

```
!pip install tiktoken
```

```

Collecting tiktoken
  Downloading tiktoken-0.9.0-cp311-cp311-manylinux_2_17_x86_64.manylinux201
Requirement already satisfied: regex>=2022.1.18 in /usr/local/lib/python3.1
Requirement already satisfied: requests>=2.26.0 in /usr/local/lib/python3.1
Requirement already satisfied: charset-normalizer<4,>=2 in /usr/local/lib/p
Requirement already satisfied: idna<4,>=2.5 in /usr/local/lib/python3.11/di
Requirement already satisfied: urllib3<3,>=1.21.1 in /usr/local/lib/python3
Requirement already satisfied: certifi>=2017.4.17 in /usr/local/lib/python3
Downloading tiktoken-0.9.0-cp311-cp311-manylinux_2_17_x86_64.manylinux2014_
1.2/1.2 MB 11.2 MB/s eta 0:00:0
Installing collected packages: tiktoken
Successfully installed tiktoken-0.9.0

```

```
from importlib.metadata import version
```

```

print("matplotlib version:", version("matplotlib"))
print("torch version:", version("torch"))
print("tiktoken version:", version("tiktoken"))

```

```

matplotlib version: 3.10.0
torch version: 2.5.1+cu124
tiktoken version: 0.9.0

```

```

GPT_CONFIG_124M = {
    "vocab_size": 50257,      # Vocabulary size
    "context_length": 1024,   # Context length
    "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
}

```

```

import torch
import torch.nn as nn

```

```

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"])

        # Use a placeholder for TransformerBlock
        self.trf_blocks = nn.Sequential(
            *[DummyTransformerBlock(cfg) for _ in range(cfg["n_layers"])]
        )

```

```

        # Use a placeholder for LayerNorm
        self.final_norm = DummyLayerNorm(cfg["emb_dim"])
        self.out_head = nn.Linear(
            cfg["emb_dim"], cfg["vocab_size"], bias=False
        )

    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

class DummyTransformerBlock(nn.Module):
    def __init__(self, cfg):
        super().__init__()
        # A simple placeholder

    def forward(self, x):
        # This block does nothing and just returns its input.
        return x

class DummyLayerNorm(nn.Module):
    def __init__(self, normalized_shape, eps=1e-5):
        super().__init__()
        # The parameters here are just to mimic the LayerNorm interface.

    def forward(self, x):
        # This layer does nothing and just returns its input.
        return x

import tiktoken

tokenizer = tiktoken.get_encoding("gpt2")

batch = []

txt1 = "Every effort moves you"
txt2 = "Every day holds a"

batch.append(torch.tensor(tokenizer.encode(txt1)))
batch.append(torch.tensor(tokenizer.encode(txt2)))
batch = torch.stack(batch, dim=0)
print(batch)

tensor([[6109, 3626, 6100, 345],
        [6100, 1110, 6622, 25711]])

```

```
[0.109, 1.110, 0.022, 2.571]]
```

```
torch.manual_seed(123)
model = DummyGPTModel(GPT_CONFIG_124M)
```

```
logits = model(batch)
print("Output shape:", logits.shape)
print(logits)
```

```
Output shape: torch.Size([2, 4, 50257])
tensor([[[[-0.9289,  0.2748, -0.7557, ..., -1.6070,  0.2702, -0.5888],
          [-0.4476,  0.1726,  0.5354, ..., -0.3932,  1.5285,  0.8557],
          [ 0.5680,  1.6053, -0.2155, ...,  1.1624,  0.1380,  0.7425],
          [ 0.0447,  2.4787, -0.8843, ...,  1.3219, -0.0864, -0.5856]],

        [[[-1.5474, -0.0542, -1.0571, ..., -1.8061, -0.4494, -0.6747],
          [-0.8422,  0.8243, -0.1098, ..., -0.1434,  0.2079,  1.2046],
          [ 0.1355,  1.1858, -0.1453, ...,  0.0869, -0.1590,  0.1552],
          [ 0.1666, -0.8138,  0.2307, ...,  2.5035, -0.3055, -0.3083]]],
        grad_fn=<UnsafeViewBackward0>)
```

```
torch.manual_seed(123)
```

```
# create 2 training examples with 5 dimensions (features) each
batch_example = torch.randn(2, 5)
```

```
layer = nn.Sequential(nn.Linear(5, 6), nn.ReLU())
out = layer(batch_example)
print(out)
```

```
tensor([[0.2260, 0.3470, 0.0000, 0.2216, 0.0000, 0.0000],
        [0.2133, 0.2394, 0.0000, 0.5198, 0.3297, 0.0000]],
        grad_fn=<ReluBackward0>)
```

```
mean = out.mean(dim=-1, keepdim=True)
var = out.var(dim=-1, keepdim=True)
```

```
print("Mean:\n", mean)
print("Variance:\n", var)
```

```
Mean:
tensor([[0.1324],
        [0.2170]], grad_fn=<MeanBackward1>)
Variance:
tensor([[0.0231],
        [0.0398]], grad_fn=<VarBackward0>)
```

```
torch.set_printoptions(sci_mode=False)
```

```

..._mean_...,
print("Mean:\n", mean)
print("Variance:\n", var)

