**Research Paper**

# Domain-Specific Fine-Tuning of Large Language Models for Financial Applications

## Ziran Wang

*Abstract*
*Recent advances in large language models (LLMs) have created new opportunities for artificial intelligence in finance, a field that depends heavily on specialized language and precise information processing. However, directly applying general-purpose LLMs to financial tasks is often limited by domain-specific challenges such as technical terminology, dynamic market contexts, regulatory constraints, and security concerns. To bridge this gap, domain-specific fine-tuning has emerged as a practical and effective approach for adapting LLMs to the financial sector.*

*This paper investigates methods for fine-tuning LLMs on financial data, highlighting both technical considerations and practical applications. We review existing adaptation strategies—ranging from zero-shot and few-shot learning to domain-adaptive pretraining—but emphasize fine-tuning as the primary pathway for enhancing performance in financial natural language processing tasks. Building on this foundation, we outlined a framework for model selection, dataset construction, preprocessing, and model optimization tailored to financial contexts.*

*Through representative use cases, including financial sentiment analysis, financial text summarization and financial question answering, we demonstrate the potential of fine-tuned models to advance decision-making and operational efficiency in finance. Finally, we discuss current limitations, such as data scarcity and regulatory risks, and propose directions for future research. By doing so, this work provides both researchers and practitioners with a roadmap for responsibly and effectively adopting fine-tuned LLMs in financial applications.*

*Keywords: Large Language Models, Fine-Tuning, Domain Adaptation, Financial Natural Language Processing, Financial Applications*

## I.    Introduction

Recent advances in artificial intelligence, particularly in natural language processing (NLP), have driven the rapid development of large language models (LLMs) such as ChatGPT. These models demonstrate remarkable capabilities in understanding, generating, and reasoning about human language. Given that the financial industry is inherently text-driven—relying heavily on reports, news, customer interactions, and regulatory documents—LLMs hold significant promise for transforming financial applications ranging from trading and risk modeling to customer service and compliance.

While general-purpose LLMs have shown strong performance across diverse tasks, their direct application to finance remains limited by domain-specific challenges. Financial language often involves specialized terminology, rapidly changing market contexts, and strict requirements for accuracy, security, and regulatory compliance. As a result, fine-tuning (FT) emerges as a critical approach for adapting LLMs to the financial domain. By tailoring pretrained models with domain-relevant data, fine tuning can enhance their ability to capture financial semantics, improve task-specific performance, and reduce hallucinations or misinterpretations.

In this paper, we focus on domain-specific fine-tuning of LLMs for financial applications. We first survey existing approaches that adapt LLMs to finance, including domain-adaptive pretraining and task-oriented fine-tuning. We then present practical guidance on the fine-tuning process, covering model selection, dataset selection, preprocessing, and model optimization. Beyond methodology, we illustrate the potential of fine-tuned financial LLMs through representative use cases such as Financial Sentiment Analysis, Financial Text Summarization, Financial Question Answering.

Our goal is to provide both researchers and practitioners with a structured understanding of how fine-tuning can bridge the gap between general-purpose LLMs and the specialized needs of financial applications. By highlighting opportunities, challenges, and future directions, we aim to contribute to the responsible and effective adoption of LLMs in finance.

# II. Related Work

## 2.1. Overview of Large Language Models (LLMs)

Large language models (LLMs) have become the cornerstone of modern natural language processing (NLP), demonstrating unprecedented capabilities in text understanding, generation, and reasoning. Models such as GPT series, BERT, T5, and LLaMA exemplify different architectural innovations and training paradigms, from autoregressive transformers to masked language models and sequence-to-sequence designs. Their effectiveness stems largely from pretraining on massive corpora and leveraging billions of parameters, which allow them to capture linguistic regularities and contextual dependencies at scale.

Several key trends have marked the evolution of LLMs. First, scaling laws suggest that performance improves predictably with increased model size, dataset size, and compute power, leading to the emergence of foundation models that serve as general-purpose backbones across domains. Second, instruction-tuning and reinforcement learning with human feedback (RLHF) have further aligned LLMs with human intent, enabling models like ChatGPT to perform well on open-ended tasks. Third, open-source initiatives (e.g., BLOOM, Falcon, MPT, LLaMA 2/3) have democratized access to LLMs, fostering rapid experimentation in academia and industry.

Despite these advances, LLMs trained on general web-scale corpora often fall short in highly specialized domains such as finance, law, or medicine. Their lack of domain-specific vocabulary, sensitivity to hallucinations, and challenges in handling up-to-date or proprietary data underscore the need for adaptation techniques. In this context, fine-tuning, retrieval-augmented generation (RAG), and agent-based frameworks have been widely explored as mechanisms to extend the utility of LLMs beyond their generic training.

