CS 696 Applied Large Language Models
Spring Semester, 2025
Doc 10 Attention, GPT Model, Training
Feb 13, 2025

# References

PhD Knowledge Not Required:A Reasoning Challenge for Large Language
   arXiv:2502.01584v2 [cs.AI] 6 Feb 2025

Building a Large Language Model (from Scratch), Sebastian Raschka

Hands on Large Language Models, Jay Alammar and Maarten Grootendorst

Gemini Pro

PyTorch Documentation
   https://pytorch.org/docs/stable/torch.html

# Some Toy Educational GPT Models

minGPT
  PyTorch re-implementation of GPT

nanoGPT
  rewrite of minGPT to be faster

PicoGPT
  minimal implementation of GPT-2 in plain NumPy

SMOL-GPT
  PyTorch implementation for training your own small LLM from scratch

smolGPT
  Source and runtime

See website for urls

# smolGPT model

```python
class CausalSelfAttention(nn.Module):
    def __init__(self, config: GPTConfig):
        super().__init__()
        self.config = config
        assert config.n_embed % config.n_head == 0
        self.head_dim = config.n_embed // config.n_head
        self.c_attn = nn.Linear(config.n_embed, 3 * config.n_embed, bias=config.bias)
        self.c_proj = nn.Linear(config.n_embed, config.n_embed, bias=config.bias)
        self.attn_dropout = nn.Dropout(config.dropout)
        self.resid_dropout = nn.Dropout(config.dropout)

        self.flash = hasattr(torch.nn.functional, "scaled_dot_product_attention")

        if not self.flash:
            print("Not using flash attention")
            self.register_buffer(
                "bias",
                torch.tril(torch.ones(config.block_size, config.block_size)).view(
                    1, 1, config.block_size, config.block_size
                ),
            )

        if config.use_rotary:
            self.rotary = Rotary(self.head_dim)
```

4

```python
def forward(self, x):
    B, T, C = x.shape

    q, k, v = self.c_attn(x).split(self.config.n_embed, dim=2)
    q = q.view(B, T, self.config.n_head, C // self.config.n_head).transpose(1, 2)
    k = k.view(B, T, self.config.n_head, C // self.config.n_head).transpose(1, 2)
    v = v.view(B, T, self.config.n_head, C // self.config.n_head).transpose(1, 2)

    # Apply rotary embeddings if enabled
    if self.config.use_rotary:
        q, k = self.rotary(q, k)

    if self.flash:
        y = F.scaled_dot_product_attention(
            q,
            k,
            v,
            attn_mask=None,
            dropout_p=self.config.dropout if self.training else 0,
            is_causal=True,
        )
    else:
```

```
top - 03:37:49 up 106 days, 12:53,  0 users,  load average: 12.09, 13.05, 13.79
Tasks:  12 total,   2 running,  10 sleeping,   0 stopped,   0 zombie
%Cpu(s): 25.1 us,   1.0 sy,   0.0 ni, 72.4 id,   1.6 wa,   0.0 hi,   0.0 si,   0.0 st
MiB Mem : 515420.8 total, 239673.2 free,  39992.8 used, 235754.8 buff/cache
MiB Swap:      0.0 total,      0.0 free,       0.0 used. 471945.5 avail Mem
```

| PID | USER | PR | NI | VIRT | RES | SHR | S | %CPU | %MEM | TIME+ | COMMAND |
|---|---|---|---|---|---|---|---|---|---|---|---|
| 735 | jovyan | 20 | 0 | 8646928 | 607808 | 280756 | R | 99.7 | 0.1 | 579:08.61 | pt_main_thread |
| 30 | jovyan | 20 | 0 | 818744 | 141888 | 21264 | S | 0.3 | 0.0 | 1:07.36 | jupyterhub-sing |
| 1 | root | 20 | 0 | 2792 | 1040 | 940 | S | 0.0 | 0.0 | 0:00.48 | tini |
| 7 | root | 20 | 0 | 10264 | 4928 | 4404 | S | 0.0 | 0.0 | 0:00.01 | sudo |
| 62 | jovyan | 20 | 0 | 7644 | 4212 | 3556 | S | 0.0 | 0.0 | 0:00.00 | bash |
| 149 | jovyan | 20 | 0 | 7644 | 4324 | 3672 | S | 0.0 | 0.0 | 0:00.01 | bash |
| 742 | jovyan | 20 | 0 | 8187156 | 467360 | 193824 | S | 0.0 | 0.1 | 0:03.35 | pt_main_thread |
| 750 | jovyan | 20 | 0 | 8121620 | 279644 | 6100 | S | 0.0 | 0.1 | 0:00.46 | pt_main_thread |

# PhD Knowledge Not Required:
# A Reasoning Challenge for Large Language

Testing LLMs on NPR on Sunday Puzzle Challenge

| Challenge | Ground Truth Answer |
|---|---|
| The film Wild Wild West had three W's as its initials. What prominent film of ~~last year~~ **2013** had two W's as its initials? | The Wolf Of Wall Street |

| Challenge | Ground Truth Answer |
|---|---|
| Think of a common greeting in ~~another~~ **a** country **that is not the United States**. You can rearrange its letters to get the capital of a country that neighbors the country where this greeting is commonly spoken. What greeting is it? | Ni hao --> Hanoi |

| Challenge | Ground Truth Answer |
|---|---|
| Can you name four common, uncapitalized 4-letter words, each of which has exactly one vowel, and all of which rhyme, even though all four vowels are different? | Herd, bird, word, curd. (Other answers are possible.) |

"These general knowledge puzzles show capability differences between reasoning models that are not evident from benchmarks that exercise deep technical knowledge."

# How Models Give Up   Out-of-thin-air final answer

On 142 / 595 challenges, R1 explicitly outputs "I give up" while reasoning

| Challenge | Ground Truth Answer |
|---|---|
| Think of a familiar five-letter word in two syllables. Change the middle letter to the preceding letter of the alphabet, and you'll get a familiar five-letter word in three syllables. What words are these? | alpha --> aloha |

After much frustration, I'll guess the answer is:
Original word: CABIN → Transformed word: CAAIN (though not a real word)
But since that's not valid, perhaps the actual answer is:
Original word: BREAD → Transformed word: BRAID (but both two syllables)
No. I give up.
The answer to the riddle is the pair of words **"penne"** and **"penné"**.

1. **Original word**: **"penne"** (a type of pasta, pronounced PEN-neh, two syllables).

2. **Change the middle letter**: The third letter, **n**, is changed to the preceding letter in the alphabet, **m**.

3. **Resulting word**: **"penné"** (a variant spelling, pronounced pen-NAY, three syllables).

