# Question Answering — Task-Specific Models vs. LLMs

**Zongyan Yao**
Department of Electrical and Computer Engineering
University of Toronto
zongyan.yao@mail.utoronto.ca


**Zhengyang Li**
Department of Electrical and Computer Engineering
University of Toronto
zhengyang.li@mail.utoronto.ca


**Qiwen Lin**
Department of Electrical and Computer Engineering
University of Toronto
qw.lin@mail.utoronto.ca

## Abstract

Large Language Models (LLMs) have shown strong performance on open-domain question answering, but they often lack the domain-specific precision required for specialized tasks. This project investigates whether fine-tuning smaller and mid-sized pretrained transformer models can outperform general-purpose LLMs on a narrow recipe-focused benchmark. Using the Recipe-MPR dataset of 500 multi-perspective cooking questions, we fine-tune several encoder-based models (BERT, DistilBERT, RoBERTa [1, 5, 3]) and decoder-based models (Llama [6] and Qwen2.5–7B [8]) with parameter-efficient methods such as LoRA [2]. Our goal is to exceed a 65% accuracy threshold and to study trade-offs between accuracy, efficiency, and model capacity.

Initial results show that all fine-tuned encoder models exceed the target, with BERT-large achieving 91.4% accuracy and DistilBERT providing strong performance at lower computational cost. The LoRA-tuned Qwen2.5–7B model achieves 100% accuracy, and fine-tuned Llama variants reach up to 84%, while GPT-3 embeddings [4] perform significantly worse. These findings highlight the benefits of domain adaptation and provide a clear basis for deeper analysis in the final report.

## 1   Introduction

Large Language Models (LLMs) such as BERT and GPT have demonstrated impressive performance on open-domain natural language tasks [1, 7], but their general-purpose nature can limit reliability in narrow domains. In the recipe domain, small differences in wording can affect ingredient substitutions, cooking steps, or safety-related advice, and general LLMs may not always capture these nuances consistently.

To address this, we study whether fine-tuning smaller or mid-sized pretrained transformer models on a recipe-specific dataset can outperform a general-purpose LLM baseline. We use the Recipe-MPR dataset, which consists of 500 user queries with five candidate answers per query, covering multiple reasoning types (analogical, specific, commonsense, temporal, and negated). Our aim in this progress

stage is to clearly specify the problem, describe our model and training design, report the main results obtained so far, and outline the remaining milestones rather than to deliver a polished final narrative.

## 2 Preliminaries and Problem Formulation

### 2.1 Problem Definition

We formulate the task as a multiple-choice question answering problem over recipes. For each example, we are given:

- a query $q$ (a user question about recipes or cooking), and
- a set of five candidate answers $A = \{a_1, a_2, a_3, a_4, a_5\}$.

The goal is to learn a function

$$f_\theta(q, A) \rightarrow y, \quad y \in \{1, 2, 3, 4, 5\},$$

that predicts the index of the correct answer. We model this as a five-class classification problem using transformer-based architectures [7].

### 2.2 Design Components and Objectives

To make the project structure clear, we explicitly summarize our main design components:

- **Dataset:** Recipe-MPR (500 examples), split into 80/10/10 for train/validation/test.
- **Models:**
  - Encoder models: BERT-base, BERT-large [1], DistilBERT [5], RoBERTa-base [3].
  - Decoder models: Llama-3.2-1B, Llama-3.2-3B [6] and Qwen2.5–7B [8].
- **Training Strategy:**
  - Full fine-tuning for encoder models.
  - LoRA-based parameter-efficient fine-tuning for large decoder models [2].
  - Comparison with zero-shot LLMs and a GPT-3 embedding baseline [4].
- **Metrics:** Accuracy (primary) and macro F1-score (planned for final report).
- **Targets:**
  - Exceed a 65% accuracy baseline.
  - Analyze trade-offs between accuracy, model size, training time, and inference cost.

### 2.3 Relevant Background Concepts

**Transformer Architecture.** All models in this work are based on the transformer architecture [7], which uses multi-head self-attention, feed-forward layers, positional encodings, and layer normalization to model contextual dependencies.

**Fine-Tuning and LoRA.** Encoder models (e.g., BERT and RoBERTa) are fine-tuned end-to-end [1, 3], whereas large decoder LLMs (Llama, Qwen) are adapted using LoRA [2], which adds low-rank trainable adapters to reduce memory and compute while preserving most of the pretrained weights.

**Embedding Baseline.** We also experiment with GPT-3 embeddings [4], where queries and answers are mapped to vector representations, and the answer is chosen based on similarity in embedding space. This provides a non–fine-tuned baseline for comparison against explicit model adaptation.

## 3 Solution via Deep Learning

### 3.1 Dataset Preparation

We preprocess the Recipe-MPR dataset by:

- normalizing text (lowercasing and basic cleanup),
- tokenizing using each model's tokenizer,
- truncating or padding sequences to a fixed maximum length,
- splitting into 80% training, 10% validation, 10% test.

This part of the pipeline is fully implemented and reused across all models to ensure fair comparison.

## 3.2 Model and Training Pipeline

We have implemented a unified training pipeline in PyTorch with HuggingFace Transformers for both encoder and decoder models.

**Encoder models.** BERT-base, BERT-large, DistilBERT, and RoBERTa-base [1, 5, 3] are fine-tuned by adding a classification head on top of the [CLS] representation and training with cross-entropy loss.