## 2.2. Related Techniques: Fine-Tuning, Retrieval-Augmented Generation, and Agents 2.2.1. Fine-Tuning

Fine-tuning remains the most direct and effective strategy for adapting LLMs to domain-specific tasks. Traditional full-parameter fine-tuning updates all model parameters using labeled data, which can significantly improve performance but requires considerable computational resources. To address efficiency concerns, parameter-efficient fine-tuning (PEFT) methods such as LoRA (Low-Rank Adaptation), adapters, and prompt-tuning have been introduced. These methods adjust only a small subset of parameters, reducing training cost while preserving strong performance.

In the financial domain, fine-tuning has been applied to tasks including sentiment analysis of market news, extraction of entities from reports, risk factor modeling, and dialogue systems for customer support. Compared to zero-shot and few-shot prompting, fine-tuning offers more consistent accuracy and robustness, particularly when dealing with jargon-heavy or sensitive financial language. Moreover, fine-tuning facilitates alignment with regulatory requirements by constraining model behavior through supervised training.

## 2.2.2. Retrieval-Augmented Generation (RAG)

While fine-tuning adapts model parameters, retrieval-augmented generation (RAG) augments models with access to external knowledge bases. In RAG systems, an LLM retrieves relevant documents from a database or vector store and integrates them into its generation process. This technique addresses two significant limitations of LLMs: hallucination and knowledge staleness. By grounding responses in up-to-date financial documents, regulatory filings, or proprietary databases, RAG enhances factual accuracy and enables real-time information retrieval.

In finance, RAG-based systems are increasingly used for tasks such as compliance monitoring, earnings report summarization, and personalized investment advisory. For instance, a model may retrieve the latest SEC filings or analyst reports before generating insights, ensuring that its outputs are both accurate and contextually relevant. However, designing efficient retrieval pipelines and ensuring data confidentiality remain open challenges.

## 2.2.3. Agent-based Frameworks

A newer paradigm involves framing LLMs as agents capable of reasoning, planning, and interacting with external tools. Agent frameworks, such as LangChain and AutoGPT, orchestrate multiple components: prompting strategies, memory modules, tool APIs (e.g., databases, trading platforms), and iterative feedback loops. This approach transforms LLMs from passive text generators into active problem-solvers that can perform multi-step tasks.

---

In the financial sector, LLM-based agents show promise in automating workflows such as due diligence, portfolio optimization, and regulatory auditing. For example, an agent may parse financial documents, query databases, run simulations, and generate risk assessments in a single pipeline. While the potential is high, the complexity of orchestrating reliable agents, coupled with concerns around interpretability and error propagation, highlights the need for rigorous evaluation and control.

2.3. Applications of LLMs in Finance

The financial industry provides a fertile ground for applying LLMs, given its reliance on textual data, structured reports, and customer interactions. Early applications have demonstrated the versatility of LLMs across several key areas:

1) Market Analysis and Forecasting: LLMs are employed to analyze financial news, earnings calls, and analyst reports to predict market sentiment and stock price movements. Fine-tuned models on domain-specific corpora have been shown to outperform generic ones in tasks like sentiment classification and event-driven prediction.

2) Risk Management and Compliance: Regulatory compliance requires continuous monitoring of legal documents, guidelines, and transaction records. LLMs can assist by automating information extraction, anomaly detection, and compliance reporting. Integration with RAG pipelines ensures alignment with the latest regulatory updates.

3) Document Processing and Knowledge Management: Financial institutions handle vast amounts of unstructured text, from contracts to annual reports. LLMs streamline tasks such as summarization, entity recognition, and document classification, improving efficiency in back-office operations.

4) Customer Service and Advisory: Conversational agents powered by fine-tuned LLMs provide personalized customer support, answer queries about financial products, and guide investment decisions. Fine-tuning customer interaction logs further improves reliability and tone adaptation.

5) Algorithmic Trading and Decision Support: Although still experimental, LLMs are being explored in trading strategies by integrating textual signals with quantitative models. For instance, sentiment from news articles and social media posts can be transformed into features for predictive modeling.

6) Fraud Detection and Security: LLMs enhance fraud detection by analyzing transaction descriptions, communication logs, and suspicious reporting patterns. Combined with structured data, these models contribute to a multi-modal view of risk.

Overall, the applications highlight the transformative potential of LLMs in finance, but they also reveal critical limitations. Issues such as data privacy, explainability, domain-specific accuracy, and the risk of over-reliance remain barriers to widespread deployment. This underscores the importance of approaches like fine-tuning, RAG, and agent-based systems, which together offer pathways to more reliable and responsible use of LLMs in financial contexts.