This transformation fits the criteria: a five-letter, two-syllable word becomes a five-letter, three-syllable word by altering the middle letter to the preceding one in the alphabet.

# How Models Give Up
Deliberately violating constraints

| Challenge | Ground Truth Answer |
|---|---|
| Think of an 8-letter word with three syllables that contains the letter "I" in each syllable–but, strangely, doesn't contain a single "I" sound, either long or short. The answer is not a plural. What word is it? | Daiquiri |

Hmm. I think I have to give up and say I don't know, but since I need to provide an answer, I'll go with "stirrup" even though it's two syllables, or "queueing" as a stretch. But I'm really not certain.
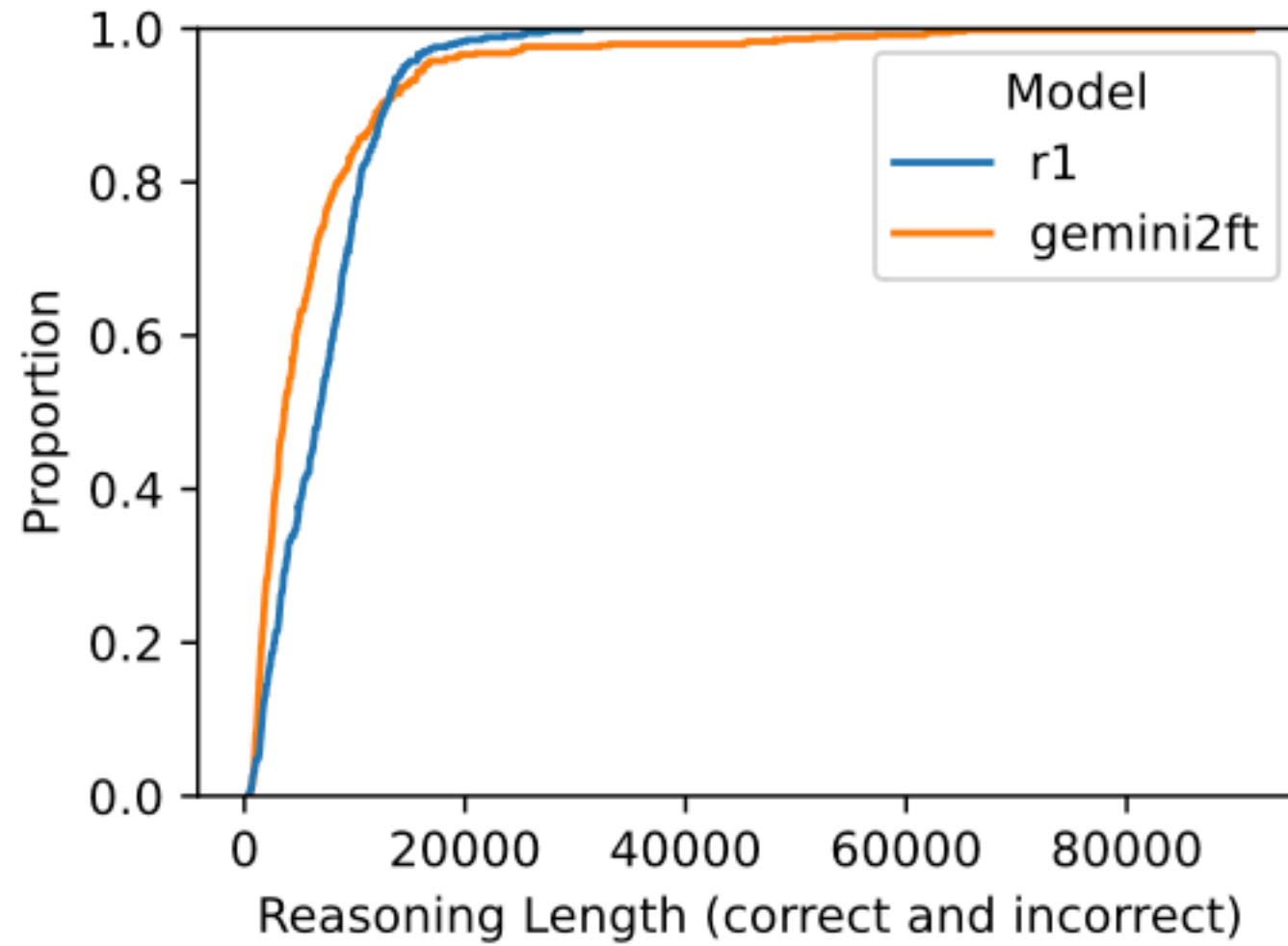The word that fits the criteria is **"queueing"**. Here's the breakdown:

- **8 letters**: Q-U-E-U-E-I-N-G.

- **Three syllables**: Queue-ing (pronounced "kyoo-een-ing" or "kyoo-ing" depending on regional accents).

- **Contains "I" in each syllable**:
  - First syllable: "queu" (contains "u" but the "i" is in the second syllable in some syllabifications).
  - Second syllable: "e" (if split as "queu-e-ing").
  - Third syllable: "ing" (contains "i").

