

Efficient Fine-Tuning of Financial Language Models with QLoRA for Alpha Signal Generation

Fusheng Luo¹, Guojun Peng¹, Jiamu Chen¹, Yining Li¹

¹ Johns Hopkins University, USA

fluo5@jh.edu, gpeng9@jh.edu, jchen471@jh.edu, yli688@jh.edu

Abstract

In this study, we investigated the effectiveness of QLoRA-based fine-tuning on Meta-LLaMA3-8B for the classification of financial sentiment and signal generation. Leveraging a diverse set of domain-specific datasets and the SWIFT framework, our model achieved superior performance over established baselines—such as FinBERT (a classic sentiment analyzer in academia), RoBERTa, and VADER—in both sentiment classification and downstream trading strategies on a separate financial news dataset.

Keywords: Financial Sentiment Analysis, Large Language Models, QLoRA, LLaMA, Fine-tuning, Market-Neutral Trading Strategy, K-means Clustering

1 Introduction

Large Language Models (LLMs) have become increasingly popular in financial text analysis, providing notable advancements over traditional dictionary-based approaches in extracting predictive insights from unstructured sources like earnings calls, news reports, and regulatory documents. Their advantage is in capturing contextual details and sentiment subtleties often overlooked by simpler methods. However, generic LLMs encounter two main challenges in finance: (1) a lack of domain-specific training, which can cause misunderstandings of financial terminology, and (2) the high computational costs

25 associated with full-model fine-tuning, making them less accessible for many
26 users.

27 Recent studies have investigated fine-tuning foundation models like LLaMA
28 with finance-focused data to enhance return predictions. While optimizing
29 the entire model yields high performance, it demands substantial resources.
30 Parameter-efficient methods such as Low-Rank Adaptation (LoRA) provide
31 a viable alternative by adjusting only a limited number of parameters, thus
32 lowering memory use. However, there is still a balance to strike between tun-
33 ing efficiency and accuracy, particularly for users with hardware limitations.

34 This study explores QLoRA, a fine-tuning method combining 4-bit quanti-
35 zation and LoRA for efficient, scalable training without performance loss.
36 QLoRA keeps the base model weights fixed in 4-bit precision and adds train-
37 able 16-bit LoRA adapters, enabling models up to 65 billion parameters to be
38 fine-tuned on a single 48GB GPU [Dettmers et al., 2023] QLoRA has shown
39 competitive results in general NLP tasks; its use in financial modeling re-
40 mains largely unexplored. This work assesses whether QLoRA can preserve
41 strong predictive capabilities for alpha signal generation based on financial
42 news sentiment.

43 Beyond the financial domain, Zhao et al. [Zhao et al., 2024] evaluated
44 LoRA-based fine-tuning across 310 LLMs, covering 31 tasks and 10 base
45 models. Their findings reveal that 4-bit LoRA models frequently outperform
46 their base models and even GPT-4 on certain tasks, while utilizing consider-
47 ably less memory. They also introduce LoRAX, an inference server able to
48 run multiple LoRA adapters simultaneously on a single GPU, demonstrating
49 the practicality of efficient parameter deployment.

50 In finance, FinLLama [Iacovides et al., 2024] customizes LLaMA-2 using LoRA
51 and 8-bit quantization on labeled financial sentiment data. It efficiently gen-
52 erates sentiment scores and trading signals while reducing computational
53 demands [Chiu and Hung, 2024]. Nonetheless, FinLLama does not explore
54 finer quantization levels or lighter infrastructure options. Conversely, QLoRA
55 addresses these limitations by incorporating NormalFloat 4-bit quantiza-
56 tion, double quantization, and paged optimizers, enabling full 16-bit training
57 quality on low-resource hardware. Models like Guanaco, developed within
58 this framework, compete with ChatGPT on benchmarks such as Vicuna
59 [Dettmers et al., 2023], demonstrating QLoRA’s potential for accessible, high-
60 performance fine-tuning in finance.