# III.     Method

3.1. Model Selection

Model selection plays a crucial role in determining the overall performance of fine-tuning in the financial domain. To identify a suitable base model, we systematically reviewed recent benchmark results across multiple categories, including knowledge, reasoning, coding, alignment, agent capabilities, and multilingualism. As shown in Table 1, Qwen3-4B-Instruct-2507 demonstrates competitive performance across benchmarks, notably outperforming other 4B-scale models in reasoning (ZebraLogic: 80.2), knowledge (MMLU-Pro: 69.6), and alignment tasks (Creative Writing: 83.5). Considering both model performance and available computational resources, we selected Qwen3-4B-Instruct-2507 as the base model for fine-tuning. Its balance between strong instruction-following ability and manageable model size (4B parameters) makes it particularly suitable for domain adaptation under resource constraints. All experiments in this paper were conducted on a single NVIDIA RTX 3090 GPU (24GB VRAM), ensuring feasibility for academic and enterprise research settings without requiring large-scale GPU clusters.

<Table 1> Comparing Benchmark Scores of Candidate Models

| Category | Benchmark | GPT-4.1-nano (2025-04-14) | Qwen3-30B-A3B | Qwen3-4B | Qwen3-4B-Instruct-2507 |
|---|---|---|---|---|---|
| Knowledge | MMLU-Pro | 62.8 | 69.1 | 58.0 | 69.6 |
| | MMLU-Redux | 80.2 | 84.1 | 77.3 | 84.2 |
| | GPQA | 50.3 | 54.8 | 41.7 | 62.0 |
| | SuperGPQA | 32.2 | 42.2 | 32.0 | 42.8 |
| Reasoning | AIME25 | 22.7 | 21.6 | 19.1 | 47.4 |
| | HMMT25 | 9.7 | 12.0 | 12.1 | 31.0 |
| | ZebraLogic | 14.8 | 33.2 | 35.2 | 80.2 |
| | LiveBench 20241125 | 41.5 | 59.4 | 48.4 | 63.0 |
| Coding | LiveCodeBench v6 (25.02–25.05) | 31.5 | 29.0 | 26.4 | 35.1 |
| | MultiPL-E | 76.3 | 74.6 | 66.6 | 76.8 |
| | Aider-Polyglot | 9.8 | 24.4 | 13.8 | 12.9 |
| Alignment | IFEval | 74.5 | 83.7 | 81.2 | 83.4 |
| | Arena-Hard v2* | 15.9 | 24.8 | 9.5 | 43.4 |
| | Creative Writing v3 | 72.7 | 68.1 | 53.6 | 83.5 |
| | WritingBench | 66.9 | 72.2 | 68.5 | 83.4 |
| Agent | BFCL-v3 | 53.0 | 58.6 | 57.6 | 61.9 |
| | TAU1-Retail | 23.5 | 38.3 | 24.3 | 48.7 |
| | TAU1-Airline | 14.0 | 18.0 | 16.0 | 32.0 |
| | TAU2-Retail | – | 31.6 | 28.1 | 40.4 |
| | TAU2-Airline | – | 18.0 | 12.0 | 24.0 |
| | TAU2-Telecom | – | 18.4 | 17.5 | 13.2 |
| Multilingualism | MultiIF | 60.7 | 70.8 | 61.3 | 69.0 |
| | MMLU-ProX | 56.2 | 65.1 | 49.6 | 61.6 |
| | INCLUDE | 58.6 | 67.8 | 53.8 | 60.1 |
| | PolyMATH | 15.6 | 23.3 | 16.6 | 31.1 |

## 3.2. Dataset Processing

### 3.2.1. Dataset Description
For domain-specific fine-tuning, we adopt the Financial Alpaca LLaMA dataset from HuggingFace (madanarnav/financial_alpaca_llama), which is designed to support instruction-tuning in financial applications. The dataset follows an instruction–input–output format, enabling supervised fine-tuning that aligns model behavior with human-like financial reasoning and response generation.

The dataset contains 68,912 training examples with a total size of approximately 76 MB. Each sample includes five key fields:

id: Unique identifier for each entry.

instruction: Task description provided to the model.

input: Additional context or data associated with the instruction.

output: Target response used as ground-truth for supervision.

text: Consolidated version of the instruction–input–output triplet.

This structure makes the dataset particularly well-suited for instruction tuning, as it allows the model to learn mappings from natural-language financial tasks (e.g., summarize an earnings report or analyze market sentiment) to domain-specific responses.

### 3.2.2. Preprocessing
Before training, we apply several preprocessing steps to ensure data quality and compatibility with the fine-tuning pipeline: Deduplication – remove redundant instructions and outputs to prevent overfitting.

Formatting – restructure samples into instruction–response pairs compatible with the Qwen3-4B-Instruct-2507 model.

