# FINTEAMEXPERTS: ROLE-SPECIALIZED MOES FOR FINANCIAL ANALYSIS

Yue Yu Ney York University yy1879@nyu.edu

Prayag Tiwari
Halmstad University, Sweden
prayag.tiwari@hh.se

## **ABSTRACT**

Large Language Models (LLMs), such as ChatGPT, Phi3 and Llama-3, are leading a significant leap in AI, as they can generalize knowledge from their training to new tasks without fine-tuning. However, their application in the financial domain remains relatively limited. The financial field is inherently complex, requiring a deep understanding across various perspectives, from macro, micro economic trend to quantitative analysis.

Motivated by this complexity, a mixture of expert LLMs tailored to specific financial domains could offer a more comprehensive understanding for intricate financial tasks. In this paper, we present the FinTeamExperts, a role-specialized LLM framework structured as a Mixture of Experts (MOEs) for financial analysis. The framework simulates a collaborative team setting by training each model to specialize in distinct roles: Macro Analysts, Micro Analysts, and Quantitative Analysts. This role-specific specialization enhances the model's ability to integrate their domain-specific expertise. We achieve this by training three 8-billion parameter models on different corpus, each dedicated to excelling in specific finance-related roles. We then instruct-tune FinTeamExperts on downstream tasks to align with practical financial tasks. The experimental results show that FinTeamExperts outperform all models of the same size and larger on two out of four datasets. On the stock prediction, which presents a more complex task, FinTeamExperts still surpass all models of the same size. This highlights the success of our role-based specialization approach and the continued training approach for FinTeamExperts.

Keywords LLMs · Financial Analysis · Mixture of Experts

# 1 Introduction

Traditional machine learning methods, such as Support Vector Machines (SVMs) [1], gradient-boosted trees [2], and logistic regression [3], have limitations when it comes to understanding, reasoning, and generalizing in language-specific tasks. This limitation was overcome with the introduction of transformer architecture [4], which relies solely on attention mechanisms, removing the need for recurrence and convolutions. By training on large datasets using autoregressive techniques, Large Language Models are able to capture context and semantic dependencies in language. They are highly effective in translating across languages, processing large volumes of data, delivering quick responses with minimal delay, and can be fine-tuned to handle specific tasks and domains.

LLMs are increasingly used in the finance industry. They improve customer service with chatbots, help summarize information, recommend relevant knowledge, and automate tasks like filling out forms. LLMs also support risk management by analyzing market and credit risks and detecting anomalies or sentiment analysis. In investment, they assist with quantitative analysis, and in document processing, they help ensure compliance with regulations. However, research in financial Large Language Models (FinLLMs) remains limited. The earlier model, FinBERT[5], adapts BERT[6] for financial sentiment analysis by pre-training on financial documents and fine-tuning with a sentiment-specific dataset, achieving superior results but struggling with tasks beyond sentiment analysis. FinMA[7] fine-tuned Llama with financial instructions, showing competitive performance but not surpassing larger models like GPT-4.

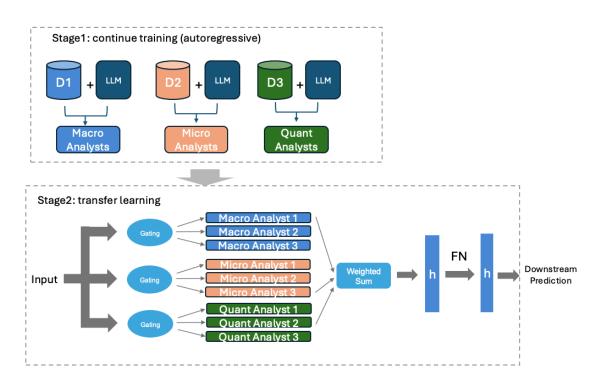


Figure 1: Architecture of FinTeamExperts with role specialized MOEs for financial analysis.

Despite these attempts in the financial domain, the intricate and domain-specific nature of financial tasks poses challenges that general-purpose LLMs are not well-suited to address. To bridge this gap, we introduce FinTeamExperts, a novel framework of LLM designed as a Mixture of Experts (MoEs) [8] to excel in finance-specific tasks. FinTeamExperts consist of three 8-billion-parameter models, each carefully tailored to comprehend and generate content relevant to the finance domain. MoEs, techniques that are used to enhance LLM performance without increasing computational demands, is used in the proposed framework to improve the domain-specific performance. This framework draws inspiration from real-world team dynamics, where each member develops expertise in distinct areas. By employing this specialized approach, we aim to create a powerful model capable of effectively handling the complex and nuanced demands of financial analysis and decision-making.