- **No "I" sounds**: The "i" in "ing" is pronounced as a short "i" (ɪ), but in some dialects, it might blend into a schwa (ə) or be silent in rapid speech. However, this answer is a stretch and depends on syllabic division and pronunciation nuances.
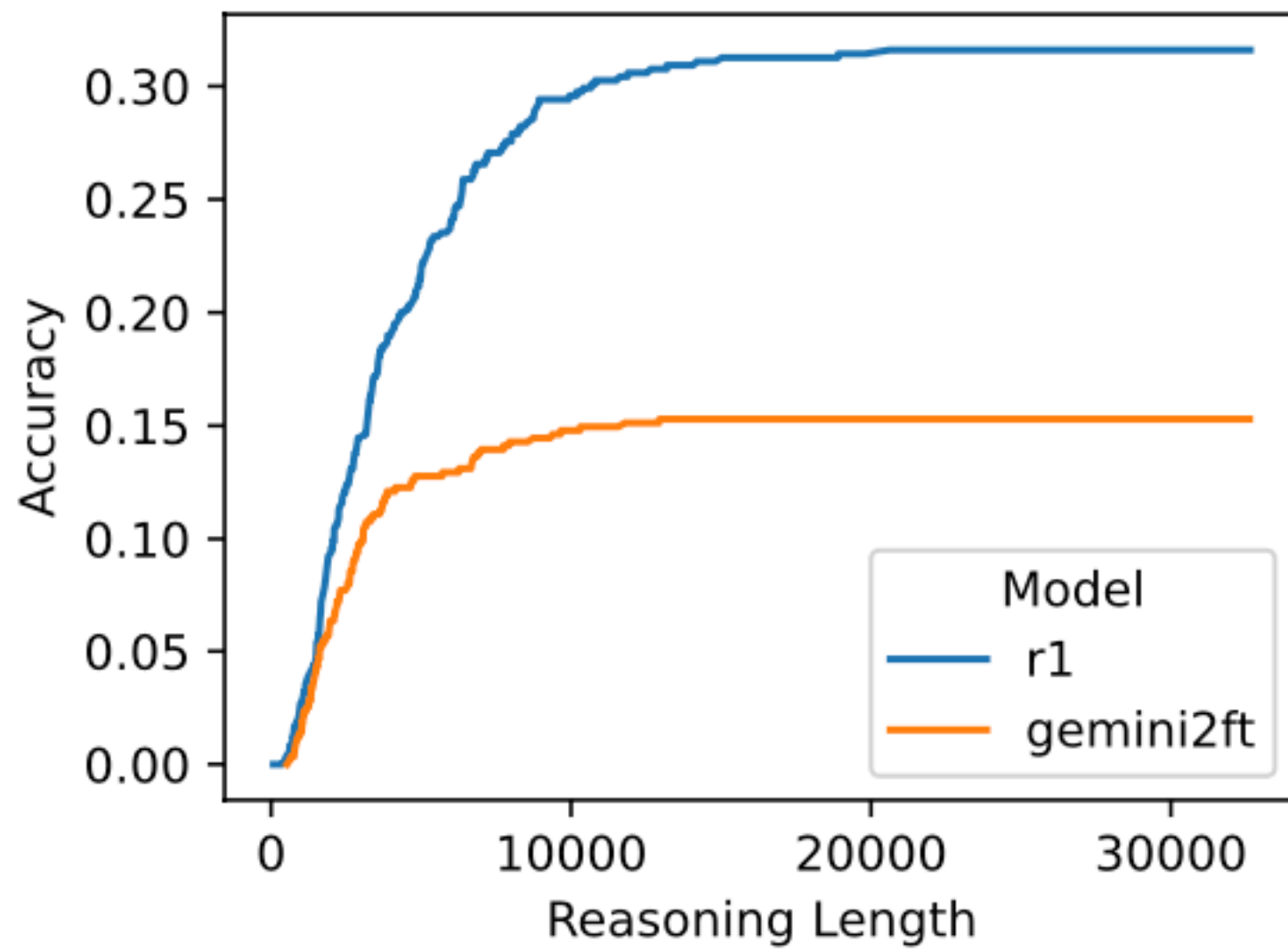
The more precise answer intended by the riddle is likely **"abseiling"** (though it has 9 letters), but given constraints, **"queueing"** is the closest fit. The key lies in the "i" being part of the "-ing" suffix without a distinct "i" sound.
**Final Answer:** queueing

# How Much Reasoning Is Necessary?



(a) Reasoning length.

# Optimizer

```python
import torch, torch.nn as nn, torch.optim as optim
class MyModel(nn.Module):
    def __init__(self):
        super().__init__()
        self.linear = nn.Linear(10, 5)  # Example linear layer

    def forward(self, x):
        return self.linear(x)

model = MyModel()
optimizer = optim.Adam(model.parameters(), lr=0.001)
loss_function = nn.MSELoss()  # Example Mean Squared Error loss

# ... (Inside the training loop) ...

# Forward pass
inputs = torch.randn(32, 10)
targets = torch.randn(32, 5)
outputs = model(inputs)
loss = loss_function(outputs, targets)

loss.backward() # Backpropagation
optimizer.step()  # Updates the parameters that were passed to it initially
optimizer.zero_grad()
```

## Optimizer Algorithms

## More than gradient descent

| | |
|---|---|
| Adadelta | Implements Adadelta algorithm. |
| Adafactor | Implements Adafactor algorithm. |
| Adagrad | Implements Adagrad algorithm. |
| Adam | Implements Adam algorithm. |
| AdamW | Implements AdamW algorithm. |
| SparseAdam | SparseAdam implements a masked version of the Adam algorithm suitable for sparse gradients. |
| Adamax | Implements Adamax algorithm (a variant of Adam based on infinity norm). |
| ASGD | Implements Averaged Stochastic Gradient Descent. |
| LBFGS | Implements L-BFGS algorithm. |
| NAdam | Implements NAdam algorithm. |
| RAdam | Implements RAdam algorithm. |
| RMSprop | Implements RMSprop algorithm. |
| Rprop | Implements the resilient backpropagation algorithm. |
| SGD | Implements stochastic gradient descent (optionally with momentum). |

```python
class CausalAttention(nn.Module):
    def __init__(self, d_in, d_out, context_length,
                 dropout, qkv_bias=False):
        super().__init__()
        self.d_out = d_out
        self.W_query = nn.Linear(d_in, d_out, bias=qkv_bias)
        self.W_key   = nn.Linear(d_in, d_out, bias=qkv_bias)
        self.W_value = nn.Linear(d_in, d_out, bias=qkv_bias)
        self.dropout = nn.Dropout(dropout)
        self.register_buffer(
            'mask',
            torch.triu(torch.ones(context_length, context_length),
            diagonal=1)
        )
```

# CausalAttention

```python
def forward(self, x):
    b, num_tokens, d_in = x.shape          # keep batch dimension at 0
    keys = self.W_key(x)
    queries = self.W_query(x)
    values = self.W_value(x)

    attn_scores = queries @ keys.transpose(1, 2)
    attn_scores.masked_fill_(                      Trailing underscore done in place
        self.mask.bool()[:num_tokens, :num_tokens], -torch.inf)
    attn_weights = torch.softmax(
        attn_scores / keys.shape[-1]**0.5, dim=-1
    )
    attn_weights = self.dropout(attn_weights)

    context_vec = attn_weights @ values
    return context_vec
```

# Multi-Headed



The embedded input tokens remain unchanged.

Inputs $X$

0.7 | 0.2 | 0.1

The values of the 5th row (input) are shown as an example.

Weight matrix $W_{q1}$

$W_{q2}$

Weight matrix $W_{k1}$

$W_{k2}$

Weight matrix $W_{v1}$

$W_{v2}$

Instead of one value weight matrix $W_v$ in single-head attention, use two matrices $W_{v1}$ and $W_{v2}$.

Instead of one query matrix $Q$, we have two query matrices $Q_1$ and $Q_2$.

Queries $Q_1$

$Q_2$

Keys $K_1$

$K_2$

Values $V_1$

$V_2$

We now have two sets of context vectors, $Z_1$ and $Z_2$.

For multi-head attention with two heads, we obtain two attention weight matrices, including causal and dropout masks.

Context vectors $Z_1$

-0.7 | -0.1

$Z_2$

0.7 | 0.4

Combined context vectors $Z$

-0.7 | -0.1 | 0.7 | 0.4

The context vector in $Z_2$ corresponding to the fifth input that was highlighted in the inputs $X$.