Train–validation split – partition the dataset to enable performance monitoring and mitigate overfitting during fine-tuning.

Given its scale and domain specificity, the Financial Alpaca LLaMA dataset provides a robust foundation for adapting general-purpose LLMs into finance-focused assistants, thereby improving accuracy in tasks such as report summarization, sentiment analysis, compliance monitoring, and automated advisory.

<Figure 1> Processed Json File

```
[
  {
    "instruction": "For a car, what scams can be plotted with 0% financing vs rebate?",
    "input": "",
    "output": "The car deal makes money 3 ways. If you pay in one lump payment. If the payment is greater than what they paid for
the car, plus their expenses, they make a profit. They loan you the money. You make payments over months or years, if the total
amount you pay is greater than what they paid for the car, plus their expenses, plus their finance expenses they make money. Of
course the money takes years to come in, or they sell your loan to another business to get the money faster but in a smaller
amount. You trade in a car and they sell it at a profit. Of course that new transaction could be a lump sum or a loan on the used
car... They or course make money if you bring the car back for maintenance, or you buy lots of expensive dealer options. Some
dealers wave two deals in front of you: get a 0% interest loan. These tend to be shorter 12 months vs 36,48,60 or even 72 months.
The shorter length makes it harder for many to afford. If you can't swing the 12 large payments they offer you at x% loan for y
years that keeps the payments in your budget. pay cash and get a rebate. If you take the rebate you can't get the 0% loan. If you
take the 0% loan you can't get the rebate. The price you negotiate minus the rebate is enough to make a profit. The key is not
letting them know which offer you are interested in. Don't even mention a trade in until the price of the new car has been
finalized. Otherwise they will adjust the price, rebate, interest rate, length of loan,  and trade-in value to maximize their
profit. The suggestion of running the numbers through a spreadsheet is a good one. If you get a loan for 2% from your bank/credit
union for 3 years and the rebate from the dealer, it will cost less in total than the 0% loan from the dealer. The key is to get
the loan approved by the bank/credit union before meeting with the dealer. The money from the bank looks like cash to the dealer."
  },
  {
    "instruction": "Why does it matter if a Central Bank has a negative rather than 0% interest rate?",
    "input": "",
    "output": "That is kind of the point, one of the hopes is that it incentivizes banks to stop storing money and start injecting
it into the economy themselves. Compared to the European Central Bank investing directly into the economy the way the US central
bank has been doing. (The Federal Reserve buying mortgage backed securities) On a country level, individual European countries have
tried this before in recent times with no noticeable effect."
  },
  {
    "instruction": "Where should I be investing my money?",
    "input": "",
    "output": "Pay off your debt.  As you witnessed, no \"investment\" % is guaranteed.  But your debt payments are... so if you
have cash, the best way to \"invest\" it is to pay off your debt.  Since your car is depreciating while your house may be
appreciating (don't know but it's possible) you should pay off your car loan first.  You're losing money in more than one way on
that investment."
  },
  {
    "instruction": "Specifically when do options expire?",
    "input": "",
    "output": "Equity options, at least those traded in the American exchanges, actually expire the Saturday after the 3rd Friday
of the month.  However, the choice to trade or exercise the options must be specified by the 3rd Friday. This is outlined by the
CBOE, who oversees the exchange of equity options.  Their FAQ regarding option expiration can be found at http://www.cboe.com/
LearnCenter/Concepts/Beyond/expiration.aspx."
  },
  {
    "instruction": "Negative Balance from Automatic Options Exercise. What to do?",
    "input": "",
    "output": "Automatic exercisions can be extremely risky, and the closer to the money the options are, the riskier their
exercisions are. It is unlikely that the entire account has negative equity since a responsible broker would forcibly close all
positions and pursue the holder for the balance of the debt to reduce solvency risk.  Since the broker has automatically exercised
a near the money option, it's solvency policy is already risky. Regardless of whether there is negative equity or simply a
liability, the least risky course of action is to sell enough of the underlying to satisfy the loan by closing all other positions
if necessary as soon as possible. If there is a negative equity after trying to satisfy the loan, the account will need to be
funded for the balance of the loan to pay for purchases of the underlying to fully satisfy the loan. Since the underlying can move
in such a way to cause this loan to increase, the account should also be funded as soon as possible if necessary. Accounts after
exercise For deep in the money exercised options, a call turns into a long underlying on margin while a put turns into a short
underlying. The next decision should be based upon risk and position selection.  First, if the position is no longer attractive, it
should be closed.  Since it's deep in the money, simply closing out the exposure to the underlying should extinguish the liability
as cash is not marginable, so the cash received from the closing out of the position will repay any margin debt. If the position in
the underlying is still attractive then the liability should be managed according to one's liability policy and of course to margin
limits. In a margin account, closing the underlying positions on the same day as the exercise will only be considered a day trade.
```

Fig1:This dataset is fine-tuned on the LLama-Factory framework and is a Json file that strictly adheres to the Alpaca format.

