv preprint arXiv:2412.15115*, 2024.

Aakanksha Chowdhery, Sharan Narang, Jacob Devlin, Maarten Bosma, Gaurav Mishra, Adam Roberts, Paul Barham, Hyung Won Chung, Charles Sutton, Sebastian Gehrmann, Parker Schuh, Kensen Shi, Sasha Tsvyashchenko, Joshua Maynez, Abhishek Rao, Parker Barnes, Yi Tay, Noam Shazeer, Vinodkumar Prabhakaran, Emily Reif, Nan Du, Ben Hutchinson, Reiner Pope, James Bradbury, Jacob Austin, Michael Isard, Guy Gur-Ari, Pengcheng Yin, Toju Duke, Anselm Levskaya, Sanjay Ghemawat, Sunipa Dev, Henryk Michalewski, Xavier Garcia, Vedant Misra, Kevin Robinson, Liam Fedus, Denny Zhou, Daphne Ippolito, David Luan, Hyeontaek Lim, Barret Zoph, Alexander Spiridonov, Ryan Sepassi, David Dohan, Shivani Agrawal, Mark Omernick, Andrew M. Dai, Thanumalayan Sankaranarayanan Pillai, Marie Pellat, Aitor Lewkowycz, Erica Moreira, Rewon Child, Aleksandr Polozov, Katherine Lee, Zongwei Zhou, Xuezhi Wang, Brennan Saeta, Mark Diaz, Orhan Firat, Michele Catasta, Jason Wei, Kathy Meier-Hellstern, Douglas Eck, Jeff Dean, Slav Petrov, and Noah Fiedel. PaLM: Scaling language modeling with pathways, 2022. arXiv:2204.02311.

Dan Hendrycks and Kevin Gimpel. Bridging nonlinearities and stochastic regularizers with gaussian error linear units, 2016. arXiv:1606.08415.

Stefan Elfving, Eiji Uchibe, and Kenji Doya. Sigmoid-weighted linear units for neural network function approximation in reinforcement learning, 2017. URL <https://arxiv.org/abs/1702.03118>.

Yann N. Dauphin, Angela Fan, Michael Auli, and David Grangier. Language modeling with gated convolutional networks, 2017. URL <https://arxiv.org/abs/1612.08083>.

Noam Shazeer. GLU variants improve transformer, 2020. arXiv:2002.05202.

Jianlin Su, Yu Lu, Shengfeng Pan, Bo Wen, and Yunfeng Liu. Roformer: Enhanced transformer with rotary position embedding, 2021.

mad Rastegari, Munish Bansal, Nandhini Santhanam, Natascha Parks, Natasha White, Navyata Bawa, Nayan Singhal, Nick Egebo, Nicolas Usunier, Nikhil Mehta, Nikolay Pavlovich Laptev, Ning Dong, Norman Cheng, Oleg Chernoguz, Olivia Hart, Omkar Salpekar, Ozlem Kalinli, Parkin Kent, Parth Parekh, Paul Saab, Pavan Balaji, Pedro Rittner, Philip Bontrager, Pierre Roux, Piotr Dollar, Polina Zvyagina, Prashant Ratanchandani, Pritish Yuvraj, Qian Liang, Rachad Alao, Rachel Rodriguez, Rafi Ayub, Raghotham Murthy, Raghu Nayani, Rahul Mitra, Rangaprabhu Parthasarathy, Raymond Li, Rebekkah Hogan, Robin Battey, Rocky Wang, Russ Howes, Ruty Rinott, Sachin Mehta, Sachin Siby, Sai Jayesh Bondu, Samyak Datta, Sara Chugh, Sara Hunt, Sargun Dhillon, Sasha Sidorov, Satadru Pan, Saurabh Mahajan, Saurabh Verma, Seiji Yamamoto, Sharadh Ramaswamy, Shaun Lindsay, Shaun Lindsay, Sheng Feng, Shenghao Lin, Shengxin Cindy Zha, Shishir Patil, Shiva Shankar, Shuqiang Zhang, Shuqiang Zhang, Sinong Wang, Sneha Agarwal, Soji Sajuyigbe, Soumith Chintala, Stephanie Max, Stephen Chen, Steve Kehoe, Steve Satterfield, Sudarshan Govindaprasad, Sumit Gupta, Summer Deng, Sungmin Cho, Sunny Virk, Suraj Subramanian, Sy Choudhury, Sydney Goldman, Tal Remez, Tamar Glaser, Tamara Best, Thilo Koehler, Thomas Robinson, Tianhe Li, Tianjun Zhang, Tim Matthews, Timothy Chou, Tzook Shaked, Varun Vontimitta, Victoria Ajayi, Victoria Montanez, Vijai Mohan, Vinay Satish Kumar, Vishal Mangla, Vlad Ionescu, Vlad Poenaru, Vlad Tiberiu Mihailescu, Vladimir Ivanov, Wei Li, Wenchen Wang, Wenwen Jiang, Wes Bouaziz, Will Constable, Xiaocheng Tang, Xiaojuan Wu, Xiaolan Wang, Xilun Wu, Xinbo Gao, Yaniv Kleinman, Yanjun Chen, Ye Hu, Ye Jia, Ye Qi, Yenda Li, Yilin Zhang, Ying Zhang, Yossi Adi, Youngjin Nam, Yu, Wang, Yu Zhao, Yuchen Hao, Yundi Qian, Yunlu Li, Yuzi He, Zach Rait, Zachary DeVito, Zef Rosnbrick, Zhaoduo Wen, Zhenyu Yang, Zhiwei Zhao, and Zhiyu Ma. The llama 3 herd of models, 2024. URL <https://arxiv.org/abs/2407.21783>.