61 To clearly illustrate our workflow, Figure 1 provides an overview of the end-

62 to-end pipeline for our sentiment-driven alpha signal generation system. We
 63 begin with data collection from four financial sentiment datasets—Financial
 64 PhraseBank, NASDAQ News Sentiment, Twitter Financial News, and FIQA2018—which
 65 are then preprocessed via label normalization and tokenization using the
 66 LLaMA3 tokenizer. The resulting data is used to fine-tune the Meta-LLaMA3-
 67 8B model using QLoRA, a 4-bit quantized parameter-efficient method. The
 68 fine-tuning process is orchestrated via the SWIFT framework, which facili-
 69 tates scalable and hardware-efficient training. The fine-tuned model is then
 70 applied to financial news headlines, producing discrete sentiment scores that
 71 are aggregated per stock. Finally, we design and evaluate three market-
 72 neutral trading strategies based on the predicted sentiments. Across all se-
 73 tups, our approach demonstrates superior performance compared to both the
 74 S&P 500 benchmark and alternative sentiment-based models, highlighting its
 75 practical value in quantitative trading.

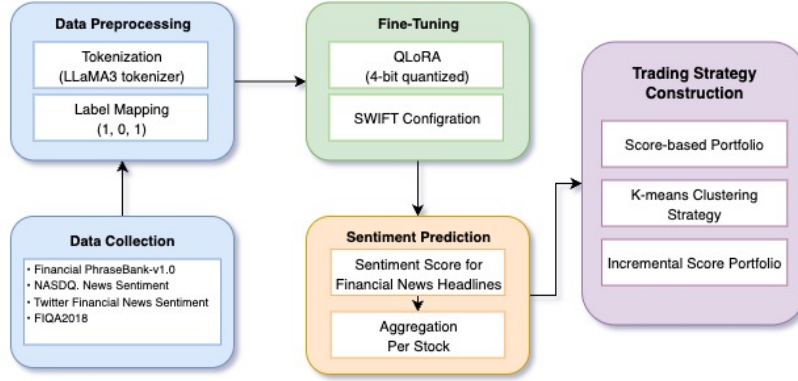


Figure 1: Workflow pipeline for financial sentiment modeling and trading.

76 1.1 Problem Setup

77 We aim to develop a sentiment classification model tailored to financial news
 78 articles using a fine-tuned LLaMA model. The task is formulated as a su-
 79 pervised learning problem where each input is a piece of financial text, and
 80 the output is a discrete sentiment label indicating the article’s sentiment
 81 regarding financial markets or assets.

82 **Problem Definition:** Let $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^N$ denote the training dataset,
 83 where each $x_i \in \mathcal{X}$ represents a financial news excerpt and $y_i \in \mathcal{Y}$ is the

84 associated sentiment label. The label set is defined as

$$\mathcal{Y} = \{-1, 0, 1\}$$

85 where:

- 86 • -1 denotes negative sentiment,
- 87 • 0 denotes neutral sentiment,
- 88 • 1 denotes positive sentiment.

89 The goal is to learn a function $f_\theta : \mathcal{X} \rightarrow \mathcal{Y}$ parameterized by θ (the model
90 parameters) that minimizes the expected classification error. Formally, we
91 seek:

$$\theta^* = \arg \min_{\theta} E_{(x,y) \sim \mathcal{D}} [\mathcal{L}(f_\theta(x), y)]$$

92 where \mathcal{L} is a loss function, here we adopt the cross-entropy loss for multi-class
93 classification.

94 **Model Setup:** We employ a pre-trained LLaMA model as the backbone,
95 fine-tuned with LoRA (Low-Rank Adaptation) techniques to enable efficient
96 parameter tuning. During fine-tuning, we adapt a small subset of paramet-
97 ers while freezing the majority of the pre-trained weights, thereby retaining
98 general language understanding while specializing the model for financial
99 sentiment detection.

100 **Assumptions:**

- 101 • The labeled dataset \mathcal{D} is representative of the financial news distribu-
102 tion we expect the model to encounter during inference.
- 103 • Sentiment expressed in financial news can be reasonably categorized
104 into the three discrete classes: negative, neutral, and positive.
- 105 • The fine-tuning process does not significantly degrade the base LLaMA
106 model’s language understanding capabilities outside the financial do-
107 main.