```

```

Mean:
  tensor([[0.1324],
          [0.2170]], grad_fn=<MeanBackward1>)
Variance:
  tensor([[0.0231],
          [0.0398]], grad_fn=<VarBackward0>)

```

```

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)
        return self.scale * norm_x + self.shift

```

```

ln = LayerNorm(emb_dim=5)
out_ln = ln(batch_example)

```

```

mean = out_ln.mean(dim=-1, keepdim=True)
var = out_ln.var(dim=-1, unbiased=False, keepdim=True)

print("Mean:\n", mean)
print("Variance:\n", var)

```

```

Mean:
  tensor([[ -0.0000],
          [ 0.0000]], grad_fn=<MeanBackward1>)
Variance:
  tensor([[1.0000],
          [1.0000]], grad_fn=<VarBackward0>)

```

```

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))
        ))

```

```
import matplotlib.pyplot as plt
```

```
gelu, relu = GELU(), nn.ReLU()
```

```

# Some sample data
x = torch.linspace(-3, 3, 100)
y_gelu, y_relu = gelu(x), relu(x)

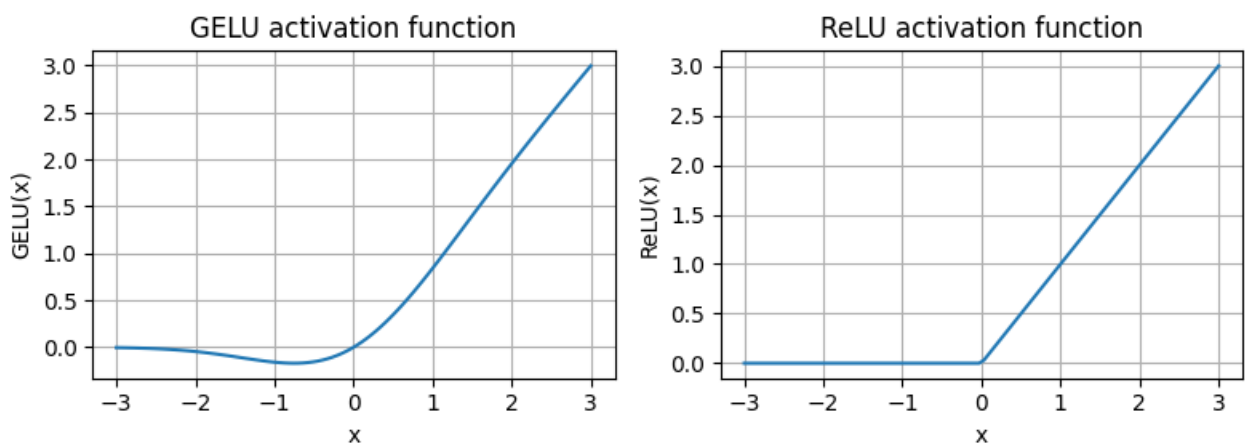
```

```

plt.figure(figsize=(8, 3))
for i, (y, label) in enumerate(zip([y_gelu, y_relu], ["GELU", "ReLU"]), 1):
    plt.subplot(1, 2, i)
    plt.plot(x, y)
    plt.title(f"{label} activation function")
    plt.xlabel("x")
    plt.ylabel(f"{label}(x)")
    plt.grid(True)

plt.tight_layout()
plt.show()

```



```

class FeedForward(nn.Module):
    def __init__(self, cfg):
        super().__init__()
        self.layers = nn.Sequential(
            nn.Linear(cfg["emb dim"] 4 * cfg["emb dim"])

```

```

        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)

print(GPT_CONFIG_124M["emb_dim"])

768

ffn = FeedForward(GPT_CONFIG_124M)

# input shape: [batch_size, num_token, emb_size]
x = torch.rand(2, 3, 768)
out = ffn(x)
print(out.shape)

torch.Size([2, 3, 768])

class ExampleDeepNeuralNetwork(nn.Module):
    def __init__(self, layer_sizes, use_shortcut):
        super().__init__()
        self.use_shortcut = use_shortcut
        self.layers = nn.ModuleList([
            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())
        ])

    def forward(self, x):
        for layer in self.layers:
            # Compute the output of the current layer
            layer_output = layer(x)
            # Check if shortcut can be applied
            if self.use_shortcut and x.shape == layer_output.shape:
                x = x + layer_output
            else:
                x = layer_output
        return x

def print_gradients(model, x):
    # Forward pass
    output = model(x)
    target = torch.tensor([[0.]])