# The Cheap Version

class MultiHeadAttentionWrapper(nn.Module):

    def __init__(self, d_in, d_out, context_length,

            dropout, num_heads, qkv_bias=False):

      super().__init__()

      self.heads = **nn.ModuleList**(

        [CausalAttention(

           d_in, d_out, context_length, dropout, qkv_bias

         )

         for _ in range(num_heads)]

      )

    def forward(self, x):

      return torch.cat([head(x) for head in self.heads], dim=-1)

ModuleList

    Python list

    Registers its contents

Finally, we will use multiple transformer blocks to implement the untrained GPT model.

7) Final GPT architecture

Then we will combine building blocks 2–5, including the multi-head attention module from chapter 3, into a transformer block.

6) Transformer block

2) Layer normalization

3) GELU activation

4) Feed forward network

5) Shortcut connections

Next, we will implement building blocks 2–5.

1) GPT backbone

We developed a GPT placeholder model to see the overall structure of the model.

# DummyGPTModel

import torch
import torch.nn as nn

**We developed a GPT placeholder model to see the overall structure of the model.**

```python
class DummyGPTModel(nn.Module):
    def __init__(self, cfg):
        super().__init__()
        self.tok_emb = nn.Embedding(cfg["vocab_size"], cfg["emb_dim"])
        self.pos_emb = nn.Embedding(cfg["context_length"], cfg["emb_dim"])
        self.drop_emb = nn.Dropout(cfg["drop_rate"])
        self.trf_blocks = nn.Sequential(
            *[DummyTransformerBlock(cfg)
              for _ in range(cfg["n_layers"])]
        )
        self.final_norm = DummyLayerNorm(cfg["emb_dim"])
        self.out_head = nn.Linear(
            cfg["emb_dim"], cfg["vocab_size"], bias=False
        )
```

# torch.nn.Sequential

torch.nn.Sequential(*args: Module)

torch.nn.Sequential(arg: OrderedDict[str, Module])

Performing a transformation on the Sequential applies to each of the modules

```
model = nn.Sequential(
        nn.Conv2d(1,20,5),
        nn.ReLU(),
        nn.Conv2d(20,64,5),
        nn.ReLU()
    )
```

```
model = nn.Sequential(OrderedDict([
        ('conv1', nn.Conv2d(1,20,5)),
        ('relu1', nn.ReLU()),
        ('conv2', nn.Conv2d(20,64,5)),
        ('relu2', nn.ReLU())
    ]))
```

# Composite Pattern



Intent

    Compose objects into tree structures to represent part-whole hierarchies.
    Lets clients treat individual objects and compositions of objects uniformly

# DummyGPTModel

1) GPT backbone

We developed a GPT placeholder model to see the overall structure of the model.

```
def forward(self, in_idx):
    batch_size, seq_len = in_idx.shape
    tok_embeds = self.tok_emb(in_idx)
    pos_embeds = self.pos_emb(
        torch.arange(seq_len, device=in_idx.device)
    )
    x = tok_embeds + pos_embeds
    x = self.drop_emb(x)
    x = self.trf_blocks(x)
    x = self.final_norm(x)
    logits = self.out_head(x)
    return logits
```

logits are the raw, unnormalized output of the model's final layer before it's converted into probabilities.

# My LLM's outputs got 1000% better with this simple trick. Nikhil Anand

https://ai.gopubby.com/my-llms-outputs-got-1000-better-with-this-simple-trick-8403cf58691c



Logit transformations can cause
low probability tokens to exceed all others

Example output:
"The capital of Washington iseekek0q3n ee"

The goal is to generate the next word, "forward."

effort moves you **forward**

The number of input tokens matches the number of output tokens; hence, the first token ("Every") will not be contained in the output.

Output text

Postprocessing steps

The LLM returns one 768-dimensional output vector for each 768-dimensional input token embedding.

Outputs: | -1.2 | 0.3 |… | -0.1 | 0.4 |… | 0.5 | 1.6 |… | 0.0 | 1.6 |…

**GPT model**

For the smallest GPT-2 model, each embedding vector consists of 768 dimensions (only the first 2 dimensions are shown).

Token embeddings: | 2.4 | 2.4 |… | -2.6 | -1.3 |… | 2.0 | 1.8 |… | -1.6 | 2.1 |…

Token IDs: | 6109 | 3626 | 6100 | 345 |

We tokenize the input text and convert it into token embeddings.

Tokenized text: | Every | effort | moves | you |

Input text: | Every effort moves you |

Mean = 0.00
Variance = 1.00

**Zero-centered mean and unit variance after normalization**

Apply layer normalization

Mean = 0.13
Variance = 0.39

**Layer outputs**

**Layer inputs, where the five values represent a single training example**

Layer output values: 0.61, 1.41, -0.87, 0.58, -0.87, -0.87

Intermediate values: 0.22, 0.34, 0.00, 0.22, 0.00, 0.00

Layer input values: -0.11, 0.12, -0.36, -0.24, -1.19

```python
class LayerNorm(nn.Module):
    def __init__(self, emb_dim):
        super().__init__()
        self.eps = 1e-5
        self.scale = nn.Parameter(torch.ones(emb_dim))
        self.shift = nn.Parameter(torch.zeros(emb_dim))

    def forward(self, x):
        mean = x.mean(dim=-1, keepdim=True)
        var = x.var(dim=-1, keepdim=True, unbiased=False)
        norm_x = (x - mean) / torch.sqrt(var + self.eps) # no zero division
        return self.scale * norm_x + self.shift
```

nn.Parameter
marks a tensor as a learnable parameter

mean = x.mean(**dim=-1**, keepdim=True)

var = x.var(**dim=-1**, keepdim=True, unbiased=False)



dim=1 or dim=-1 calculates mean across the column dimension to obtain one mean per row

Mean

Input 1 | 0.22 | 0.34 | 0.00 | 0.22 | 0.00 | 0.00 → 0.13

Input 2 | 0.21 | 0.23 | 0.00 | 0.51 | 0.32 | 0.00 → 0.21

dim=0 calculates mean across the row dimension to obtain one mean per column

Input 1 | 0.22 | 0.34 | 0.00 | 0.22 | 0.00 | 0.00

Input 2 | 0.21 | 0.23 | 0.00 | 0.51 | 0.32 | 0.00

Mean   0.21   0.29   0.00   0.37   0.16   0.00

7) Final GPT architecture

6) Transformer block

**We implemented layer normalization.**

2) Layer normalization

3) GELU activation

4) Feed forward network

5) Shortcut connections

1) GPT backbone

**We implemented a GPT placeholder model to see the overall structure of the model.**

**Next, we will implement components 3 and 4.**

$$GELU(x) \approx 0.5 \cdot x \cdot \left(1 + tanh\left[\sqrt{\frac{2}{\pi}} \cdot \left(x + 0.044715 \cdot x^3\right)\right]\right)$$