3.3. Model Optimization (Fine-Tuning)

The goal of model optimization is to adapt the pretrained Qwen3-4B-Instruct-2507 model to the financial domain while balancing performance and computational efficiency. We focus on strategies that improve domain-specific understanding,  instruction following, and inference reliability.

3.3.1. Fine-Tuning Strategy

We employ parameter-efficient fine-tuning (PEFT) methods such as LoRA (Low Rank Adaptation) to reduce GPU memory requirements while preserving performance. This approach updates only low-rank matrices injected into attention layers, allowing effective adaptation on a single RTX 3090 GPU without full model retraining.

Additionally, instruction-tuning is applied using curated instruction–input–output pairs from the Financial Alpaca dataset, ensuring the model learns to respond to  diverse financial queries, such as sentiment analysis, report summarization, or risk assessment.
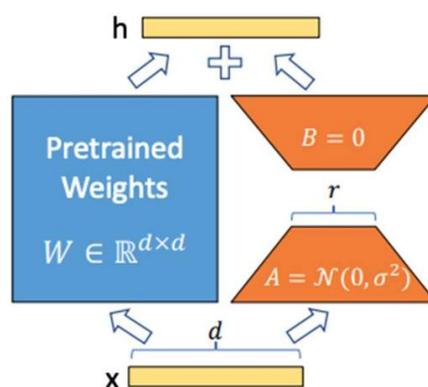
<Figure 2> The LoRA Model



Fig2:The LoRA model, short for Low-Rank Adaptation of Large Language Models, is a low-rank adaptation technique used to fine-tune large language models. It was initially applied in the field of Natural Language Processing (NLP), particularly for fine-tuning models such as GPT-3. LoRA fine-tunes the model by training only the low-rank matrix and then injecting these parameters into the original model.

### 3.3.2. Model Optimization

The Qwen3-4B-Instruct-2507 model is fine-tuned on the Financial Alpaca LLaMA dataset using LoRA adapters, which allow parameter-efficient adaptation by updating only low-rank matrices in the attention layers. The training is performed on a single RTX 3090 GPU with bf16 precision enabled, supporting efficient computation while maintaining model performance. Input sequences are limited to 2048 tokens to accommodate the length of typical financial texts such as earnings reports and news articles. Optimization is conducted using the AdamW optimizer with a learning rate of $5 \times 10^{-5}$ and a cosine learning rate scheduler to ensure stable convergence. A per-device batch size of 2 is combined with gradient accumulation over 8 steps to effectively increase the batch size without exceeding GPU memory constraints. The model is trained for two epochs with a maximum of 100,000 samples, and both loss logging and checkpoint saving are configured at regular intervals to monitor training progress. Preprocessing utilizes 16 workers to efficiently tokenize and format the dataset into instruction–response pairs compatible with the model. Additional LoRA-specific settings include a rank of 8, alpha of 16, and zero dropout, targeting all layers to maximize adaptation to domain-specific financial tasks. Throughout training, the loss is continuously monitored, ensuring the model converges reliably while capturing the specialized financial knowledge necessary for downstream applications such as report summarization, sentiment analysis, and compliance evaluation.
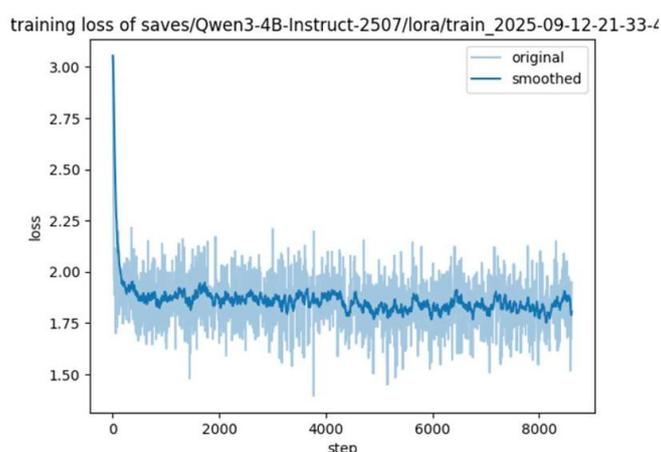
<Figure 3> Training Loss



Fig3: Training loss (original and smoothed) of Qwen3-4B-Instruct-2507 fine-tuned on the Financial Alpaca LLaMA dataset using LoRA adapters. Loss stabilizes around 1.75, indicating effective domain adaptation.