108 2 Data Sources and preprocessing

109 2.1 Dataset

110 We have divided the process of our model building into two parts: fine-tuning
111 the pre-trained Large Language Models (LLM) and using the post-trained
112 model to predict the sentiment of our News dataset and conducting trading
113 strategies combined with the dataset of stocks price in the US equity markets.

114 **Part 1: Fine-tuning Dataset** We generally use four datasets here

- 115 (i) Financial PhraseBank-v1.0
- 116 (ii) NASDAQ news sentiment
- 117 (iii) Twitter Financial News Sentiment
- 118 (iv) FIQA2018

119 **(i)**
120 Financial Phrase Bank dataset covers a collection of 4840 sentences.
121 The selected collection of phrases was annotated by 16 people with ad-
122 equate background knowledge on financial markets. Three of the anno-
123 tators were researchers and the remaining 13 annotators were master’s
124 students at the School of Business at Aalto University with majors pri-
125 marily in finance, accounting, and economics. They have labeled the
126 News with sentiments of positive (+1), negative (-1) and neutral (0).
127 E.g., The news ”The operating margin came down to 2.4 from 5.7.” is
128 labeled as negative.

129 **(ii)**
130 The ”NASDAQ news sentiment dataset” consists of synthetic data points
131 generated using GPT-4o, and can be used to train models for analysing
132 sentiment from the title and meta-description of an article. This can be
133 used for tasks such as analyzing market sentiment regarding a partic-
134 ular NASDAQ-listed stock from search results. GPT has also labeled
135 the news articles with positive, negative and neutral sentiments.

136 **(iii)**
137 The Twitter Financial News dataset is an English-language dataset

138 containing an annotated corpus of finance-related tweets. This data set
139 is used to classify financial-related tweets according to their sentiment.
140 The data set contains 11,932 documents annotated with 3 labels: label
141 0: Bearish; label 1: Bullish, label 2: Neutral.

142 (iv)
143 FIQA2018 dataset (often stylized as FiQA or FiQA: A Challenge on
144 Financial Opinion Mining and Question Answering) refers to a dataset
145 and shared task introduced for the WSDM Cup 2018. It is designed to
146 benchmark financial domain sentiment analysis. It also classifies news
147 sentiments into positive and negative. But it differs in that they give
148 a sentiment score from -1 to +1 to measure the sentiment of the News
149 (" +1" refers to highly positive, and "-1" refers to highly negative). A
150 score of 0 indicates neutral or without clear polarity.

151 **Part 2: Predicting Real-Time News** We use two datasets: the news
152 dataset, which includes news articles and alerts (with alerts representing
153 news with only headlines), and the corresponding stock price dataset. We
154 collected news text data from Refinitiv (now known as LSEG), CRSP, and
155 Yahoo Finance for stock prices. Then, we create a news dataset covering
156 a list of stocks within a time range of 16 months (2024.01.01 - 2025.05.30).
157 This news dataset includes 502 stocks [S&P Dow Jones Indices, 2025] and
158 77000 news headlines. Additionally, for these 502 companies, we also have
159 stock prices and volume data covering the same period from Yahoo Finance.
160 This is the fundamental dataset we use to evaluate our trading strategies in
161 comparison to the benchmark strategies discussed in the following sections.

162 2.2 Data Preprocessing

163 The preprocessing steps applied to the raw news data obtained from Refinitiv
164 (LSEG) are as follows. Due to the data extraction constraint of 100 records
165 per request, we adopt a recursive collection approach by dividing the full
166 time horizon into consecutive two-week intervals. This strategy allows us to
167 systematically retrieve news for each stock of interest over the entire sample
168 period. After compiling the data across all intervals, we concatenate them
169 into a unified dataset, append the relevant stock ticker symbols, and set the
170 date as the index to match the structure required for our daily long-short
171 trading strategy. To ensure consistency and prevent overcounting, we also
172 remove any duplicated records.