Our methodology follows a two-phase training process, as shonw in the Figure 1. In the first phase, the models are pre-trained on a curated corpus of role-related data to establish a strong foundational understanding of financial concepts. In the second phase, models are combined using routing gates and undergo instruct-tuning. This fine-tuning process helps the models produce outputs that are well-aligned with real-world tasks in the finance domain. The supervised instruct-tuning ensures that FinTeamExperts are skilled at applying knowledge from different experts to different real-world scenarios, making them effective for complex financial analysis and decision-making tasks.

We evaluate FinTeamExperts on four public datasets, where it consistently outperforms models of the same size. These results demonstrate the model's superior performance across a range of finance-specific tasks, highlighting the effectiveness of the proposed team experts settling and the specialized two phase training framework.

FinTeamExperts represent a significant advancement in the MOE framework of large language models to the financial domain. By offering domain-specific roles and training methodology, the framework provides a valuable tool for financial professionals, enhancing the performance and reliability of financial analysis and decision making in complex real-world scenarios.

We summarize our contributions as follows:

- A novel real-world role-based framework, FinTeamExperts, where each expert is specialized in a distinct financial role.
- An adapted MoE framework that enables team-based specialization, enhancing task-specific performance.
- The framework's effectiveness is validated through extensive experiments, showcasing superior performance on finance-related tasks.

An ablation study that demonstrates the individual contributions of each expert model within the framework.

# 2 Methodologies

In this section, we present the framework, outline the roles within the team, explore the adaptation of learning processes, and discuss the architectures of MOEs.

The FinTeamExperts architecture consists of two stages designed for financial analysis using a MoE approach. In stage one of continue autoregressive training, LLMs are fine-tuned with different datasets (marked as D1, D2, D3 in Figure 1) to create specialized models, referred to as Experts. Each expert model is trained on its respective dataset, allowing it to specialize in a specific domain or task relevant to financial analysis. This phase focuses on autoregressive training, where the experts continue with predicting the next token in the given corpus.

In stage two of transfer learning, the trained experts are applied together to new financial tasks. The input is passed through a gating mechanism, which selects the most appropriate experts for the task at hand. The outputs from these selected experts are combined via a weighted manner and processed by hidden layers and a feedforward network, which refine the aggregated information. This setup enables the model to generalize and make downstream predictions for financial applications, such as trend analysis or risk assessment.

#### 2.1 Framework Formulation

The core of the MOE is the routing mechanism of inputs to Experts. We set up a group of experts, corresponding to macro experts, micro experts, and quant experts. A hard gating mechanism is used to select one expert from each group, and their outputs are then weighted via another soft gating mechanism. The routing gate can be expressed as  $g(x_i)g(x_i) = (g_1(x_i), g_2(x_i), \dots, g_K(x_i))$  representing a probability distribution over the experts.

$$g(x_i) = \text{Softmax}(W_q \cdot x_i) \tag{1}$$

where  $x_i$  is the input token and  $W_g$  is the gating network weight matrix, which learns to route inputs to appropriate experts.

$$y_{i} = \sum_{j=1}^{K} g_{j}(x_{i})E_{j}(x_{i})$$
 (2)

In a simplified three-role gating framework, we can denote the output  $y_i$  for an input token  $x_i$ , as a weighted sum of the outputs from the selected experts. Each expert  $E_j(x_i)$  of expert j produces a unique output for teh input  $x_i$ . K is the total number of experts available, for instance, three in our study. Typically, only a subset of these experts is used for each input but we have it as a weighted contribution.

To build a adaptable framework where a flexible number (e.g., K) of experts caon contribute to each expert category. The weighted output  $y_i$  is thus adjusted as a sum over all roles and their respective experts, with each category weighted by a corresponding gating functions. This is defined as:

$$y_i = \sum_{r \in \{\text{Macro,Micro,Quant}\}} \sum_{j=1}^K g_r(x_i) \cdot h_{r,j}(x_i) \cdot E_{r,j}(x_i)$$
(3)

where r denotes the role (e.g., Macro, Micro, Quant), and  $g_{r,j}(x_i)$  is the gating function that assigns weights to each expert. In particular, we set

$$h_{r,j} = \text{GumbelSoftmax}(W_g \cdot x_i), \tag{4}$$

for all r to select one expert from each group of experts.

For inference, we define the inference head  $\hat{y}_t = C(y_{j < t})$ , which processes outputs from previous tokens to generate the next prediction.