An Yang, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chengyuan Li, Dayiheng Liu, Fei Huang, Haoran Wei, Huan Lin, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin Yang, Jiayi Yang, Jingren Zhou, Junyang Lin, Kai Dang, Keming Lu, Keqin Bao, Kexin Yang, Le Yu, Mei Li, Mingfeng Xue, Pei Zhang, Qin Zhu, Rui Men, Runji Lin, Tianhao Li, Tingyu Xia, Xingzhang Ren, Xuancheng Ren, Yang Fan, Yang Su, Yichang Zhang, Yu Wan, Yuqiong Liu, Zeyu Cui, Zhenru Zhang, and Zihan Qiu. Qwen2.5 technical report. *arXiv preprint arXiv:2412.15115*, 2024.

Aakanksha Chowdhery, Sharan Narang, Jacob Devlin, Maarten Bosma, Gaurav Mishra, Adam Roberts, Paul Barham, Hyung Won Chung, Charles Sutton, Sebastian Gehrmann, Parker Schuh, Kensen Shi, Sasha Tsvyashchenko, Joshua Maynez, Abhishek Rao, Parker Barnes, Yi Tay, Noam Shazeer, Vinodkumar Prabhakaran, Emily Reif, Nan Du, Ben Hutchinson, Reiner Pope, James Bradbury, Jacob Austin, Michael Isard, Guy Gur-Ari, Pengcheng Yin, Toju Duke, Anselm Levskaya, Sanjay Ghemawat, Sunipa Dev, Henryk Michalewski, Xavier Garcia, Vedant Misra, Kevin Robinson, Liam Fedus, Denny Zhou, Daphne Ippolito, David Luan, Hyeontaek Lim, Barret Zoph, Alexander Spiridonov, Ryan Sepassi, David Dohan, Shivani Agrawal, Mark Omernick, Andrew M. Dai, Thanumalayan Sankaranarayanan Pillai, Marie Pellat, Aitor Lewkowycz, Erica Moreira, Rewon Child, Aleksandr Polozov, Katherine Lee, Zongwei Zhou, Xuezhi Wang, Brennan Saeta, Mark Diaz, Orhan Firat, Michele Catasta, Jason Wei, Kathy Meier-Hellstern, Douglas Eck, Jeff Dean, Slav Petrov, and Noah Fiedel. PaLM: Scaling language modeling with pathways, 2022. arXiv:2204.02311.

Dan Hendrycks and Kevin Gimpel. Bridging nonlinearities and stochastic regularizers with gaussian error linear units, 2016. arXiv:1606.08415.

Stefan Elfving, Eiji Uchibe, and Kenji Doya. Sigmoid-weighted linear units for neural network function approximation in reinforcement learning, 2017. URL <https://arxiv.org/abs/1702.03118>.

Yann N. Dauphin, Angela Fan, Michael Auli, and David Grangier. Language modeling with gated convolutional networks, 2017. URL <https://arxiv.org/abs/1612.08083>.