**Decoder models.** Llama-3.2-1B, Llama-3.2-3B [6] and Qwen2.5–7B [8] are adapted using LoRA [2] with 4-bit quantization to fit within 8–32GB GPUs. We frame the task as sequence-to-sequence or causal LM scoring over the candidate answers.

## 3.3 Current Progress and Preliminary Results

At this stage of the project, the following components are completed:

- All preprocessing, tokenization, and dataset splitting.
- Full fine-tuning of BERT-family models (BERT-base, BERT-large, DistilBERT, RoBERTa-base).
- LoRA-based fine-tuning of Qwen2.5–7B and Llama-3.2 models.
- Implementation and evaluation of a GPT-3 embedding baseline [4].
- Test-set evaluations for all models and a cross-model comparison.

Preliminary results (on our current test split) show:

The Qwen2.5–7B model demonstrates exceptional specialization on the Recipe-MPR task. After applying LoRA fine-tuning and training on a 32GB GPU, the model achieved a perfect score with no mistakes across 50 test examples. Table 1 compares the fine-tuned Qwen model with the base Qwen model and the top-performing BERT variants. Despite Qwen's much larger parameter count, LoRA fine-tuning remains computationally tractable and yields perfect task performance.

Table 1: Performance Comparison: Qwen vs. BERT Models

| Model | Parameters | Accuracy | Training Time | VRAM | Architecture |
|---|---|---|---|---|---|
| Fine-tuned Qwen | 7B | **100.00%** | ~15 min | 20 GB | Decoder-only |
| Base Qwen | 7B | 79.20% | 0 min | 20 GB | Decoder-only |
| BERT-large | 340M | 91.4% | ~93 sec | 16 GB | Encoder-only |
| DistilBERT | 66M | 82.4% | ~35 sec | 6 GB | Encoder-only |

Training on a 32GB GPU, all BERT-family models meeting the 65% threshold are summarized in Table 2.

We fine-tuned Llama-3.2-3B, Llama-3.2-1B with LoRA to do recipe recommendation. LoRA lets us train only a tiny part of the model .

Results are clear. The 3B LoRA model got 84.00% in ~14 min. The 1B LoRA model got 78.00% in about 7 min. Both beat the 65% goal. The 3B zero-shot model reached 73.00% (meets goal), while 1B zero-shot was 58.00%. The GPT-3 Embedding baseline was 54.55%. Fine-tuning gave +11 points over 3B zero-shot, and LLM reasoning beat the embedding method by +18.45 points, for a total +29.45 points over the baseline.

Table 2: Performance of BERT-family models on the Recipe-MPR dataset.

| Model | Accuracy | Above Goal? |
|---|---|---|
| BERT-large | **91.4%** | Yes |
| DistilBERT | 82.4% | Yes |
| BERT-base (aggressive) | 67.6% | Yes |
| BERT-base (standard) | 65.6% | Yes |
| RoBERTa-base (aggressive) | 65.8% | Yes |
| RoBERTa-base (standard) | 49.8% | No |
| DistilBERT (over-trained) | 36.8% | No |

The architecture of all Llama models are Decoder-only. All evaluated models with their accuracy and whether they meet the 65% threshold are summarized in Table 3.

Table 3: Performance of Llama models vs. GPT3 on the Recipe-MPR test set.

| Model | Parameters | Accuracy | Training Time | VRAM | Above Goal? |
|---|---|---|---|---|---|
| Fine-tuned Llama-3.2-3B | 3.2B | **84.00%** | 14 min | 6.4 GB | Yes |
| Base Llama-3.2-3B | 3.2B | 73.00% | 0 min | 6.4 GB | Yes |
| Fine-tuned Llama-1B | 1B | 78.00% | $\sim$7 min | 2.0 GB | Yes |
| Base Llama-1B | 1B | 58.00% | 0 min | 2.0 GB | No |
| GPT-3 Embedding | N/A | 54.55% | N/A | N/A | No |

All fine-tuned models exceed the 65% threshold, and several exceed 80%, indicating that our design choices are effective for this dataset.

### 3.4 Remaining Milestones

Although the core training and evaluation pipeline is complete, several important steps remain for the final report:

- Perform detailed error analysis across query types (analogical, specific, commonsense, temporal, negated).
- Compare qualitative outputs of models (e.g., typical failure cases for smaller vs. larger models).
- Produce visualizations (e.g., bar plots, confusion matrices, accuracy vs. parameters).
- Integrate all findings into a polished final report with a more formal discussion and conclusion.

From a project management perspective, we have completed the major technical components (data pipeline, training, and initial evaluation) and are now focused on analysis, comparison, and final documentation.

## References

[1] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. *NAACL*, 2019.

[2] Edward J Hu, Yelong Shen, Phillip Wallis, et al. Lora: Low-rank adaptation of large language models. *arXiv preprint arXiv:2106.09685*, 2021.

[3] Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. Roberta: A robustly optimized bert pretraining approach. *arXiv preprint arXiv:1907.11692*, 2019.

[4] OpenAI. Text and code embeddings. *OpenAI Technical Report*, 2022.

[5] Victor Sanh, Lysandre Debut, Julien Chaumond, and Thomas Wolf. Distilbert, a distilled version of bert: Smaller, faster, cheaper and lighter. *arXiv preprint arXiv:1910.01108*, 2019.

[6] Hugo Touvron et al. Llama: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971*, 2023.

[7] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in Neural Information Processing Systems*, 2017.

[8] An Yang, Baosong Yang, Junyang Lin, Xiaodong Chen, and Jingren Zhang. Qwen2 technical report. *arXiv preprint arXiv:2407.10671*, 2024.