The overall loss for the FinTeamMOE model,  $\mathcal{L}_{task}$ , includes both the task-specific loss and a regularization term. The task-specific loss, such as cross-entropy, is used for training the model on the primary task:

$$\mathcal{L}_{\text{MOE}} = \mathcal{L}_{\text{task}} + \lambda \cdot H(q(x)) \tag{5}$$

where  $\lambda$  is a regularization coefficient controlling the influence of the entropy term H(g(x)). This entropy term promotes balanced utilization of experts by encouraging the gate function to distribute weights across experts effectively.0

#### 2.2 Teamed Roles

To optimize performance and achieve strategic goals, FinTeamExperts focus on three pivotal roles within the investment setting:

**Macro Analysts**: Macro analysts study broader economic trends, geopolitical events, and global market movements to inform investment strategies. Their insights are crucial for understanding the larger economic context in which specific investments operate, enabling more informed decision-making.

**Micro Analyst**: Portfolio managers are responsible for making investment decisions and managing a portfolio of assets. They analyze market trends, assess risk, and determine the best investment strategies to maximize returns while adhering to the company's risk tolerance and investment objectives. Their role is vital in driving the overall investment strategy and ensuring alignment with the company's financial goals.

**Quantitative Analysts**: Quants develop mathematical models and algorithms to analyze financial data and identify trading opportunities. They work on creating predictive models, optimizing trading strategies, and automating trading processes using statistical and computational techniques. In the age of data-driven decision-making, quants provide the technical expertise needed to enhance trading efficiency, manage risks, and uncover alpha-generating opportunities.

Together, these roles form a comprehensive team within FinTeamExperts that balances strategic decision-making, technical innovation, and economic analysis, driving success in the competitive landscape of trading companies.

## 2.3 Adapting LLM Knowledge

We continue training the LLM checkpoints (such as GPT-2 and LLaMA-3-8B) using typical autoregressive learning with the next token prediction as the learning objective, defined as:

$$P(x_1, x_2, \dots, x_T) = \prod_{t=1}^{T} P(x_t \mid x_1, x_2, \dots, x_{t-1})$$
(6)

where  $x_1, x_2, \dots, x_T$  are the tokens from the corpus. The model generates each token  $x_t$  conditioned on the previous tokens.

Adapting large language models to specialize in finance-domain tasks involves curating a comprehensive financial corpus and implementing role-specific training. The initial phase includes pretraining the model on a curated dataset of market reports, economic analyses, and financial statements, building a robust understanding of financial terminology and concepts. Each model within FinTeamExperts is then trained to specialize in one of the three roles, developing expertise in macroeconomic analysis, asset management, or statistical trading techniques. This specialized training enables the models to integrate their expertise, forming a comprehensive financial analysis tool.

## 2.4 Mixture-of-Experts Architecture

The FinTeamExperts framework leverages the Mixture of Experts (MOEs) architecture to optimize performance in finance-related tasks. MOEs utilize multiple expert models, each specializing in different aspects of financial analysis, and dynamically route queries to the most relevant expert. This architecture allows for efficient resource allocation and enhances the model's ability to handle diverse and complex financial scenarios.

We employ Dynamic Routing: Queries are dynamically routed to the most appropriate expert model based on the specific financial task, improving accuracy and efficiency.

Specialized Training: Each expert model undergoes specialized training focused on its designated role (Macro Analysts, Portfolio Managers, Quantitative Analysts), ensuring depth of knowledge and proficiency in specific areas.

Hierarchical Expertise: The architecture supports hierarchical expertise, where higher-level experts oversee and refine the outputs of lower-level models, ensuring coherent and high-quality analysis. This MOEs-based methodology allows FinTeamExperts to leverage the strengths of individual expert models while maintaining flexibility and adaptability in addressing a wide range of financial tasks. The integration of these innovations demonstrates the potential of advanced LLMs in transforming financial analysis and decision-making.

# 3 Experiments

## 3.1 Role-specific Continue Training Datasets

In this study, we utilize three categories of datasets to continue pretrain models for financial roles: Macro of financial News, Financial Statements and Reports, and Market and Transactional Data. Each category serves a distinct role in equipping the model with the necessary knowledge to support the varied analytical needs of macroeconomic, microeconomic, and quantitative investment strategies.

**Financial News Datasets:** This category includes datasets that capture the sentiment and content of financial news articles, providing valuable insights into market trends and investor sentiment. One notable resource is the Thomson Reuters News Analytics (TRNA) dataset, which offers sentiment data derived from a vast archive of news articles. These datasets are crucial for developing models that can interpret the influence of global events and market news on investment decisions.