173 2.2.1 Label Mapping

174 To enable consistent sentiment classification across diverse datasets, we apply
175 a standardized label mapping function that transforms the original dataset-
176 specific sentiment annotations into a unified label space of $\{-1, 0, 1\}$, cor-
177 responding to *negative*, *neutral*, and *positive* sentiments, respectively. The
178 mapping rules for each dataset are as follows:

- 179 • Financial PhraseBank: The original labels are textual strings: "negative",
180 "neutral", and "positive". These are mapped directly via:

$$\text{"negative"} \rightarrow -1, \quad \text{"neutral"} \rightarrow 0, \quad \text{"positive"} \rightarrow 1$$

181 Any unexpected label is defaulted to 0 (neutral).

- 182 • FIQA 2018: The dataset provides a continuous sentiment score in the
183 range $[-1.0, 1.0]$. We discretize it as follows:

$$\text{score} < -0.2 \rightarrow -1, \quad \text{score} > 0.2 \rightarrow 1, \quad \text{otherwise} \rightarrow 0$$

- 184 • Twitter Financial Sentiment: Labels are originally encoded as integers:
185 0 (negative), 1 (positive), and 2 (neutral). These are remapped to:

$$0 \rightarrow -1, \quad 2 \rightarrow 0, \quad 1 \rightarrow 1$$

- 186 • NASDAQ News: The sentiment is expressed as an integer score, typi-
187 cally ranging from 0 to 5. The mapping is defined as:

$$\text{score} \leq 1 \rightarrow -1, \quad \text{score} \geq 4 \rightarrow 1, \quad \text{otherwise} \rightarrow 0$$

188 This label normalization procedure ensures all training samples conform to
189 the same target format, enabling consistent training and evaluation across
190 datasets.

191 2.2.2 Tokenizer

192 The tokenizer uses the official vocabulary and tokenization rules for LLaMA
193 3. It maps raw input to token IDs aligned with the model’s embedding space.

194 3 Methodology

195 3.1 QLoRA

196 **LoRA Adaptation.** Low-Rank Adaptation (LoRA) [Hu et al., 2022] is a
197 parameter-efficient fine-tuning method that freezes the original pretrained
198 weight matrix $W_0 \in R^{m \times n}$ and introduces two small trainable matrices $A \in$
199 $R^{m \times r}$ and $B \in R^{r \times n}$, where $r \ll n$. The weight update is parameterized as

$$\Delta W = AB, \quad (1)$$

200 so that the forward computation during fine-tuning becomes

$$y = (W_0 + \frac{\alpha}{r} AB)x, \quad (2)$$

201 where α is a scaling factor controlling the contribution of the LoRA update.

202 In our implementation, we set $r = 8$ and $\alpha = 32$, applying adapters to all
203 linear layers (including attention and feed-forward projections). This con-
204 figuration introduces less than 0.05% additional trainable parameters while
205 maintaining full model expressiveness.

206 **Quantile.** Using the LoRA components described above, we define QLoRA
207 [Dettmers et al., 2023] for a single linear layer in the quantized base model
208 with a LoRA adapter as follows:

$$Y^{\text{BF16}} = X^{\text{BF16}} \cdot \text{doubleDequant}(c_1^{\text{FP32}}, c_2^{\text{FP8}}, W^{\text{NF4}}) + X^{\text{BF16}} \cdot \frac{\alpha}{r} A^{\text{BF16}} B^{\text{BF16}}. \quad (3)$$

209 where $\text{doubleDequant}(\cdot)$ reconstructs the quantized weight matrix from its
210 two-level quantization scheme:

$$\text{doubleDequant}(c_1, c_2, W) = \text{dequant}(\text{dequant}(c_1, c_2, W)) = W^{\text{BF16}}. \quad (4)$$

211 In our configuration, we use the **NF4** quantization format for W , which pro-
212 vides a non-uniform 4-bit representation optimized for model performance,
213 and we apply *double quantization* to compress the scaling factors c_1 and
214 c_2 into low-bit formats. All matrix multiplications and LoRA updates are
215 performed in **bfloat16** precision to ensure numerical stability and efficient
216 utilization of modern GPUs.

During training, only the LoRA adapter parameters A and B are updated, while the quantized backbone weights remain frozen. The gradients with respect to the LoRA parameters are computed using the reconstructed W^{BF16} obtained from the double dequantization process. This design significantly reduces memory consumption while maintaining training stability and model quality.