Noam Shazeer. GLU variants improve transformer, 2020. arXiv:2002.05202.

Jianlin Su, Yu Lu, Shengfeng Pan, Bo Wen, and Yunfeng Liu. Roformer: Enhanced transformer with rotary position embedding, 2021.

Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In *Proc. of ICLR*, 2015.

Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. In *Proc. of ICLR*, 2019.

Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Nee-lakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. Language models are few-shot learners. In *Proc. of NeurIPS*, 2020.

Jared Kaplan, Sam McCandlish, Tom Henighan, Tom B. Brown, Benjamin Chess, Rewon Child, Scott Gray, Alec Radford, Jeffrey Wu, and Dario Amodei. Scaling laws for neural language models, 2020. arXiv:2001.08361.

Jordan Hoffmann, Sebastian Borgeaud, Arthur Mensch, Elena Buchatskaya, Trevor Cai, Eliza Rutherford, Diego de Las Casas, Lisa Anne Hendricks, Johannes Welbl, Aidan Clark, Tom Hennigan, Eric Noland, Katie Millican, George van den Driessche, Bogdan Damoc, Aurelia Guy, Simon Osindero, Karen Simonyan, Erich Elsen, Jack W. Rae, Oriol Vinyals, and Laurent Sifre. Training compute-optimal large language models, 2022. arXiv:2203.15556.

Ari Holtzman, Jan Buys, Li Du, Maxwell Forbes, and Yejin Choi. The curious case of neural text degeneration. In *Proc. of ICLR*, 2020.

Yao-Hung Hubert Tsai, Shaojie Bai, Makoto Yamada, Louis-Philippe Morency, and Ruslan Salakhutdinov. Transformer dissection: An unified understanding for transformer’s attention via the lens of kernel. In Kentaro Inui, Jing Jiang, Vincent Ng, and Xiaojun Wan, editors, *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 4344–4353, Hong Kong, China, November 2019. Association for Computational Linguistics. doi: 10.18653/v1/D19-1443. URL <https://aclanthology.org/D19-1443/>.

Amirhossein Kazemnejad, Inkit Padhi, Karthikeyan Natesan, Payel Das, and Siva Reddy. The impact of positional encoding on length generalization in transformers. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023. URL <https://openreview.net/forum?id=Drrl2gcjz1>.

Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In *Proc. of ICLR*, 2015.

Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. In *Proc. of ICLR*, 2019.

Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Nee-lakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. Language models are few-shot learners. In *Proc. of NeurIPS*, 2020.

Jared Kaplan, Sam McCandlish, Tom Henighan, Tom B. Brown, Benjamin Chess, Rewon Child, Scott Gray, Alec Radford, Jeffrey Wu, and Dario Amodei. Scaling laws for neural language models, 2020. arXiv:2001.08361.

Jordan Hoffmann, Sebastian Borgeaud, Arthur Mensch, Elena Buchatskaya, Trevor Cai, Eliza Rutherford, Diego de Las Casas, Lisa Anne Hendricks, Johannes Welbl, Aidan Clark, Tom Hennigan, Eric Noland, Katie Millican, George van den Driessche, Bogdan Damoc, Aurelia Guy, Simon Osindero, Karen Simonyan, Erich Elsen, Jack W. Rae, Oriol Vinyals, and Laurent Sifre. Training compute-optimal large language models, 2022. arXiv:2203.15556.

Ari Holtzman, Jan Buys, Li Du, Maxwell Forbes, and Yejin Choi. The curious case of neural text degeneration. In *Proc. of ICLR*, 2020.

Yao-Hung Hubert Tsai, Shaojie Bai, Makoto Yamada, Louis-Philippe Morency, and Ruslan Salakhutdinov. Transformer dissection: An unified understanding for transformer’s attention via the lens of kernel. In Kentaro Inui, Jing Jiang, Vincent Ng, and Xiaojun Wan, editors, *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 4344–4353, Hong Kong, China, November 2019. Association for Computational Linguistics. doi: 10.18653/v1/D19-1443. URL <https://aclanthology.org/D19-1443/>.

Amirhossein Kazemnejad, Inkit Padhi, Karthikeyan Natesan, Payel Das, and Siva Reddy. The impact of positional encoding on length generalization in transformers. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023. URL <https://openreview.net/forum?id=Drrl2gcjz1>.