**Financial Statements and Reports Datasets:** This category encompasses structured financial data extracted from corporate reports such as balance sheets, income statements, and cash flow statements. The SEC EDGAR database <sup>1</sup> provides access to detailed financial disclosures filed by publicly traded companies in the U.S., including annual (10-K) and quarterly (10-Q) reports. Similarly, the WRDS database <sup>2</sup> offer comprehensive financial statement data from global companies. These datasets enable models to perform detailed microeconomic analysis, assessing the financial health and performance of individual companies.

**Market and Transactional Data Datasets:** This category includes datasets containing historical and real-time market data, such as stock prices, trading volumes, and other transactional data. The NASDAQ Data <sup>3</sup> provides a wide range of financial, economic, and alternative data, while the Alpha Vantage API offers access to real-time and historical market data across various asset classes. Additionally, Yahoo Finance <sup>4</sup> offers extensive historical market data and financials, which are essential for developing quantitative models that optimize trading strategies and identify market opportunities.

By integrating these diverse datasets, we create a robust foundation for training models that can navigate the complexities of financial markets, supporting the analytical needs of macro analysts, micro analysts, and quantitative analysts alike.

#### 3.2 Downstream Datasets

FPB dataset [9] contains 4,840 sentences from English-language financial news articles, each labeled by sentiment. Sentences are categorized based on the level of agreement among 5 to 8 annotators, providing a clear indication of consensus in sentiment labeling.

The FLARE-FIQASA dataset [10] is a labeled collection of financial texts, such as social media posts and headlines, used for sentiment analysis. It categorizes each text as positive, negative, or neutral to help analyze market sentiment in finance-focused content.

FinQA is designed for numerical reasoning over financial data, integrating both textual and tabular information. It contains examples with pre-text, post-text, and tables, along with associated questions and reasoning programs to compute numerical answers.

FOMC dataset [11] is a hawkish-dovish classification task involves categorizing statements based on sentiment toward monetary policy. 'Hawkish' statements signal a preference for tightening policy to control inflation, while 'dovish' ones indicate support for accommodative measures to boost growth. This classification is challenging due to the nuanced language in FOMC statements, which influences financial markets and economic expectations.

## 3.3 Settings

We build our expert model, on top of GPT-3-Large (with less than 1B parameters), noted as E-1B, and LLaMA-3-8B, noted as E-8B.

E-1B consists of 24 layers, each with 16 attention heads and a hidden size of 1,536, while E-8B comprises 36 layers, 32 attention heads, and a hidden size of 4,096. Using pretrained weights as the backbone for both models, we further train them on their respective corpora as specified in section 3.1.

<sup>1</sup>https://www.sec.gov/

<sup>&</sup>lt;sup>2</sup>https://wrds-www.wharton.upenn.edu/

<sup>3</sup>https://data.nasdaq.com/

<sup>&</sup>lt;sup>4</sup>https://www.yahoofinanceapi.com/

Speaker	Dialog			
User	What is the sentiment of the following financial post: Positive, Negative, or Neutral Text: What's up with \$LULU? Numbers looked good, not great, but good. I thir conference call will instill confidence.			
Assistant	The sentiment is neutral.			
User	Examine the data and tweets to deduce if the closing price of \$cvx will boost or lower at 2017-02-21. Kindly confirm either Rise or Fall.  Context: date,open,high,low,close,adj-close,inc-5,inc-10,inc-15,inc-20,inc-25,inc-30 2017-02-06,0.3,0.5,-0.5,-0.5,-0.5,-0.7,0.6,1.2,1.6,2.1,2.5 2017-02-07,1.4,1.5,-0.3,-1.4,-1.4,0.8,1.6,2.4,2.8,3.3,3.7 2017-02-08,-0.4,0.2,-1.1,0.2,0.2,0.7,0.9,2.0,2.5,2.9,3.3			
	2017-02-06: #dividend growth investing at #work - the streak continues! #investing #dividends \$cvx  rt AT_USER chevron corporation \$cvx shares bought by los angeles capital management & equity research inc.  2017-02-07: rt AT_USER 4 energy stocks that are ticking time bombs \$cvx \$oxy \$cop rt AT_USER mike liss of AT_USER century (\$twvlx) talked about \$cvx \$apc \$mdlz \$mdt & \$mrk during "hold it or fold it"  2017-02-08: rt AT_USER dynamic capital management ltd reduces stake in chevron corporation \$cvx dynamic capital management ltd reduces stake in chevron corporation \$cvx rt AT_USER #chevron partners zinox for co 2017-02-09: Answer:			
Assistant	The prediction is Rise.			

Table 1: Example of a dialog between the user and the assistant.

