

# 1. Model Choice:

## Why did you choose the Qwen 2.5 model?

- **Qwen 2.5 Model Overview:**
  - In this section, you should describe the Qwen 2.5 model, explaining its origin, the specific version (in this case, "Qwen2.5-3B-Instruct"), and why you selected it over other models. You can mention that it is based on a strong architecture like GPT (Generative Pre-trained Transformer) which has been fine-tuned for instruction-based tasks. You can also explain its potential advantages such as:
    - **Large-scale pretrained language model:** It has been trained on vast amounts of data and can handle a variety of natural language understanding and generation tasks.
    - **Instruct-tuned model:** It's specialized for instruction-following tasks, making it suitable for your specific use case where the model is required to answer questions or follow instructions.
    - **Size and capability:** You may also mention the model's size, in this case, the 3 billion parameter variant, as it provides a balance between performance and resource constraints.

## Why fine-tune the Qwen 2.5 model?

- **Task Specificity:** Pre-trained models like Qwen 2.5 are often generalized for many tasks, but for specific applications (e.g., specialized Q&A or domain-specific language), fine-tuning can improve the model's performance in those areas.
- **Transfer Learning Benefits:** Fine-tuning allows the model to leverage the knowledge it has learned in pre-training while adapting to the unique requirements of your specific dataset or problem.
- **Performance Gains:** Fine-tuning on your custom dataset may lead to improved accuracy, specificity, and adaptability compared to using the pre-trained model directly.

## 2. Data Preprocessing:

### How was the data processed?

- **Dataset Overview:**
  - Start by describing the dataset you used for training (or the type of data if you have created your dataset). For instance, if your task is question-answering, describe how questions and answers are structured and the source of this data (e.g., custom data, public datasets).
- **Preprocessing Techniques:**
  - **Tokenization:** Explain why you tokenized the data using the tokenizer from the same model family. Tokenization is the first step in preparing text for input into a transformer model. It involves splitting text into smaller units (tokens), which the model can understand. You may mention:
    - The tokenizer used (e.g., AutoTokenizer).
    - How tokenization is important to ensure consistent input format for the model.
  - **Truncation and Padding:** Explain that inputs to the model need to have a consistent length, so you used truncation and padding techniques.
    - **Truncation:** You truncated longer sequences to fit a predefined `max_length`. This ensures that the model processes only the relevant information and does not run out of memory with excessively long texts.
    - **Padding:** For shorter sequences, padding ensures that all inputs are of the same length. It avoids errors and helps the model handle batches efficiently.
    - Mention the padding strategy, such as `padding="max_length"`, which pads to the maximum length of the model, and the reason why truncation was necessary.
  - **Why these preprocessing steps were chosen:**
    - Tokenization ensures the model can interpret the data correctly.
    - Truncation ensures the inputs don't exceed the model's capacity.
    - Padding makes sure all sequences are processed in batches without errors.

### 3. Training Setup:

#### What hyperparameters were used and why?

- **Batch Size:**
  - Explain the choice of batch size, such as `per_device_train_batch_size=2`. Larger batch sizes provide more training examples in parallel but require more memory. Smaller batches can reduce memory usage and improve generalization. In this case, you might have chosen 2 for a balance between computational efficiency and stability.
- **Number of Epochs:**
  - Explain the number of epochs you chose (e.g., 3 epochs). Epochs define how many times the model sees the entire training data. A higher number of epochs might help the model learn better, but it can also cause overfitting. You can justify the number based on the size of your data or experimentation.
- **Learning Rate:**
  - Describe your choice of learning rate (e.g.,  $5e-5$ ). The learning rate controls how much the model's weights are updated during each step. A too-high learning rate might make the model diverge, and a too-low learning rate could lead to slow learning. You might have chosen a lower learning rate for fine-tuning to prevent drastic changes to the pretrained weights.
- **Evaluation Strategy and Logging:**
  - **Logging:** Discuss your logging frequency (`logging_steps=10`), and why you track progress during training.
  - **Save Steps and Evaluation Strategy:** Explain the need to evaluate and save the model at certain intervals. For example, saving every 100 steps (`save_steps=100`) allows you to store intermediate models and avoid losing progress in case of issues.
- **Use of Trainer:**
  - The Trainer class simplifies training, so explain that using it allows you to leverage built-in features like logging, saving, and evaluation without writing a lot of custom training loop code.

## 4. Quantization:

### Why use 4-bit quantization?

- **Overview of Quantization:**
  - Quantization reduces the number of bits required to represent the model weights, allowing it to fit into smaller memory spaces while maintaining a reasonable level of accuracy.
- **Impact on Performance and Size:**
  - **Performance:** 4-bit quantization is a tradeoff between model size and accuracy. It significantly reduces the memory footprint (in this case, by about 4x) and enables faster inference, especially on specialized hardware like GPUs and TPUs.
  - **Model Size:** By reducing the precision of the model weights from 16 or 32 bits to 4 bits, you reduce the model size, which can be helpful when deploying models on resource-constrained devices or for large-scale inference.
  - **Tradeoffs:** You might explain that 4-bit quantization comes at the cost of some model performance (slight reduction in accuracy), but the benefits in terms of deployment efficiency and memory usage outweigh these concerns, especially for real-time or edge applications.
- **Why use 4-bit and not higher or lower?**
  - **Balance:** 4-bit precision strikes a good balance between speed, memory efficiency, and maintaining acceptable performance.

## 5. Evaluation Metrics:

### What metrics were used to evaluate the model?

- **Accuracy:**
  - If the task involves classification or question-answering, accuracy is often a key metric. You should explain how you defined accuracy in the context of your task (e.g., percentage of correct answers).
- **BLEU or ROUGE Score (if applicable):**
  - For generative tasks like text generation or summarization, BLEU (Bilingual Evaluation Understudy) or ROUGE (Recall-Oriented Understudy for Gisting Evaluation) scores are often used to measure the quality of generated text compared to human-written text.
- **Perplexity (if used):**
  - Perplexity is a common metric for language modeling tasks. It measures how well the model predicts a sample. Lower perplexity means better performance.
- **F1 Score or Precision/Recall (if used):**
  - For tasks like information extraction or question answering, the F1 score, precision, and recall might be more appropriate as they measure how well the model extracts relevant information.
- **Custom Evaluation Framework:**
  - If you implemented a custom evaluation framework, explain how you structured it. For example, you could have used a combination of automatic metrics like BLEU, human evaluation for more subjective tasks, or domain-specific evaluation criteria.