

# LLM Fine-tuning Challenge: Enhancing Qwen 2.5-3B for AI Research QA

## Project Overview

### Technical Approach

Document Processing

Synthetic QA Generation

Fine-tuning Pipeline

Retrieval-Augmented Generation (RAG) System

Model Quantization

Model Inference

Evaluation Framework

### Implementation Details

Data Processing and Generation

Fine-tuning Configuration

RAG Implementation

Quantization Process

### Conclusion

This report documents our approach to enhancing Qwen 2.5-3B for answering questions about technical AI research papers and documents.

## Project Overview

We implemented a comprehensive solution to create a specialized QA model for technical AI research domains. Our approach focused on five key components:

1. Document Processing: Converting technical research papers into structured formats
2. Synthetic QA Generation: Creating high-quality training data from processed documents

3. Model Fine-tuning: Optimizing Qwen 2.5-3B using efficient parameter-efficient methods
4. Retrieval-Augmented Generation: Enhancing response quality with relevant context
5. Evaluation Framework: Assessing performance improvements with multiple metrics

## Technical Approach

### Document Processing

Our document processing pipeline extracts and structures information from technical markdown documents. The system handles specialized formatting common in research papers and segments text into meaningful chunks with optimal overlap to preserve context between segments. This approach ensures that complex technical concepts remain coherent and properly linked during the QA generation phase.

### Synthetic QA Generation

To create high-quality training data, we developed a synthetic QA generation system using larger language models as teachers. The system employs carefully crafted instruction templates optimized for technical content, focusing on generating diverse question types including factual, analytical, comparative, and hypothetical questions. This method produces training and validation datasets that reflect the complexity and specificity of AI research questions.

### Fine-tuning Pipeline

We implemented Quantized Low-Rank Adaptation (QLoRA) for efficient fine-tuning of the Qwen 2.5-3B model. This approach dramatically reduces memory requirements while preserving model

quality. Our hyperparameter selection was optimized specifically for technical domain adaptation, with particular attention to learning rate scheduling, adapter rank, and quantization settings. Training was monitored using Weights & Biases integration for comprehensive tracking of model performance.

## **Retrieval-Augmented Generation (RAG) System**

Our RAG implementation uses a FAISS-based vector store for semantic document retrieval. The system employs domain-optimized embeddings specifically tuned for technical content. The retrieval component prioritizes both semantic similarity and information density when selecting context for the model, ensuring relevant technical details are available during generation while avoiding information overload.

## **Model Quantization**

To meet the competition requirements, we developed a robust quantization pipeline that converts the fine-tuned model to an efficient 4-bit GGUF format. Our quantization approach preserves model capabilities while significantly reducing size and enabling deployment on resource-constrained environments. The process includes calibration steps to minimize accuracy loss during quantization, particularly for technical terminology.

## **Model Inference**

The inference system combines the quantized model with the RAG architecture for optimal performance. It manages context length constraints effectively by using intelligent chunking and summarization techniques. The system can operate in both RAG-enabled and standalone modes, allowing for flexibility in deployment scenarios with different computational resources.

## **Evaluation Framework**

Our comprehensive evaluation framework uses multiple metrics to assess model performance. These include both automatic metrics (ROUGE, BLEU) and customized technical accuracy assessments. The evaluation process compares the model performance with and without RAG enhancement, providing insights into the contribution of each component to overall system quality.

## **Implementation Details**

### **Data Processing and Generation**

We processed a collection of technical AI research papers focused on distributed systems and performance optimization. The document processing pipeline preserved formatting and technical terminology while splitting content into semantically meaningful chunks. Using these chunks, we generated approximately 2,500 synthetic QA pairs with diverse question types and comprehensive answers, which were split into training (80%), validation (10%), and test (10%) sets.

### **Fine-tuning Configuration**

Our fine-tuning process used the Qwen 2.5-3B-Instruct model with QLoRA adapters targeting attention layers. We set the LoRA rank to 16 with an alpha of 32, finding this balance optimal for technical domain adaptation. Training used mixed precision (fp16) with a learning rate of  $2e-4$ , batch size of 4, and gradient accumulation over 4 steps. The training ran for 1,000 steps with evaluation checkpoints every 200 steps.

### **RAG Implementation**

The RAG system uses the BAAI/bge-small-en-v1.5 embedding model for efficient document indexing. Document chunks are stored in a FAISS vector database for rapid similarity searching. During inference, the system retrieves the 3 most relevant document

chunks and formats them into a structured context prompt that guides the model's responses, significantly improving accuracy on technical questions.

## Quantization Process

The 4-bit quantization process was implemented using llama.cpp's conversion tools, with calibration on a representative dataset sample to ensure minimal quality degradation. The resulting GGUF file is approximately 1.9GB, representing an 80% reduction from the original model size while maintaining core capabilities for technical question answering.

## Conclusion

Our implementation successfully enhances Qwen 2.5-3B for the specialized task of answering questions about technical AI research. The combination of efficient fine-tuning techniques, synthetic data generation, and retrieval augmentation creates a system well-suited to the challenges of interpreting and responding to complex technical queries. The 4-bit quantized model provides an accessible deployment option while maintaining strong performance on technical AI domain questions.

The approach demonstrates that relatively small models (3B parameters) can achieve strong performance in specialized domains when properly optimized and augmented with retrieval systems. Future work could focus on expanding the document corpus and refining the RAG system for handling conflicting information in research papers.