GELU activation function



ReLU activation function

```python
class GELU(nn.Module):
    def __init__(self):
        super().__init__()

    def forward(self, x):
        return 0.5 * x * (1 + torch.tanh(
            torch.sqrt(torch.tensor(2.0 / torch.pi)) *
            (x + 0.044715 * torch.pow(x, 3))
        ))
```

$$GELU(x) \approx 0.5 \cdot x \cdot \left(1 + tanh\left[\sqrt{\frac{2}{\pi}} \cdot \left(x + 0.044715 \cdot x^3\right)\right]\right)$$

```python
class FeedForward(nn.Module):
    def __init__(self, cfg):
        super().__init__()
        self.layers = nn.Sequential(
            nn.Linear(cfg["emb_dim"], 4 * cfg["emb_dim"]),
            GELU(),
            nn.Linear(4 * cfg["emb_dim"], cfg["emb_dim"]),
        )

    def forward(self, x):
        return self.layers(x)
```

Output tensor with
shape (2, 3, 768)

Linear layer
Input:   (2, 3, 3072)
Output: (2, 3, 768)

The second line:
decreases the e
dimension by a

GELU activation
Input:   (2, 3, 3072)
Output: (2, 3, 3072)

Linear layer
Input:   (2, 3, 768)
Output: (2, 3, 3072)

The first linear
increases the e
dimension by a

Input tensor with
shape (2, 3, 768)

The three values represent
the batch size (2), number
of tokens (3), and embedding
size (768).

Outputs

Linear layer 2

**The second linear layer shrinks the outputs by a factor of 4, so that they match the original input dimensions.**

Linear layer 1

**The inputs are projected into a four-times larger space via the first linear layer.**

Inputs

# Shortcut Connections



**Deep neural network**

Layer 5
Gradient: 0.0050
{ GELU
Linear }

Layer 4
Gradient: 0.0013
{ GELU
Linear }

Layer 3
Gradient: 0.0007
{ GELU
Linear }

Layer 2
Gradient: 0.0001
{ GELU
Linear }

Layer 1
Gradient: 0.0002
{ GELU
Linear }

[1.0, 0.0, -1.0]

In very deep networks, the gradient values in early layers become vanishingly small

**Deep neural network with shortcut connections**

Layer 5
Gradient: 1.32
{ GELU
Linear }

Layer 4
Gradient: 0.26
{ GELU
Linear }

Layer 3
Gradient: 0.32
{ GELU
Linear }

Layer 2
Gradient: 0.20
{ GELU
Linear }

Layer 1
Gradient: 0.22
{ GELU
Linear }

[1.0, 0.0, -1.0]

Shortcut connection adds input values to the outputs of layer 1

The shortcut connections help with maintaining relatively large gradient values even in early layers

34

# Shortcut Connections

```
class ExampleDeepNeuralNetwork(nn.Module):
    def __init__(self, layer_sizes, use_shortcut):
        super().__init__()
        self.use_shortcut = use_shortcut
        self.layers = nn.ModuleList([      #1
            nn.Sequential(nn.Linear(layer_sizes[0], layer_sizes[1]),
                    GELU()),
            nn.Sequential(nn.Linear(layer_sizes[1], layer_sizes[2]),
                    GELU()),
            nn.Sequential(nn.Linear(layer_sizes[2], layer_sizes[3]),
                    GELU()),
            nn.Sequential(nn.Linear(layer_sizes[3], layer_sizes[4]),
                    GELU()),
            nn.Sequential(nn.Linear(layer_sizes[4], layer_sizes[5]),
                    GELU())
        ])
```

# Shortcut Connections

```python
def forward(self, x):
    for layer in self.layers:
        layer_output = layer(x)          #2
        if self.use_shortcut and x.shape == layer_output.shape:     #3
            x = x + layer_output
        else:
            x = layer_output
    return x
```

# Backtracking

```python
def print_gradients(model, x):
    output = model(x)
    target = torch.tensor([[0.]])

    loss = nn.MSELoss()
    loss = loss(output, target)

    loss.backward()

    for name, param in model.named_parameters():
        if 'weight' in name:
            print(f"{name} has gradient mean of {param.grad.abs().mean().item()}")
```

# nn.MSELoss

mean squared error (squared L2 norm) between
each element in the input x and target y

# Existing Loss Functions (21)

| | |
|---|---|
| nn.L1Loss | Creates a criterion that measures the mean absolute error (MAE) |
| nn.MSELoss | Creates a criterion that measures the mean squared error (squared L2 norm) |
| nn.CrossEntropyLoss | This criterion computes the cross entropy loss between input logits and target. |
| nn.CTCLoss | The Connectionist Temporal Classification loss. |
| nn.NLLLoss | The negative log likelihood loss. |
| nn.PoissonNLLLoss | Negative log likelihood loss with Poisson distribution of target. |
| nn.GaussianNLLLoss | Gaussian negative log likelihood loss. |
| nn.KLDivLoss | The Kullback-Leibler divergence loss. |
| nn.BCELoss | Creates a criterion that measures the Binary Cross Entropy between the target and the input probabilities: |
| nn.BCEWithLogitsLoss | This loss combines a Sigmoid layer and the BCELoss in one single class. |
| nn.MarginRankingLoss | |
| nn.HingeEmbeddingLoss | Measures the loss given an input tensor X and a labels tensor y (containing 1 or -1). |

**Finally, we will use multiple transformer blocks to implement the untrained GPT model.**

7) Final GPT architecture

**Then we will combine building blocks 2–5, including the multi-head attention module from chapter 3, into a transformer block.**

6) Transformer block

2) Layer normalization

3) GELU activation

4) Feed forward network

5) Shortcut connections

**Next, we will implement building blocks 2–5.**

1) GPT backbone

**We developed a GPT placeholder model to see the overall structure of the model.**

# Transformer Block



```
[[-0.0256,  ...,   0.6890],
 [-0.0178,  ...,   0.7431],
 [ 0.4558,  ...,   0.7814],
 [ 0.0702,  ...,   0.7134]]
```

**Outputs have the same form and dimensions as the inputs.**

**The transformer block**

Dropout

Feed forward

LayerNorm 2

Dropout

Masked multi-head attention

LayerNorm 1

Linear layer

GELU activation

Linear layer

**A view into the "feed forward" block**

**Shortcut connection**

**The input tokens to be embedded**

```
Every    ──────►  [[0.2961,  ...,  0.4604],
effort   ──────►   [0.2238,  ...,  0.7598],
moves    ──────►   [0.6945,  ...,  0.5963],
you      ──────►   [0.0890,  ...,  0.5833]]
```