```

```

# Calculate loss based on how close the target
# and output are
loss = nn.MSELoss()
loss = loss(output, target)

# Backward pass to calculate the gradients
loss.backward()

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

```

```
layer_sizes = [3, 3, 3, 3, 3, 1]
```

```
sample_input = torch.tensor([[1., 0., -1.]])
```

```
torch.manual_seed(123)
```

```
model_without_shortcut = ExampleDeepNeuralNetwork(
    layer_sizes, use_shortcut=False
)
```

```
print_gradients(model_without_shortcut, sample_input)
```

```

layers.0.0.weight has gradient mean of 0.00020173584925942123
layers.1.0.weight has gradient mean of 0.00012011159560643137
layers.2.0.weight has gradient mean of 0.0007152040489017963
layers.3.0.weight has gradient mean of 0.0013988736318424344
layers.4.0.weight has gradient mean of 0.005049645435065031

```

```
torch.manual_seed(123)
```

```
model_with_shortcut = ExampleDeepNeuralNetwork(
    layer_sizes, use_shortcut=True
)
```

```
print_gradients(model_with_shortcut, sample_input)
```

```

layers.0.0.weight has gradient mean of 0.22169791162014008
layers.1.0.weight has gradient mean of 0.20694105327129364
layers.2.0.weight has gradient mean of 0.32896995544433594
layers.3.0.weight has gradient mean of 0.2665732204914093
layers.4.0.weight has gradient mean of 1.3258540630340576

```

```
class MultiHeadAttention(nn.Module):
```

```

    def __init__(self, d_in, d_out, context_length, dropout, num_heads, qkv_bia
        super().__init__()
        assert d_out % num_heads == 0, "d_out must be divisible by num_heads"

```

```

        self.d_out = d_out
        self.num heads = num heads

```

```

self.head_dim = d_out // num_heads # Reduce the projection dim to match

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.out_proj = nn.Linear(d_out, d_out) # Linear layer to combine heads
self.dropout = nn.Dropout(dropout)
self.register_buffer('mask', torch.triu(torch.ones(context_length, context_length), diagonal=1))

def forward(self, x):
    b, num_tokens, d_in = x.shape

    keys = self.W_key(x) # Shape: (b, num_tokens, d_out)
    queries = self.W_query(x)
    values = self.W_value(x)

    # We implicitly split the matrix by adding a `num_heads` dimension
    # Unroll last dim: (b, num_tokens, d_out) -> (b, num_tokens, num_heads, head_dim)
    keys = keys.view(b, num_tokens, self.num_heads, self.head_dim)
    values = values.view(b, num_tokens, self.num_heads, self.head_dim)
    queries = queries.view(b, num_tokens, self.num_heads, self.head_dim)

    # Transpose: (b, num_tokens, num_heads, head_dim) -> (b, num_heads, num_tokens, head_dim)
    keys = keys.transpose(1, 2)
    queries = queries.transpose(1, 2)
    values = values.transpose(1, 2)

    # Compute scaled dot-product attention (aka self-attention) with a causal mask
    attn_scores = queries @ keys.transpose(2, 3) # Dot product for each head

    # Original mask truncated to the number of tokens and converted to bool
    mask_bool = self.mask.bool()[:num_tokens, :num_tokens]

    # Use the mask to fill attention scores
    attn_scores.masked_fill_(mask_bool, -torch.inf)

    attn_weights = torch.softmax(attn_scores / keys.shape[-1]**0.5, dim=-1)
    attn_weights = self.dropout(attn_weights)

    # Shape: (b, num_tokens, num_heads, head_dim)
    context_vec = (attn_weights @ values).transpose(1, 2)

    # Combine heads, where self.d_out = self.num_heads * self.head_dim
    context_vec = context_vec.contiguous().view(b, num_tokens, self.d_out)
    context_vec = self.out_proj(context_vec) # optional projection

    return context_vec

```

```

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 connection for attention block
    shortcut = x
    x = self.norm1(x)
    x = self.att(x) # Shape [batch_size, num_tokens, emb_size]
    x = self.drop_shortcut(x)
    x = x + shortcut # Add the original input back

    # Shortcut connection for feed forward block
    shortcut = x
    x = self.norm2(x)
    x = self.ff(x)
    x = self.drop_shortcut(x)
    x = x + shortcut # Add the original input back

    return x

