I.      Model architecture and parameters

The MiniGPT model is a simplified version of a Generative Pre-trained Transformer (GPT) model. It consists of the following main components:

1. **Token Embedding Table (token_embedding_table):** This layer converts input token IDs into dense vector representations (embeddings). The size of the embedding is D_MODEL.
2. **Position Embedding Table (position_embedding_table):** This layer provides positional information for the tokens in the sequence. It adds a learned embedding based on the token's position to the token embedding. The maximum sequence length is T.
3. **Transformer Blocks (blocks):** This is the core of the model, consisting of a sequence of N_LAYER TransformerBlocks. Each TransformerBlock contains:
   - **Multi-Head Attention (sa):** This layer allows the model to attend to different parts of the input sequence simultaneously. It uses N_HEAD parallel CausalSelfAttention heads. Each head has a size of D_MODEL // N_HEAD. A causal mask is applied to prevent attending to future tokens.
   - **Feed-Forward Network (ffwd):** A simple two-layer linear network with a ReLU activation.
   - **Layer Normalization (ln1, ln2):** Applied before the attention and feed-forward layers to help stabilize training.
4. **Final Layer Normalization (ln_f):** A layer normalization applied after the transformer blocks.
5. **Language Model Head (lm_head):** A linear layer that maps the output of the transformer blocks to the vocabulary size (VOCAB_SIZE) to produce logits for the next token prediction.


The key parameters defining the model's size and training configuration are:
•       **B (Batch Size):** 16 - The number of sequences processed in parallel during training.
•       **T (Sequence Length):** 64 - The maximum length of input sequences the model can handle.
•       **D_MODEL (Embedding Dimension):** 128 - The size of the token and positional embeddings, and the hidden size throughout the transformer blocks.
•       **N_HEAD (Number of Attention Heads):** 4 - The number of parallel attention heads in the MultiHeadAttention layer.
•       **N_LAYER (Number of Transformer Layers):** 2 - The number of TransformerBlocks in the model.
•       **LEARNING_RATE:** 5e-4 - The learning rate for the AdamW optimizer.
•       **NUM_EPOCHS:** 5 - The number of times the training loop will iterate over the entire dataset.
•       **VOCAB_SIZE:** 30522 - The size of the vocabulary, determined by the tokenizer (bert-base-uncased).
In essence, the MiniGPT model processes input sequences by embedding tokens and their positions, then applying multiple layers of self-attention and feed-forward networks to learn relationships between tokens and predict the next token in the sequence.

II.     Dataset details

For dataset, we use the dataset collected and processed from Assignment 1, Wikipedia English 20231101.en, the first 150k rows, about 1.14G data.