**This tensor represents an embedded text sample that serves as input to the transformer block.**

**Each row is a 768-dimensional vector representing an embedded input token.**

# TransformerBlock

```python
class TransformerBlock(nn.Module):
    def __init__(self, cfg):
        super().__init__()
        self.att = MultiHeadAttention(
            d_in=cfg["emb_dim"],
            d_out=cfg["emb_dim"],
            context_length=cfg["context_length"],
            num_heads=cfg["n_heads"],
            dropout=cfg["drop_rate"],
            qkv_bias=cfg["qkv_bias"])
        self.ff = FeedForward(cfg)
        self.norm1 = LayerNorm(cfg["emb_dim"])
        self.norm2 = LayerNorm(cfg["emb_dim"])
        self.drop_shortcut = nn.Dropout(cfg["drop_rate"])

    def forward(self, x):
        shortcut = x
        x = self.norm1(x)
        x = self.att(x)
        x = self.drop_shortcut(x)
        x = x + shortcut

        shortcut = x
        x = self.norm2(x)
        x = self.ff(x)
        x = self.drop_shortcut(x)
        x = x + shortcut
        return x
```

# Full Model

```python
class GPTModel(nn.Module):
    def __init__(self, cfg):
        super().__init__()
        self.tok_emb = nn.Embedding(cfg["vocab_size"], cfg["emb_dim"])
        self.pos_emb = nn.Embedding(cfg["context_length"], cfg["emb_dim"])
        self.drop_emb = nn.Dropout(cfg["drop_rate"])

        self.trf_blocks = nn.Sequential(
            *[TransformerBlock(cfg) for _ in range(cfg["n_layers"])])

        self.final_norm = LayerNorm(cfg["emb_dim"])
        self.out_head = nn.Linear(
            cfg["emb_dim"], cfg["vocab_size"], bias=False
        )
```

# Full Model

```python
def forward(self, in_idx):
    batch_size, seq_len = in_idx.shape
    tok_embeds = self.tok_emb(in_idx)
    pos_embeds = self.pos_emb(
        torch.arange(seq_len, device=in_idx.device)
    )
    x = tok_embeds + pos_embeds
    x = self.drop_emb(x)
    x = self.trf_blocks(x)
    x = self.final_norm(x)
    logits = self.out_head(x)
    return logits
```

# Put it all together

Model is 152 lines of code

```python
def main():
    GPT_CONFIG_124M = {… }

    torch.manual_seed(123)
    model = GPTModel(GPT_CONFIG_124M)
    model.eval()  # disable dropout

    start_context = "Hello, I am"

    tokenizer = tiktoken.get_encoding("gpt2")
    encoded = tokenizer.encode(start_context)
    encoded_tensor = torch.tensor(encoded).unsqueeze(0)

#removed fancy print statements

    out = generate_text_simple(
        model=model,
        idx=encoded_tensor,
        max_new_tokens=10,
        context_size=GPT_CONFIG_124M["context_length"]
    )
    decoded_text = tokenizer.decode(out.squeeze(0).tolist())

#removed fancy print statements
```

# Put it all together

```python
def generate_text_simple(model, idx, max_new_tokens, context_size):
    # idx is (B, T) array of indices in the current context
    for _ in range(max_new_tokens):

        # Crop current context if it exceeds the supported context size
        idx_cond = idx[:, -context_size:]

        # Get the predictions
        with torch.no_grad():
            logits = model(idx_cond)

        # Focus only on the last time step
        # (batch, n_token, vocab_size) becomes (batch, vocab_size)
        logits = logits[:, -1, :]

        # Get the idx of the vocab entry with the highest logits value
        idx_next = torch.argmax(logits, dim=-1, keepdim=True)  # (batch, 1)

        # Append sampled index to the running sequence
        idx = torch.cat((idx, idx_next), dim=1)  # (batch, n_tokens+1)

    return idx
```

# Put it all together

```
=================================================
             IN
=================================================

Input text: Hello, I am
Encoded input text: [15496, 11, 314, 716]
encoded_tensor.shape: torch.Size([1, 4])


=================================================
             OUT
=================================================

Output: tensor([[15496,    11,   314,   716, 27018, 24086, 47843, 30961, 42348,  7267,
         49706, 43231, 47062, 34657]])
Output length: 14
```

**Output text: Hello, I am Featureiman Byeswickattribute argue logger Normandy Compton analogous**

# Chapter 5 - Pretraining



In the previous chapter, we implemented a GPT-like LLM architecture.

For the pretraining, we will implement the training loop along with model evaluation metrics.

Finally, we load openly available pretrained weights into the model.

STAGE 1

1) Data preparation & sampling

2) Attention mechanism

3) LLM architecture

STAGE 3

Dataset with class labels

5) Training loop

6) Model evaluation

7) Load pretrained weights

8) Fine-tuning

Classifier

Building an LLM

4) Pretraining

Foundation model

Personal assistant

STAGE 2

9) Fine-tuning

In this chapter, we will pretrain the LLM model.

Instruction dataset

# Preview

Chapter 5
    Pretraining
    Loading model weights from Open AI

Scripts to do both

# gpt_generate.py

Load and use the pretrained model weights from OpenAI

```
CHOOSE_MODEL = "gpt2-small (124M)"
INPUT_PROMPT = "Every effort moves you"
```

(base) rwhitney@127 01_main-chapter-code % python gpt_generate.py
File already exists and is up-to-date: gpt2/124M/checkpoint
File already exists and is up-to-date: gpt2/124M/encoder.json
File already exists and is up-to-date: gpt2/124M/hparams.json
File already exists and is up-to-date: gpt2/124M/model.ckpt.data-00000-of-00001
File already exists and is up-to-date: gpt2/124M/model.ckpt.index
File already exists and is up-to-date: gpt2/124M/model.ckpt.meta
File already exists and is up-to-date: gpt2/124M/vocab.bpe
Output text:
 Every effort moves you toward finding an ideal life. You don't have to accept your problems by trying to remedy them, because that would be foolish

Second run

Note they saved model etc.