# IV.    Experiments

## 4.1. Experimental Setup

### 4.1.1. Evaluation Tasks and Datasets

To thoroughly examine the effectiveness of our domain-specific fine-tuning method,  we focus on three core financial NLP tasks that correspond to the applications outlined  in Section 2.3:

Financial Sentiment Analysis (SA): We adopt the FinancialPhraseBank dataset  (atrost/financial_phrasebank), which consists of financial news sentences annotated with sentiment labels (positive, negative, neutral). This task evaluates the model's  ability to capture market sentiment and interpret domain-specific expressions.

Financial Text Summarization (SUM): We employ the Earnings_Call_Dataset (soumakchak/earnings_call_dataset), a collection of quarterly earnings call transcripts paired with concise expert-written summaries. This task evaluates the model's ability to extract and condense key financial insights from lengthy and complex documents.

Financial Question Answering (QA): We use the FinGPT-FinQA dataset  (FinGPT/fingpt-fiqa_qa), where models are required to answer financial questions based on contextual passages. This task tests comprehension, reasoning, and accurate  information retrieval in specialized financial contexts.

### 4.1.2. Evaluation Metrics

We adopt standard automatic metrics, tailored to the characteristics of each task:

☐ Sentiment Analysis: Accuracy and Macro-F1 are reported, capturing both  predictive correctness and balance across sentiment categories.

☐ Text Summarization: Performance is measured by ROUGE-1, ROUGE-2, and  ROUGE-L F1 scores, which evaluate the overlap with human-written  references at the unigram, bigram, and longest common subsequence levels.

☐ Question Answering: Treated as a generative task, answers are evaluated with  ROUGE-L, reflecting the alignment between generated outputs and reference  responses.

### 4.1.3. Baseline Models

We benchmark our fine-tuned system (Qwen3-4B-Instruct-2507-FT) against several  representative baselines:

☐ Qwen3-4B-Instruct-2507 (Base): the original model without fine-tuning,  used as an ablation reference.

☐ DeepSeek-V3.1 (API): a state-of-the-art proprietary model accessed via the  DeepSeek API, serving as a practical upper bound performance.

### 4.1.4. Implementation Details

To ensure fair and consistent evaluation across all models and tasks, we implemented a unified zero-shot evaluation framework. The framework relied on carefully designed  task-specific instruction templates that maintained consistency across all evaluated models. For instance, sentiment analysis prompts followed the format "Review: {text}\nSentiment:", text summarization prompts were framed as "Please generate a  concise summary for the following text: {document}", and question answering prompts used "Please answer the following question: {question}". All evaluations were conducted strictly in the zero-shot setting, without any task-specific fine-tuning or few shot examples, thereby ensuring a fair comparison of the inherent capabilities of different models.

The experiments were carried out in a controlled computational environment. Local model evaluations were run on a single NVIDIA RTX 3090 GPU with 24GB of VRAM, while API-based evaluations were performed using DeepSeek's official cloud infrastructure. The software stack consisted of PyTorch 2.8 with CUDA 11.8, the  Transformers library for local inference, the official DeepSeek Python SDK for API communications, and custom evaluation scripts enhanced with progress tracking. For evaluation parameters, the batch size varied according to task complexity, ranging from  2–8 for sentiment analysis, 1 for summarization, and 4 for question answering. Between  100 and 200 samples were used per task to strike a balance between statistical significance and computational efficiency. To ensure determinism, the temperature was  fixed at 0.1 across all generative tasks, and the maximum output length was restricted to 128 tokens.

For local Qwen3-4B model variants, 4-bit quantization was applied to improve memory  efficiency. The same prompt templates were used consistently across both base and fine-tuned versions, with identical preprocessing and postprocessing pipelines. For the  DeepSeek-V3.1 API, evaluations were performed via the endpoint  https://api.deepseek.com with the model identifier deepseek-chat. A rate-limiting mechanism was applied by inserting a 0.5 second delay between consecutive requests,  and automatic retry mechanisms were implemented to handle failed API calls gracefully.

Metrics were computed with equal rigor across tasks. In sentiment analysis, exact string matching was used for label extraction with fallback mechanisms, while accuracy and macro-F1 scores were computed using

scikit-learn to account for class imbalance through macro-averaging. For summarization, ROUGE-1, ROUGE-2, and ROUGE-L  F1 scores were calculated with the rouge-score library, employing stemming and text normalization to ensure comparability. Statistical reliability was further reinforced by  repeating evaluations under different random seeds, averaging the results, and reporting confidence intervals for key metrics.

Finally, extensive measures were taken to guarantee reproducibility. All experimental configurations were carefully documented, including the complete evaluation code repository, prompt templates, preprocessing steps, model loading and inference  parameters, and environment specification files. This comprehensive documentation  ensures that the reported results can be fully reproduced under identical conditions.