3.2 Swift

SWIFT (Scalable lightWeight Infrastructure for Fine-Tuning) is an open-source [Zhao et al., 2025], unified framework designed to support the full lifecycle of Large Language Models (LLMs) including training, fine-tuning, evaluation, and deployment. A key feature of SWIFT is its modular and extensible architecture, which enables users to train large-scale models efficiently on limited hardware.

To facilitate efficient fine-tuning of the LLaMA 3 model, we incorporate components from the SWIFT framework. While our overall pipeline is based on Transformers, we specifically use SWIFT’s ”Swift” and ”TrainingArguments” modules to simplify configuration and the integration of parameter-efficient adapters. These modules guarantee reliable and streamlined LoRA integration and training setup for complex models, especially in 4-bit quantized form, without needing manual patching or specifying target modules.

3.3 Benchmark Comparison

We benchmark our QLoRA-fine-tuned LLaMA3-8B model against a suite of established financial sentiment analysis models, including FinBERT, VADER, and Financial RoBERTa, to assess performance across both classification accuracy and downstream trading signal generation. All models are evaluated on the same dataset to ensure consistency (our own test/out-of-sample dataset).

FinBERT is a BERT-based encoder-only model developed by ProsusAI [Araci, 2019], fine-tuned specifically for sentiment classification tasks in the financial domain. It was trained on a subset of the Financial PhraseBank, which contains manually annotated sentences from financial news. The model outputs a softmax distribution over three sentiment classes—positive, neutral, and negative—and is widely used for traditional text classification pipelines.

250 VADER (Valence Aware Dictionary and sEntiment Reasoner) [Hutto and Gilbert, 2014]
251 is a rule-based sentiment analysis tool optimized for social media and short
252 financial text. Unlike transformer-based models, VADER relies on a lexi-
253 con of sentiment-labeled words and empirically tuned heuristics to produce
254 a compound sentiment score ranging from -1 (most negative) to $+1$ (most
255 positive). For benchmarking purposes, we discretize this continuous output
256 into three sentiment categories using empirically chosen thresholds.

257 Financial RoBERTa, developed by Soleimani et al. [Soleimanian, 2023], is
258 a regression-based RoBERTa variant trained on financial text to predict a
259 continuous sentiment score in the range of -1 to $+1$. It is fine-tuned on
260 curated financial datasets using a mean-squared-error objective and demon-
261 strates strong correlation with market-relevant sentiment. Like VADER, its
262 output is also discretized into three sentiment classes for comparison with
263 other models.

264 This diverse set of models—spanning dictionary-based, regression-based, classification-
265 based, and generative paradigms—allows for a robust benchmarking of our
266 QLoRA-LLaMA3 model’s generalization and domain understanding capabil-
267 ities.

268 3.4 Trading Strategy

269 In this section, we focus on generating market-neutral strategies. Such
270 strategies isolate the alpha from the beta, and any subsequent returns cap-
271 tured are attributed to the model’s skills rather than the market movements.
272 Market-neutral strategies are also more robust under stress test and regu-
273 latory scrutiny. We also maintain a dollar-neutral position, which means
274 that our methods generally long-short the same amount of money (USD) to
275 lower the systematic risk (portfolio beta is near zero). From these stances,
276 we continue to build two trading mechanisms.

277 3.4.1 Score-based portfolio

278 To assess the practical profitability of our sentiment classification model, we
279 have developed a straightforward portfolio strategy based on daily news sen-
280 timent predictions [Chen et al., 2022]: score-based portfolio. For each stock,
281 we aggregate all predicted sentiment labels from the daily news to compute an
282 average sentiment score. Utilizing this score, we create a zero-net-investment

283 portfolio by taking equal-weighted long positions in the top quintile (15%) of
 284 stocks exhibiting the most positive sentiment, while simultaneously shorting
 285 the bottom quintile (15%) with the most negative sentiment. The portfolio
 286 is not rebalanced daily, as news coverage for all stocks may not be available
 287 every day. To address this issue, we establish a minimum threshold of 30 for
 288 the number of sentiment scores generated. If the total number of sentiment
 289 scores produced in a day exceeds this threshold, we consider that a trading
 290 day is effective for reconstructing our trading portfolio. This approach allows
 291 us to adjust our positions based on market conditions rather than making
 292 arbitrary changes, which could lead to increased explicit and implicit trading
 293 costs.