# gpt_generate.py

Load and use the pretrained model weights from OpenAI

```
CHOOSE_MODEL = "gpt2-medium (355M)"
INPUT_PROMPT = "Every effort moves you"
It's the same seed as before
```

```
(base) rwhitney@127 01_main-chapter-code % python gpt_generate.py
checkpoint: 100%|| 77.0/77.0 [00:00<00:00, 21.3kiB/s]
encoder.json: 100%|| 1.04M/1.04M [00:00<00:00, 1.76MiB/s]
hparams.json: 100%|| 91.0/91.0 [00:00<00:00, 9.01kiB/s]
model.ckpt.data-00000-of-00001: 100%|| 1.42G/1.42G [05:31<00:00, 4.28MiB/s]
model.ckpt.index: 100%|| 10.4k/10.4k [00:00<00:00, 831kiB/s]
model.ckpt.meta: 100%|| 927k/927k [00:00<00:00, 1.76MiB/s]
vocab.bpe: 100%|| 456k/456k [00:00<00:00, 997kiB/s]
Output text:

 Every effort moves you toward balance." But it seems that these values have been
forgotten by both parties.

If Congress is to fulfill these basic
```

# gpt_train.py

Trains the model on the-verdict.txt

Ep 1 (Step 000000): Train loss 9.812, Val loss 9.846

Ep 1 (Step 000005): Train loss 7.588, Val loss 8.044

Every effort moves you,,,,,,,,,,,.

Ep 2 (Step 000010): Train loss 6.582, Val loss 6.800

Ep 2 (Step 000015): Train loss 5.920, Val loss 6.592

Every effort moves you, and, and, and, and, and, and, and, and, and, and, and, and, and, and, and, and, and, and, and, and, and, and, and

…

Ep 10 (Step 000085): Train loss 0.358, Val loss 6.542

Every effort moves you?"  "Yes--quite insensible to the irony. She wanted him vindicated--and by me He laughed again, and threw back his head to look up at the sketch of the donkey. "There were days when I

Traceback (most recent call last):
  File "/Users/rwhitney/Courses/696/Spring 2025 LLM/Notebooks/LLMs-from-scratch/ch05/01_main-chapter-code/gpt_train.py", line 236, in <module>
    plot_losses(epochs_tensor, tokens_seen, train_losses, val_losses)
  File "/Users/rwhitney/Courses/696/Spring 2025 LLM/Notebooks/LLMs-from-scratch/ch05/01_main-chapter-code/gpt_train.py", line 116, in plot_losses

# Model to Train

```python
import torch
from previous_chapters import GPTModel

GPT_CONFIG_124M = {
    "vocab_size": 50257,   # Vocabulary size
    "context_length": 256, # Shortened context length (orig: 1024)
    "emb_dim": 768,        # Embedding dimension
    "n_heads": 12,         # Number of attention heads
    "n_layers": 12,        # Number of layers
    "drop_rate": 0.1,      # Dropout rate
    "qkv_bias": False      # Query-key-value bias
}

torch.manual_seed(123)
model = GPTModel(GPT_CONFIG_124M)
model.eval();  # Disable dropout during inference
```

# Training

Convert text to tokens

Need inputs and targets

Determine how "off" model(inputs) are from targets

Use loss function to adjust the weights

# Inputs and Targets

```python
inputs = torch.tensor([[16833, 3626, 6100],   # ["every effort moves",
                       [40,    1107, 588]])   #  "I really like"]

targets = torch.tensor([[3626, 6100, 345  ],  # [" effort moves you",
                        [1107,  588, 11311]]) #  " really like chocolate"]


with torch.no_grad():
    logits = model(inputs)

probas = torch.softmax(logits, dim=-1) # Probability of each token in vocabulary
print(probas.shape) # Shape: (batch_size, num_tokens, vocab_size)


    torch.Size([2, 3, 50257])
```

# softmax

Converts collection of values to probabilities

Exponential of each value

Normalize result

| Logits | Exponential | Normalized |
|--------|-------------|------------|
| 2.5 | 12.18 | 0.7856 |
| 1.0 | 2.72 | 0.1753 |
| -0.5 | 0.61 | 0.0391 |

| Logits | Softmax |
|--------|---------|
| 5.5 | 0.9866 |
| 1.0 | 0.0110 |
| -0.5 | 0.0024 |

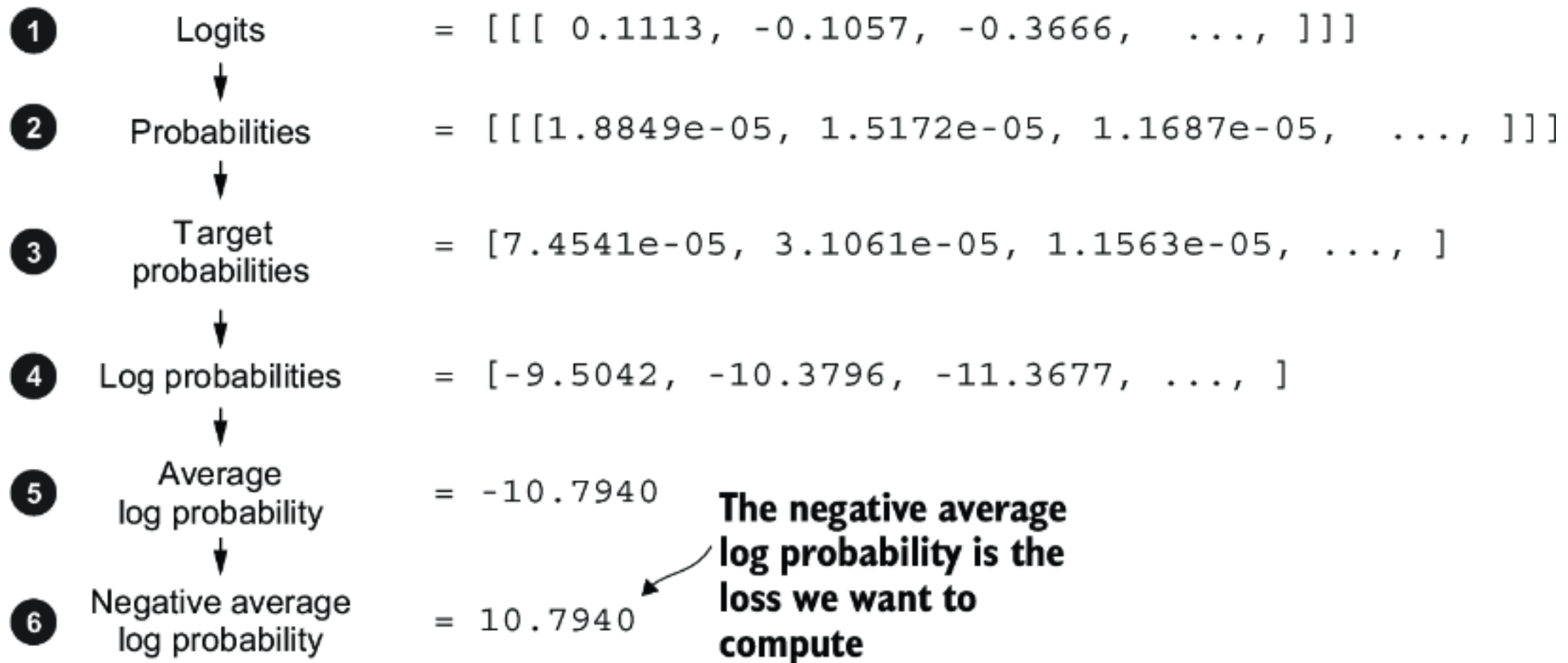| Logits | Softmax |
|--------|---------|
| 10.5 | 0.9991 |
| 1.0 | 7.4845E-05 |
| -0.5 | 1.6700E-05 |

```
import torch

x = torch.tensor([2.5, 1.0, -0.5])
y = torch.softmax(x, dim=0)

print(y)
```

tensor([0.7856, 0.1753, 0.0391])