III.    Training setup and hyperparameter experiments

**Training Setup:**
1.      **Device:** The code attempts to use a GPU (cuda) if available, otherwise it falls back to the CPU (cpu).
2.      **Dataset:** A custom PyTorch Dataset called NextTokenPredictionDataset is used to load data from a processed dataset stored locally.
3.      **DataLoader:** A DataLoader is used to iterate over the dataset in batches. A custom_collate_fn is used to handle padding of sequences to a maximum length of T + 1. The input X is the sequence up to the second to last token, and the target Y is the sequence shifted by one token to the right.
4.      **Optimizer:** The AdamW optimizer is used for training.
5.      **Loss Function:** Cross-entropy loss is used for the next token prediction task.
6.      **Metrics:** The code tracks and reports the average training loss and perplexity at specified intervals. Perplexity is calculated as the exponentiation of the loss.
7.      **Checkpointing:** The model's state dictionary is saved to a file (mini_gpt_checkpoint.pt) after each epoch.
8.      **Metric Saving:** Training losses and perplexities are saved to a .npz file (training_metrics.npz) for later visualization.
**Hyperparameters:**
The following hyperparameters are defined and used in the training process:
•       **B (Batch Size):** 16
•       **T (Sequence Length):** 64
•       **D_MODEL (Embedding Dimension):** 128
•       **N_HEAD (Number of Attention Heads):** 4
•       **N_LAYER (Number of Transformer Layers):** 2
•       **LEARNING_RATE:** 5e-4
•       **NUM_EPOCHS:** 5
•       **VOCAB_SIZE:** 30522 (derived from the tokenizer)
•       **EVAL_INTERVAL:** 100 (Although defined, evaluation is not explicitly implemented in the provided training loop; this variable seems to be a placeholder or for future use).
•       **LOG_INTERVAL:** 10 - The frequency (in steps) at which training loss and perplexity are logged.
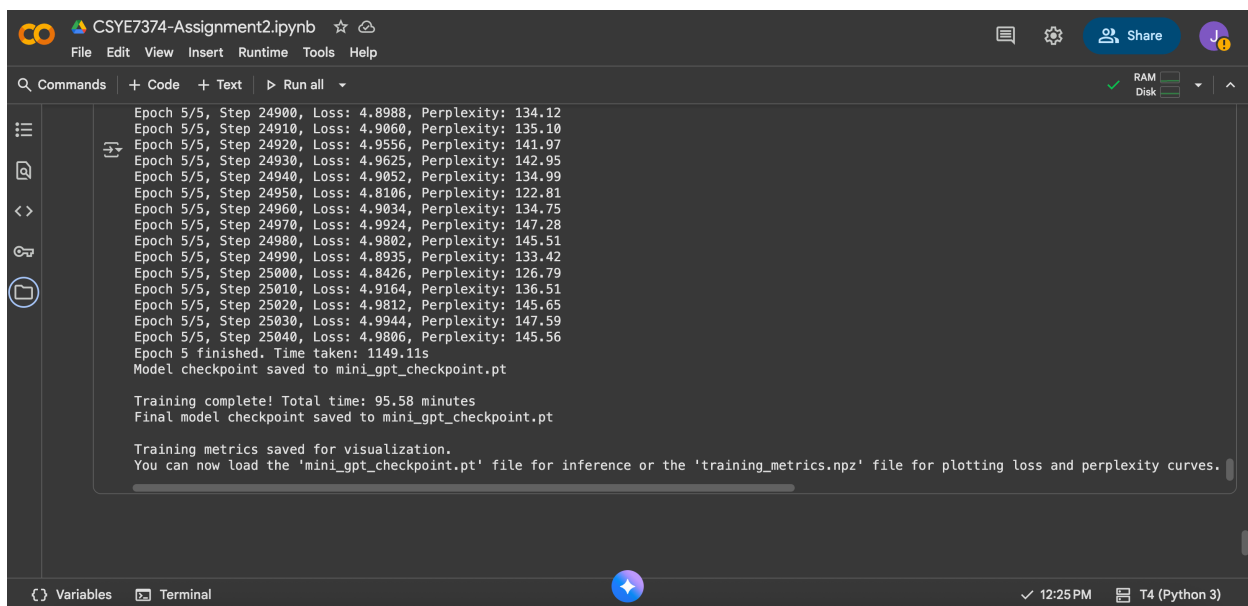•       **CHECKPOINT_PATH:** "mini_gpt_checkpoint.pt" - The file path for saving model checkpoints.
The training loop iterates for a fixed number of epochs, processing data in batches, calculating the loss, performing backpropagation, and updating the model's weights using the optimizer. Training progress is logged periodically, and checkpoints are saved after each epoch.

## IV.    Observations and challenges

I deliberately used both GPU on Colab and CPU on my own MacBook to train. The difference is striking. Despite the powerful M1Max chip, it still took 800 minutes to train locally, while it took only less than two hours to train on GPU, thanks to its significant parallelism.

Also, since the vectors are initiated randomly each time, the training outcome is different for each session. Notice that in the following images, the loss and perplexity are lower on Colab compared to local computer.

On Colab



On vscode