294 **3.4.2 K-means clustering trading strategy**

295 We design a market-neutral, static long-short trading strategy by clustering
 296 stocks based solely on their time-series sentiment patterns, without using
 297 any price or return information. Specifically, we construct a matrix of our-
 298 model's sentiment scores with tickers as rows and dates as columns, then
 299 apply KMeans clustering across the entire time period. This groups stocks
 300 into clusters based on how their sentiment evolves over time. We then com-
 301 pute an aggregate sentiment score for each ticker (e.g., average or cumulative
 302 sentiment), and rank entire clusters based on the average sentiment of their
 303 members.

304 From this ranking, we long the tickers in the top 20% of clusters (those
 305 with the most positive sentiment trajectories) and short the tickers in the
 306 bottom 20% of clusters (those with the most negative sentiment). These
 307 long and short positions are held unchanged throughout the entire period
 308 (no rebalancing). In particular, since the number of tickers may differ for the
 309 top and bottom clusters, we just long and short \$ 1 total for each position, i.e.
 310 if we have 120 tickers for short position, we only sell \$ 1, the same happens for
 311 the long position. Portfolio performance is evaluated by computing the daily
 312 spread return (long minus short), accumulating it over time, and analyzing
 313 metrics such as Sharpe ratio, max drawdown, and annualized volatility —
 314 all while strictly avoiding look-ahead bias.

315 3.4.3 Incremental score-based portfolio

316 Furthermore, [Kim et al., 2023] suggests that the (nearly) daily news can
317 often be noisy and unreliable for constructing our trading strategies. To
318 address this issue and in line with their analysis, we propose an incremental
319 score-based portfolio. The concept revolves around measuring changes in
320 sentiment scores between two consecutive effective trading days, allowing us
321 to capture the "spike" in market sentiment. To hit this level of thinking,
322 we forward fill all the null values in the sentiments column, reflecting the
323 unchanged market emotions towards any particular stock. This approach
324 helps us distinguish between noise and genuine market trends. Our ranking
325 method is straightforward; we subtract the sentiment score of a particular
326 stock on trading day i from that of trading day $i - 1$. We then rank these
327 increments by their values and select the top 20% for long positions and the
328 bottom 20% for short positions.

329 Additionally, we account for the lag in the societal impact of news; therefore,
330 we base our position on the news from yesterday and make trades today. We
331 will hold this position for 2 days to capture the directional return from our
332 factors.

333 3.4.4 Benchmark trading strategy

334 As the cornerstone of our trading strategy evaluation, we adopt a long-only
335 investment approach in the S&P 500 index as our benchmark. This bench-
336 mark serves as a representative proxy for broad market performance, en-
337 abling us to assess the relative value added by our sentiment-driven model.
338 By comparing against a passive, well-diversified index strategy, we establish
339 a robust baseline to evaluate whether our active strategy generates excess
340 return (alpha) beyond general market trends.

Algorithm 1 Score-Based Portfolio Construction

```
1: for each day  $t$  do
2:   Extract all sentiment signals: sent[day t]
3:   if number of valid (non-NaN) scores  $\geq 30$  then
4:     Sort stocks by sentiment score
5:     Select:
      • Top 20%  $\rightarrow$  long $1
      • Bottom 20%  $\rightarrow$  short $1
6:     Update current_longs and current_shorts
7:   end if
8:   Calculate P&L on day  $t + 1$ 
9: end for
```

Algorithm 2 Incremental Score-Based Portfolio Construction

```
1: for each trading day  $t$  do
2:   if both  $t$  and  $t-1$  are effective trading days (i.e., sentiment signals  $\geq 30$ ) then
3:     Forward-fill missing sentiment values up to day  $t$ 
4:     Extract sentiment scores: sent[day t] and sent[day t-1]
5:     Compute sentiment increment: delta = sent[day t] - sent[day t-1]
6:     Sort stocks by delta
7:     Select:
      • Top 20%  $\rightarrow$  long 5% of capital
      • Bottom 20%  $\rightarrow$  short 5% of capital
8:     Execute positions based on sentiment from day  $t-1$  on day  $t$ 
9:     Hold positions for 2 trading days
10:    Calculate P&L for the 2-day holding period
11:   end if
12: end for
```