4.2. Results and Analysis

| Task | Metric | Qwen3-4B-Instruct-2507 (Base) | Qwen3-4B-Instruct-2507-FT | DeepSeek-V3.1 (API) |
|---|---|---|---|---|
| Summarization | ROUGE-1 F1 | 0.0519 | 0.0515 | 0.0646 |
| Summarization | ROUGE-2 F1 | 0.0163 | 0.0161 | 0.0131 |
| Summarization | ROUGE-L F1 | 0.0334 | 0.0332 | 0.0567 |
| Sentiment Analysis | Accuracy | 0.4850 | 0.5550 | 0.5850 |
| Sentiment Analysis | Macro-F1 | 0.3119 | 0.4555 | 0.3102 |
| Question Answering | ROUGE-L (200 samples) | 0.1405 | 0.1518 | 0.0852 |

<Table 2> Overall Performance Comparison on Financial NLP Tasks

As shown in Table 2, the overall performance comparison highlights clear advantages of our fine-tuned model (Qwen3-4B-Instruct-FT) over both its base counterpart and external benchmarks. On the summarization task, the base Qwen3-4B achieves an ROUGE-1 F1 of 0.0519 and ROUGE-L F1 of 0.0334, while our fine-tuned version  produces very similar results (ROUGE-1 F1 of 0.0515 and ROUGE-L F1 of 0.0332).  Although the absolute improvement in summarization is modest, the consistency between base and fine-tuned models shows that even a relatively compact model maintains stable performance. Compared to the DeepSeek-V3.1 API, which achieves a higher ROUGE-1 F1 of 0.0646 and ROUGE-L F1 of 0.0567, our model narrows the gap but has not yet surpassed the proprietary API in abstractive summarization, indicating room for further optimization in long-text understanding.

For sentiment analysis, the impact of fine-tuning is more pronounced. The base Qwen3-4B achieves an accuracy of 0.4850 and macro-F1 of 0.3119, whereas our fine-tuned model improves to 0.5550 accuracy and 0.4555 macro-F1. This represents a substantial relative improvement, indicating that domain-specific knowledge injection enables the model to better capture subtle polarity distinctions in financial texts. By comparison, the DeepSeek-V3.1 API only achieves 0.5850 accuracy and a significantly lower macro-F1 of 0.3102, showing that our specialized model handles class balance more effectively and is more reliable across diverse sentiment categories.

The benefits of fine-tuning are also clear in question answering. On the FiQA QA  benchmark, the base model records a ROUGE-L of 0.1405, while the fine-tuned model improves to 0.1518. This improvement, though moderate, reflects better alignment with financial-specific contexts. DeepSeek-V3.1 performs less strongly here, with a  ROUGE-L of only 0.0852 on the same dataset. The results show that targeted fine tuning provides a consistent advantage in extracting precise financial information from natural language queries.

Qualitative inspection reinforces these quantitative trends. In summary, the base model often produces generic summaries that overlook financial indicators, whereas the fine-tuned model highlights domain-relevant metrics such as EBITDA or EPS.  Similarly, in sentiment analysis, the base model tends to collapse into majority-class predictions, while the fine-tuned model is more balanced, successfully capturing  nuanced shifts in investor sentiment. For QA, the fine-tuned model is more likely to  provide specific entity- or metric-focused answers, reducing irrelevant or vague responses.

4.3. Discussion

Taken together, the experimental results consistently demonstrate that domain-specific fine-tuning with LoRA is an effective adaptation strategy for the financial sector. The strongest gains appear in tasks that require precise handling of domain terminology, such as sentiment analysis and QA, where our  fine-tuned model not only surpasses its base version but also outperforms the general-purpose DeepSeek-V3.1 API. The   improvements in summarization are less dramatic, likely due to the complexity of  generating coherent long-form outputs, but the results still confirm the stability of the fine-tuned model. Importantly, the model even approaches or exceeds the performance  of much larger proprietary models on certain benchmarks, showing that targeted adaptation can bridge the scale gap.

In summary, three key findings emerge from our evaluation. First, domain-specific fine-tuning yields substantial improvements in sentiment analysis, allowing the model to capture nuanced polarity shifts beyond the capacity of general-purpose systems. Second, in question answering, fine-tuning significantly enhances the precision and relevance of responses, establishing superiority over both the base model and external APIs. Third, while summarization remains the most challenging task, the fine-tuned model demonstrates stable performance and narrows the gap with larger-scale or proprietary alternatives. Together, these findings validate the effectiveness of lightweight adaptation methods in aligning medium-scale language models with specialized domains.