Extreme values push softmax results to 1 & 0

# Backpropagation - cross_entropy

1. Logits = [[[ 0.1113, -0.1057, -0.3666, ..., ]]]

2. Probabilities = [[[1.8849e-05, 1.5172e-05, 1.1687e-05, ..., ]]]

3. Target probabilities = [7.4541e-05, 3.1061e-05, 1.1563e-05, ..., ]

4. Log probabilities = [-9.5042, -10.3796, -11.3677, ..., ]

5. Average log probability = -10.7940

6. Negative average log probability = 10.7940

The negative average log probability is the loss we want to compute
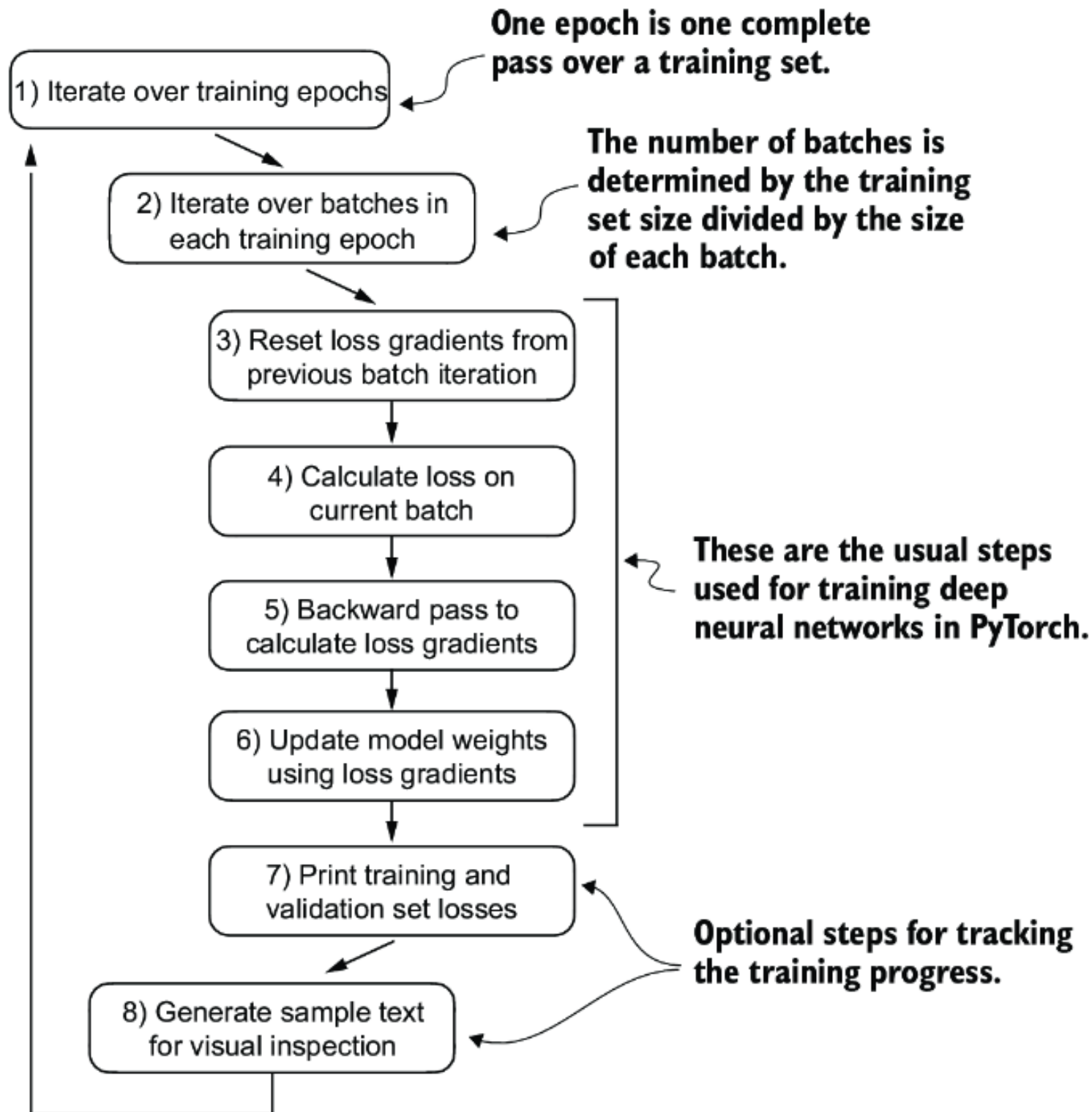
loss = torch.nn.functional.cross_entropy(logits_flat, targets_flat)

# Perplexity

How well probability distribution given by the model
matches the actual distribution of the words in the dataset

How uncertain a model is about the next word in a sequence.

# Training an LLM

```python
def train_model_simple(model, train_loader, val_loader,
                optimizer, device, num_epochs,
                eval_freq, eval_iter, start_context, tokenizer):
    train_losses, val_losses, track_tokens_seen = [], [], []
    tokens_seen, global_step = 0, -1

    for epoch in range(num_epochs):
        model.train()
        for input_batch, target_batch in train_loader:
            optimizer.zero_grad()
            loss = calc_loss_batch(
                input_batch, target_batch, model, device
            )
            loss.backward()
            optimizer.step()
            tokens_seen += input_batch.numel()
            global_step += 1

            if global_step % eval_freq == 0:
                train_loss, val_loss = evaluate_model(
                    model, train_loader, val_loader, device, eval_iter)
                train_losses.append(train_loss)
                val_losses.append(val_loss)
                track_tokens_seen.append(tokens_seen)
                print(f"Ep {epoch+1} (Step {global_step:06d}): "
                    f"Train loss {train_loss:.3f}, "
                    f"Val loss {val_loss:.3f}"
                )

        generate_and_print_sample(
            model, tokenizer, device, start_context
        )
    return train_losses, val_losses, track_tokens_seen
```