341 4 Experiments and Results

342 4.1 Experiments

343 We use Google Colab to train and test our model. We train our model based
344 on four financial text datasets, which cost about 18 hours, with 90 computer
345 units. Then we compare the classification precision with another financial
346 sentiment llm, FinBert. The results shows that our model can classify the
347 sentiment more precisely. The portfolio return also reflects that our model
348 is profitable.

349 4.1.1 Parameters

350 **Low-Rank Adaptation (LoRA) Configuration:** LoRA parameters are
351 set with a rank $r = 16$, `lora_alpha = 32`, and `lora_dropout = 0.1`.

352 **Learning Rate and Optimizer:** The model is fine-tuned using a learning
353 rate of 1×10^{-4} with the AdamW optimizer.

354 **Batch Size and Epochs:** Training is conducted over 1 epoch with a batch
355 size of 8 and gradient accumulation set to 4.

356 **Training Details:** We set the ratio 9:1 of the train set to the test set;
357 Gradient checkpointing is enabled to reduce memory usage during backprop-
358 agation, which facilitates training with limited GPU resources. A cosine
359 learning rate scheduler is used to gradually decay the learning rate over time,
360 enhancing convergence and stability.

361 4.2 Evaluation Matrices

362 4.2.1 Classification evaluation metrics

363 We use the micro-level F1 score, precision, recall to measure the model’s
364 predictability on a given class. As the model varies, the sentiment dataset
365 may be highly imbalanced, this is where we introduce the macro-level average
366 and weighted average indices in the classification reports in tables 1 and 2.

367 4.2.2 Portfolio’s performance evaluation metrics

368 How do we evaluate the performance of our trading strategies? For our
369 market-neutral approaches, it is useful to compare metrics such as the Sharpe
370 ratio, maximum drawdown, annualized return, and volatility. Since we gen-
371 erate directional factors, adding the hit ratio as an evaluation metric is also
372 beneficial. We rebalance our long and short positions daily. If a long position
373 gains or a short position (including reversed positions) profits, it counts as a
374 ”hit”; otherwise, it’s a ”miss.” We then compute the ratio of hits to total hits
375 and misses. This provides a clear measure of the model’s ability to predict
376 market directions.

377 4.3 Results

378 4.3.1 Models’ performances on sentiments classification tasks

379 We evaluated our fine-tuned Meta-LLaMA3-8B sentiment classifier on a test
380 set of 9,064 financial news samples. From table 1, our model achieved an
381 impressive overall accuracy of 92.18%, with a macro-average $F1$ score of
382 0.5787 and a micro-average $F1$ score of 0.9218. The class-wise performance
383 reveals that the model performs exceptionally well on Positive and Neutral
384 samples, with $F1$ scores of 0.9636 and 0.7172, respectively. However, its
385 performance on the Negative class remains modest ($F1 = 0.0552$), likely due
386 to class imbalance within the training data.

Table 1: Classification Performance of Our QLoRA-tuned LLaMA3-8B Model

Class	Precision	Recall	F1-Score	Support
Negative	0.14	0.03	0.06	144
Neutral	0.61	0.88	0.72	983
Positive	0.99	0.94	0.96	7937
Accuracy	0.92			
Macro Avg	0.58	0.62	0.58	9064
Weighted Avg	0.93	0.92	0.92	9064

387 For baseline comparison (in Table 2), we evaluated FinBERT [?] on the same
388 test set. Although FinBERT achieved high precision for the Positive class
389 (0.8232), its overall recall was significantly lower, resulting in a $F1$ score of