Model	Size	FPB	FiQA-SA	FOMC
FinMA-Full [7]	7B	87.0	79.0	-
GPT4[12]	-	78.0	80.0	71.0
BLOOM [13]	176B	50.0	53.0	-
BloombergGPT [14]	50B	51.0	75.0	-
Llama-3-8B [15]	8B	84.95	65.1	59.2
Qwen-2-7B [16]	7B	52.0	57.0	63.0
Mistral-8×7B-v1 [17]	7B	29.0	16.0	37.0
FinTeamExpert	$3\times1B$	<u>89.3</u>	76.3	64.2
FinTeamExpert	$3\times 8B$	90.5	81.0	<u>66.5</u>

Table 2: Main Results of FinTeamExperts for sentiment tasks

For E-1B, we use a learning rate of  $1 \times 10^{-4}$  with 500 warmup steps, a cosine learning rate scheduler, and the Adam optimizer [19], with  $\beta_1 = 0.9$ ,  $\beta_2 = 0.95$ , and a weight decay of 0.1. The model is trained for one epoch.

For E-8B, we use a learning rate of  $1 \times 10^{-4}$  with 1,000 warmup steps, a cosine scheduler, and the Adam optimizer with  $\beta_1 = 0.9$ ,  $\beta_2 = 0.95$ , and an epsilon value of  $1 \times 10^{-6}$ . We apply a weight decay of 0.1 and train the model for one epoch.

Both models are trained using 16-bit mixed precision on four 24GB A10G GPUs with fully sharded data parallelism, implemented via the Accelerate library<sup>5</sup>.

<sup>&</sup>lt;sup>5</sup>https://huggingface.co/docs/accelerate

Model	Size	CIKM18
ChatGPT	=	55.0
GPT4 [12]	-	57.0
Gemini [18]	-	54.0
Qwen2-7B [16]	7B	52.0
FinMA 7B-full [7]	7B	53.0
Llama-3-8B [15]	8B	51.9
Mistral-8×7B-v1 [17]	7B	42.0
FinTeamExpert	3×1B	50.3
FinTeamExpert	$3\times 8B$	<u>56.3</u>

Table 3: Results on stock prediction dataset

## 3.4 Reseults

The FinTeamExpert models, particularly the 3×8B version, achieve the best performance across all sentiment tasks, with the highest scores in FPB (90.5), FiQA-SA (81.0), and FOMC (65.5). This suggests that increasing model size and using an ensemble of models enhances performance in financial sentiment analysis. Even the smaller FinTeamExpert (3×1B) performs competitively, especially in FPB (89.3) and FOMC (64.2), demonstrating that well-structured ensembles of smaller models can rival larger models in specialized tasks.

In comparison, FinMA-Full (7B) performs well on FPB (87.0) and FiQA-SA (79.0), though it falls slightly behind the FinTeamExpert models. On the other hand, GPT-4, while a powerful general-purpose model, lags behind in financial tasks with a lower FPB score (78.0). Large models like BLOOM (176B) and BloombergGPT (50B) underperform in FPB, highlighting that model size alone does not guarantee effectiveness without domain-specific adaptation.

LLaMA-3-8B (8B) shows strong results, particularly on FPB (84.95), demonstrating solid performance compared to similarly sized models. Overall, FinTeamExperts are the most effective across all tasks, indicating that specialized, ensemble-based models outperform general-purpose models in the financial sentiment domain.

Table 3 shows the result on stock prediction dataset. GPT-4 leads the CIKM18 task with a score of 57.0, closely followed by the domain-specific FinTeamExpert (3×8B) at 56.3, showing that both general-purpose and specialized ensemble models can excel in financial tasks. FinMA (7B) and LLaMA-3-8B achieve scores of 53.0 and 51.9, respectively, demonstrating solid performance but falling short of the top two models.

The smaller FinTeamExpert (3×1B) scores 50.3, highlighting that even smaller ensemble models remain competitive, though not as strong as their larger counterparts. Overall, GPT-4 and FinTeamExpert (3×8B) dominate, while other models offer reasonable results in the CIKM18 task.

## 3.5 Ablation Study

We conduct a perplexity analysis on three role-oriented training models and our FinTeamExperts model, using mini-test samples to assess their perplexity scores. Figure ?? illustrates these results, with a benchmark perplexity score of approximately 15 for the vanilla Llama-3-8B model before it was trained with role-specific adaptations. Perplexity trends are plotted for the macro, micro, and quant models, while downstream learning results are shown for FinTeamExperts. As training progressed, perplexity scores for all three role-oriented models decreased and eventually converged around 6.5. Among these, the macro-role model achieved the lowest perplexity, with the quant-role model slightly higher, though the differences were minor.