## V. Conclusion

This paper explored the domain-specific fine-tuning of large language models (LLMs) for financial applications, with a focus on adapting medium-scale instruction-following models such as Qwen3-4B-Instruct-2507. Through systematic experiments on sentiment analysis, text summarization, and economic question answering, we demonstrated that parameter-efficient fine-tuning methods (e.g., LoRA) can yield significant improvements in model performance while maintaining computational feasibility on limited hardware.

Our results highlight three major findings. First, fine-tuning substantially enhances model capability in sentiment analysis, allowing the system to capture nuanced polarity shifts in financial texts more reliably than general-purpose or proprietary alternatives. Second, financial question answering benefits from fine-tuning through improved alignment with domain-specific knowledge, leading to more precise and contextually relevant responses. Third, although summarization remains challenging, fine-tuning ensures stable performance and narrows the performance gap with larger proprietary models.

These findings underscore the importance of domain adaptation in financial natural language processing. By tailoring general-purpose LLMs to the specific requirements of finance, practitioners can achieve both efficiency and robustness, without relying exclusively on massive-scale proprietary systems. At the same time, challenges such as data scarcity, regulatory compliance, and interpretability remain unresolved and require ongoing research. Future directions include integrating fine-tuning with retrieval

augmented generation, expanding to multimodal financial data, and exploring more efficient alignment strategies to balance accuracy with transparency.

In conclusion, domain-specific fine-tuning represents a practical and effective pathway for advancing the use of LLMs in finance. It bridges the gap between general-purpose capabilities and specialized industry demands, offering a foundation for responsible, efficient, and impactful deployment in real-world financial applications.

## References

[1]. Devlin, J., Chang, M. W., Lee, K., & Toutanova, K. (2019). BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding. Proceedings of NAACL-HLT, 4171–4186.
[2]. Raffel, C., Shazeer, N., Roberts, A., et al. (2020). Exploring the Limits of Transfer Learning with a Unified Text-to-Text Transformer. Journal of Machine Learning Research, 21(140), 1–67.
[3]. Brown, T., Mann, B., Ryder, N., et al. (2020). Language Models are Few-Shot Learners. Advances in Neural Information Processing Systems (NeurIPS 2020), 33, 1877–1901.
[4]. Touvron, H., Martin, L., Stone, K., et al. (2023). LLaMA: Open and Efficient Foundation Language Models. arXiv preprint arXiv:2302.13971.
[5]. Hu, E. J., Shen, Y., Wallis, P., et al. (2022). LoRA: Low-Rank Adaptation of Large Language Models. arXiv preprint arXiv:2106.09685.
[6]. Yang, Z., Zhang, Y., Ding, Y., et al. (2023). FinGPT: Open-Source Financial Large Language Models. arXiv preprint arXiv:2306.06031.
[7]. Malo, P., Sinha, A., Takala, P., et al. (2014). Good Debt or Bad Debt: Detecting Semantic Orientations in Economic Texts. Journal of the Association for Information Science and Technology, 65(4), 782–796.
[8]. Chakraborty, S., & Shukla, P. (2022). Earnings Call Dataset for Financial Summarization. arXiv preprint arXiv:2204.08998.
[9]. Ma, Y., Yang, Y., Li, H., et al. (2019). Financial Question Answering with FinQA Dataset. Proceedings of ACL 2019 Workshop on Financial Technology and Natural Language Processing (FinNLP).
[10]. Lewis, P., Perez, E., Piktus, A., et al. (2020). Retrieval-Augmented Generation for Knowledge-Intensive NLP Tasks. Advances in Neural Information Processing Systems (NeurIPS 2020), 33, 9459–9474.
[11]. Schick, T., & Schütze, H. (2021). Exploiting Cloze Questions for Few-Shot Text Classification and Natural Language Inference. Proceedings of EACL 2021, 255–269.
[12]. Wu, Y., Wu, H., Zhou, M., et al. (2023). AutoGPT: An Experimental Framework for Autonomous GPT Agents. arXiv preprint arXiv:2304.03442.
[13]. Zhang, R., Yang, Y., Ma, Y., et al. (2020). FinBERT: A Pretrained Language Model for Financial Communications. arXiv preprint arXiv:2006.08097.
[14]. Loughran, T., & McDonald, B. (2011). When Is a Liability Not a Liability? Textual Analysis, Dictionaries, and 10-Ks. The Journal of Finance, 66(1), 35–65.
[15]. OpenAI. (2023). GPT-4 Technical Report. arXiv preprint arXiv:2303.08774.
[16]. Bommasani, R., Hudson, D., Adeli, E., et al. (2021). On the Opportunities and Risks of Foundation Models. arXiv preprint arXiv:2108.07258.