```

```
torch.manual_seed(123)
```

```

x = torch.rand(2, 4, 768) # Shape: [batch_size, num_tokens, emb_dim]
block = TransformerBlock(GPT_CONFIG_124M)
output = block(x)

```

```

print("Input shape:", x.shape)
print("Output shape:", output.shape)

```

```

Input shape: torch.Size([2, 4, 768])
Output shape: torch.Size([2, 4, 768])

```

```

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
        )

    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  # Shape [batch_size, num_tokens, emb_size]
        x = self.drop_emb(x)
        x = self.trf_blocks(x)
        x = self.final_norm(x)
        logits = self.out_head(x)
        return logits

```

```

torch.manual_seed(123)
model = GPTModel(GPT_CONFIG_124M)

```

```

out = model(batch)
print("Input batch:\n", batch)
print("\nOutput shape:", out.shape)
print(out)

```

Input batch:

```

tensor([[6109, 3626, 6100, 345],
        [6109, 1110, 6622, 257]])

```

Output shape: torch.Size([2, 4, 50257])

```

tensor([[[ 0.1381,  0.0077, -0.1963, ..., -0.0222, -0.1060,  0.1717],
          [ 0.3865, -0.8408, -0.6564, ..., -0.5163,  0.2369, -0.3357],
          [ 0.6989, -0.1829, -0.1631, ...,  0.1472, -0.6504, -0.0056],
          [-0.4290,  0.1669, -0.1258, ...,  1.1579,  0.5303, -0.5549]],

        [[ 0.1094, -0.2894, -0.1467, ..., -0.0557,  0.2911, -0.2824],
          [ 0.0882, -0.3552, -0.3527, ...,  1.2930,  0.0053,  0.1898],
          [ 0.6091,  0.4702, -0.4094, ...,  0.7688,  0.3787, -0.1974],
          [-0.0612, -0.0737,  0.4751, ...,  1.2463, -0.3834,  0.0609]]],
        grad_fn=<UnsafeViewBackward0>)

```

```

total_params = sum(p.numel() for p in model.parameters())
print(f"Total number of parameters: {total_params:,}")

```

Total number of parameters: 163,009,536

```

print("Token embedding layer shape:", model.tok_emb.weight.shape)
print("Output layer shape:", model.out_head.weight.shape)

```

Token embedding layer shape: torch.Size([50257, 768])

Output layer shape: torch.Size([50257, 768])

```
output_layer_shape: torch.Size([50257, 768])
```

```
total_params_gpt2 = total_params - sum(p.numel() for p in model.out_head.parameters())
print(f"Number of trainable parameters considering weight tying: {total_params_gpt2}")
```

```
Number of trainable parameters considering weight tying: 124,412,160
```

```
# Calculate the total size in bytes (assuming float32, 4 bytes per parameter)
total_size_bytes = total_params * 4
```

```
# Convert to megabytes
```

```
total_size_mb = total_size_bytes / (1024 * 1024)
```

```
print(f"Total size of the model: {total_size_mb:.2f} MB")
```

```
Total size of the model: 621.83 MB
```

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

        # Crop current context if it exceeds the supported context size
        # E.g., if LLM supports only 5 tokens, and the context size is 10
        # then only the last 5 tokens are used as context
        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_tokens, vocab_size) becomes (batch, vocab_size)
        logits = logits[:, -1, :]

        # Apply softmax to get probabilities
        probas = torch.softmax(logits, dim=-1) # (batch, vocab_size)

        # Get the idx of the vocab entry with the highest probability value
        idx_next = torch.argmax(probas, 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
```

```
start_context = "Hello, I am"
```

```
encoded = tokenizer.encode(start_context)
print("encoded:", encoded)
```

```
encoded_tensor = torch.tensor(encoded).unsqueeze(0)
print("encoded_tensor.shape:", encoded_tensor.shape)
```

```
encoded: [15496, 11, 314, 716]
encoded_tensor.shape: torch.Size([1, 4])
```

```
model.eval() # disable dropout
```

```
out = generate_text_simple(
    model=model,
    idx=encoded_tensor,
    max_new_tokens=6,
    context_size=GPT_CONFIG_124M["context_length"]
)
```

```
print("Output:", out)
print("Output length:", len(out[0]))
```

```
Output: tensor([[15496,    11,   314,   716, 27018, 24086, 47843, 30961, 42
Output length: 10
```

```
decoded_text = tokenizer.decode(out.squeeze(0).tolist())
print(decoded_text)
```

```
Hello, I am Featureiman Byeswickattribute argue
```