FinTeamExperts, routed and fine-tuned for downstream financial tasks, began with an initial perplexity below 8 and converged to a value under 6. This indicates that FinTeamExperts is better adapted to the financial dataset, demonstrating enhanced understanding and alignment with the financial tasks.

The ablation study, as shown in Figure 3, reveals how removing individual roles—Macro, Micro, or Quant—impacts model performance across different tasks. The results clearly show that the full FinTeamExperts, without any role removed, consistently outperforms all ablated versions, particularly in complex tasks like FPB. This highlights the critical importance of each role in achieving top performance. For InFiQA-SA, dropping any role yields similar performance, suggesting that no single role is indispensable for this task. However, in FPB, the complete model far exceeds any version with a role removed, underscoring the collaboration required for optimal results. For FOMC,

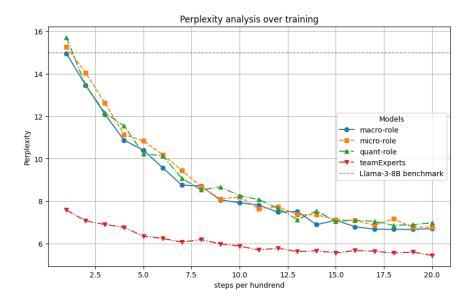


Figure 2: Perplexity analysis

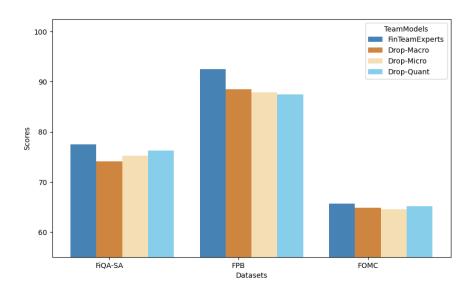


Figure 3: Ablation study of drop-one analysis

droping Micro lead to lower performance, indicating the importance of micro knowledge of this task. Overall, these findings reinforce the value of having all roles intact for the most challenging tasks.

# 4 Related Works

# 4.1 Financial LLMs

Prior to the rise of LLMs, deep learning played a pivotal role in the financial domain, particularly in tasks like portfolio management [20] and risk management [21]. In portfolio management, deep learning can be utilized to enhance the optimization process by directly maximizing metrics such as the Sharpe ratio. This approach simplifies asset selection by effectively capturing correlations across different asset classes, thereby streamlining decision-making in constructing a balanced portfolio.

In the era of advanced LLMs, research specifically targeting financial language models (FinLLMs) remains relatively limited, as highlighted by a recent survey [22]. One notable example is FinBERT [5], which adapts the general-purpose BERT model for financial sentiment analysis using a two-step approach. First, it is pre-trained on financial texts, utilizing a subset of the Reuters TRC2 dataset to enhance its understanding of financial terminology. Next, a dense layer is added to the final hidden state of the classification token (CLS), and the model is fine-tuned using the Financial PhraseBank (FPB) dataset. FinBERT achieves outstanding performance in financial sentiment analysis, surpassing state-of-the-art models in this specific task. However, it remains limited in scope and does not perform as effectively on other financial tasks.

FinGPT [23] provides an end-to-end framework for training and applying financial large language models in the finance industry. It uses the efficient Low-rank Adaptation (LoRA) [24] technique to fine-tune open-source models like LLaMA and ChatGLM with around 50,000 samples. However, its evaluation is currently limited to finance classification tasks only.

BloombergGPT [14], a 50-billion parameter model based on BLOOM's architecture, is one of the first decoder-only LLMs trained specifically for finance. Trained from scratch on 363 billion tokens from financial documents and 345 billion from general datasets, it predicts the next token in documents without fine-tuning on instructions. However, its results lag behind those of other models, including some much smaller ones, as detailed in [25].

FinMA [7] curated an instruction dataset for a financial LLM and fine-tuned it on Llama, resulting in strong performance among similar-sized LLMs, though not surpassing larger models like GPT-4.

# 4.2 Mixture of Experts

Mixture of Experts (MoEs), derived from Gaussian Mixture models, are employed to boost the performance of large language models (LLMs) without increasing computational resource requirements.

SwithTransformer [26], places MoE layers after the multi-head attention mechanism in each transformer block to select the feedforward layers. Unlike LLMs where all parameters are used for every input, SwitchTransformer activates only a small subset of experts for each input, significantly reducing computational costs. A routing network determines which expert should handle each input, ensuring that only the most relevant parts of the model are utilized for each task.

GLaM [27], a model with 1.2 trillion parameters, is approximately 7 times larger than GPT-3, yet it achieves this scale with only a third of the energy consumption required to train GPT-3. GLaM's architecture integrates MoE layers with Transformer layers. For each input token, a gating module selects the two most relevant experts, and a weighted average of their outputs is passed to the next Transformer layer. The gating network, which is trained to identify the optimal experts for each token in the input sequence, ensures efficient use of computational resources while maintaining high performance.

ST-MoE [28], is a sparse Transformer model with 269B parameters, offering performance at a computational cost similar to that of a dense 32B parameter encoder-decoder Transformer. The model recommends a top-2 routing mechanism, with each input token routed to the two most relevant experts, ensuring computational efficiency. A capacity factor of 1.25 is recommended, which controls the number of tokens processed by each expert, and this factor can be adjusted during evaluation to meet changing memory or computational requirements. Additionally, quality improvements are achieved through dense layer stacking and the introduction of a multiplicative bias.

# 5 Conclusion

In this paper, we introduced FinTeamExperts, a novel framework of role-specialized large language models (LLMs) designed as a Mixture of Experts (MOEs) to excel in financial analysis tasks. By mimicking a real-world team setting, each model in FinTeamExperts specializes in one of three critical roles: Macro Analysts, Portfolio Managers, and Quantitative Analysts. This specialization allows the models to integrate their expertise effectively, forming a comprehensive and robust financial analysis tool.

Our approach includes several innovative contributions. Firstly, we are the first to implement role-based teams of LLMs as MOEs, aiming to mimic practical implementation scenarios within the finance domain. This method leverages the strengths of individual expert models while maintaining flexibility and adaptability in handling a wide range of financial tasks. Secondly, we introduced advancements in the MOEs architecture, such as dynamic routing, specialized training, and hierarchical expertise, which significantly enhance the model's performance in downstream financial tasks.

Through instruct-tuning and rigorous experiments, we demonstrated that FinTeamExperts outperform existing LLMs in finance-related tasks, underscoring the effectiveness of our pretraining methodology and specialized training strategies.

These contributions showcase the potential of advanced LLMs in transforming financial analysis and decision-making, paving the way for more sophisticated and practical AI applications in the finance industry.

There are several directions for future exploration and enhancement. First, expanding the diversity of role-based teams within the MOEs framework could further refine task-specific expertise, particularly by incorporating specialized models trained on emerging financial topics, such as ESG (Environmental, Social, and Governance) criteria or digital asset analytics. Additionally, investigating the effects of cross-domain transfer learning may yield insights into how financial LLMs can benefit from knowledge in adjacent fields, such as legal or regulatory compliance, which often intersect with financial analysis. Another promising direction is the optimization of dynamic routing strategies to allow for more granular control over model selection based on real-time task complexity and data characteristics. This could involve developing adaptive routing algorithms that leverage reinforcement learning or other self-learning techniques, allowing the MOEs framework to continuously improve task assignment efficiency and performance.

## References

- [1] Hui Xue, Qiang Yang, and Songcan Chen. Svm: Support vector machines. In *The top ten algorithms in data mining*, pages 51–74. Chapman and Hall/CRC, 2009.
- [2] Jerry Ye, Jyh-Herng Chow, Jiang Chen, and Zhaohui Zheng. Stochastic gradient boosted distributed decision trees. In *Proceedings of the 18th ACM conference on Information and knowledge management*, pages 2061–2064, 2009.
- [3] Michael P LaValley. Logistic regression. Circulation, 117(18):2395–2399, 2008.
- [4] Cem Subakan, Mirco Ravanelli, Samuele Cornell, Mirko Bronzi, and Jianyuan Zhong. Attention is all you need in speech separation. In *ICASSP 2021-2021 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 21–25. IEEE, 2021.
- [5] Zhuang Liu, Degen Huang, Kaiyu Huang, Zhuang Li, and Jun Zhao. Finbert: A pre-trained financial language representation model for financial text mining. In *Proceedings of the twenty-ninth international conference on international joint conferences on artificial intelligence*, pages 4513–4519, 2021.
- [6] Jacob Devlin. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv* preprint *arXiv*:1810.04805, 2018.
- [7] Qianqian Xie, Weiguang Han, Xiao Zhang, Yanzhao Lai, Min Peng, Alejandro Lopez-Lira, and Jimin Huang. Pixiu: A large language model, instruction data and evaluation benchmark for finance, 2023.
- [8] Weilin Cai, Juyong Jiang, Fan Wang, Jing Tang, Sunghun Kim, and Jiayi Huang. A survey on mixture of experts. *arXiv preprint arXiv:2407.06204*, 2024.
- [9] P. Malo, A. Sinha, P. Korhonen, J. Wallenius, and P. Takala. Good debt or bad debt: Detecting semantic orientations in economic texts. *Journal of the Association for Information Science and Technology*, 65, 2014.
- [10] TheFinAI and ChanceFocus. FLARE-FIQASA dataset. https://huggingface.co/datasets/TheFinAI/flare-fiqasa, 2023. Financial sentiment analysis dataset for text classification tasks in finance.
- [11] Agam Shah, Suvan Paturi, and Sudheer Chava. Trillion dollar words: A new financial dataset, task & market analysis. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 6664–6679, Toronto, Canada, July 2023. Association for Computational Linguistics.
- [12] Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*, 2023.
- [13] Teven Le Scao, Angela Fan, Christopher Akiki, Ellie Pavlick, Suzana Ilić, Daniel Hesslow, Roman Castagné, Alexandra Sasha Luccioni, François Yvon, Matthias Gallé, et al. Bloom: A 176b-parameter open-access multilingual language model. 2023.
- [14] Shijie Wu, Ozan Irsoy, Steven Lu, Vadim Dabravolski, Mark Dredze, Sebastian Gehrmann, Prabhanjan Kambadur, David Rosenberg, and Gideon Mann. Bloomberggpt: A large language model for finance. *arXiv preprint arXiv:2303.17564*, 2023.
- [15] Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*, 2024.
- [16] An Yang, Baosong Yang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Zhou, Chengpeng Li, Chengyuan Li, Dayiheng Liu, Fei Huang, et al. Qwen2 technical report. *arXiv preprint arXiv:2407.10671*, 2024.

- [17] Albert Q. Jiang, Alexandre Sablayrolles, Antoine Roux, Arthur Mensch, Blanche Savary, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Emma Bou Hanna, Florian Bressand, Gianna Lengyel, Guillaume Bour, Guillaume Lample, Lélio Renard Lavaud, Lucile Saulnier, Marie-Anne Lachaux, Pierre Stock, Sandeep Subramanian, Sophia Yang, Szymon Antoniak, Teven Le Scao, Théophile Gervet, Thibaut Lavril, Thomas Wang, Timothée Lacroix, and William El Sayed. Mixtral of experts, 2024.
- [18] Gemini Team, Rohan Anil, Sebastian Borgeaud, Jean-Baptiste Alayrac, Jiahui Yu, Radu Soricut, Johan Schalkwyk, Andrew M Dai, Anja Hauth, Katie Millican, et al. Gemini: a family of highly capable multimodal models. *arXiv* preprint arXiv:2312.11805, 2023.
- [19] Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization, 2017.
- [20] Zihao Zhang, Stefan Zohren, and Stephen Roberts. Deep learning for portfolio optimization. arXiv preprint arXiv:2005.13665, 2020.
- [21] Akib Mashrur, Wei Luo, Nayyar A Zaidi, and Antonio Robles-Kelly. Machine learning for financial risk management: a survey. *Ieee Access*, 8:203203–203223, 2020.
- [22] Jean Lee, Nicholas Stevens, Soyeon Caren Han, and Minseok Song. A survey of large language models in finance (finllms). arXiv preprint arXiv:2402.02315, 2024.
- [23] Hongyang Yang, Xiao-Yang Liu, and Christina Dan Wang. Fingpt: Open-source financial large language models. *arXiv preprint arXiv:2306.06031*, 2023.
- [24] Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. Lora: Low-rank adaptation of large language models. *arXiv preprint arXiv:2106.09685*, 2021.
- [25] Pau Rodriguez Inserte, Mariam Nakhlé, Raheel Qader, Gaetan Caillaut, and Jingshu Liu. Large language model adaptation for financial sentiment analysis. *arXiv e-prints*, pages arXiv–2401, 2024.
- [26] William Fedus, Barret Zoph, and Noam Shazeer. Switch transformers: Scaling to trillion parameter models with simple and efficient sparsity. *Journal of Machine Learning Research*, 23(120):1–39, 2022.
- [27] Nan Du, Yanping Huang, Andrew M Dai, Simon Tong, Dmitry Lepikhin, Yuanzhong Xu, Maxim Krikun, Yanqi Zhou, Adams Wei Yu, Orhan Firat, et al. Glam: Efficient scaling of language models with mixture-of-experts. In *International Conference on Machine Learning*, pages 5547–5569. PMLR, 2022.
- [28] Barret Zoph, Irwan Bello, Sameer Kumar, Nan Du, Yanping Huang, Jeff Dean, Noam Shazeer, and William Fedus. St-moe: Designing stable and transferable sparse expert models. *arXiv preprint arXiv:2202.08906*, 2022.