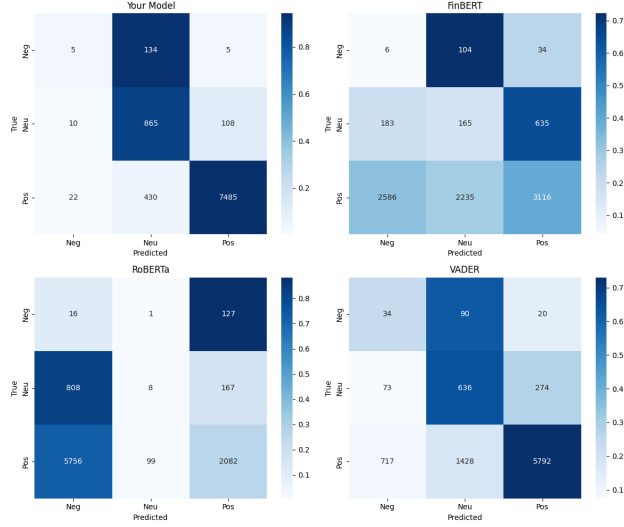


Figure 2: Confusion Matrices of Sentiment Classifier.

only 0.5316 for that class. Performance in the negative and neutral classes was notably poor, with an overall accuracy of just 36.26% and a macro $F1$ of 0.2101.

For baseline comparison, we evaluated the rule-based sentiment analyzer VADER [Hutto and Gilbert, 2014] on the same test set. VADER achieved strong precision on the Positive class (0.9466), but significantly underperformed on the Negative class ($F1$ score of only 0.0659) and the Neutral class ($F1$ of 0.4101). The model exhibited a tendency to overpredict Positive sentiment, resulting in an imbalanced performance. Overall accuracy was 71.05%, while macro $F1$ remained modest at 0.4303.

We also evaluated the Financial RoBERTa model [Soleimanian, 2023], which outputs continuous sentiment scores subsequently discretized into three classes. Although the model showed high precision for the Positive class (0.8800), its recall was relatively poor (0.2579), leading to an $F1$ score of 0.3991 for that class. Performance on the Negative and Neutral classes was notably weak, with near-zero $F1$ scores. The overall accuracy reached just 23.11%, and macro $F1$ dropped to 0.1436, suggesting poor generalization beyond majority sentiment prediction.

Additionally, we schematically present the confusion matrices of the four models in Figure 2. From the figure, it is evident that our LLaMA3-based classifier substantially outperforms the baselines across all three sentiment classes. Notably, our model achieves strong performance in the neutral

Table 2: Performance Comparison of Baseline Models (without Support column)

Model	Class	Precision	Recall	F1-Score
VADER	Negative	0.04	0.24	0.07
	Neutral	0.30	0.65	0.41
	Positive	0.95	0.73	0.83
	Accuracy		0.71	
	Macro Avg	0.43	0.54	0.43
	Weighted Avg	0.87	0.71	0.77
FinBERT	Negative	0.00	0.04	0.00
	Neutral	0.07	0.17	0.09
	Positive	0.82	0.39	0.53
	Accuracy		0.36	
	Macro Avg	0.30	0.20	0.21
	Weighted Avg	0.73	0.36	0.48
RoBERTa	Negative	0.00	0.11	0.00
	Neutral	0.07	0.01	0.01
	Positive	0.88	0.26	0.40
	Accuracy		0.23	
	Macro Avg	0.32	0.13	0.14
	Weighted Avg	0.78	0.23	0.36

class, with minimal confusion between neutral and positive instances, a challenge that significantly impacts FinBERT and RoBERTa. FinBERT shows a strong bias toward predicting the positive class, with over 2,500 Negative and 600+ Neutral examples misclassified as Positive. Similarly, RoBERTa heavily overpredicts the positive class while failing to recognize Neutral and Negative sentiments. VADER performs comparatively better than FinBERT and RoBERTa in distinguishing Neutral from Positive, but still exhibits class imbalance and reduced granularity. Although the Negative class remains underrepresented and challenging for all models, our fine-tuned LLaMA3 classifier demonstrates higher robustness and generalization, making it a reliable choice for downstream financial sentiment analysis.

4.3.2 Models' performances on S&P 500 companies' news

Histogram of news sentiments Following Figure 3, we can have a glimpse of the general sentiment prediction distribution for all four models.

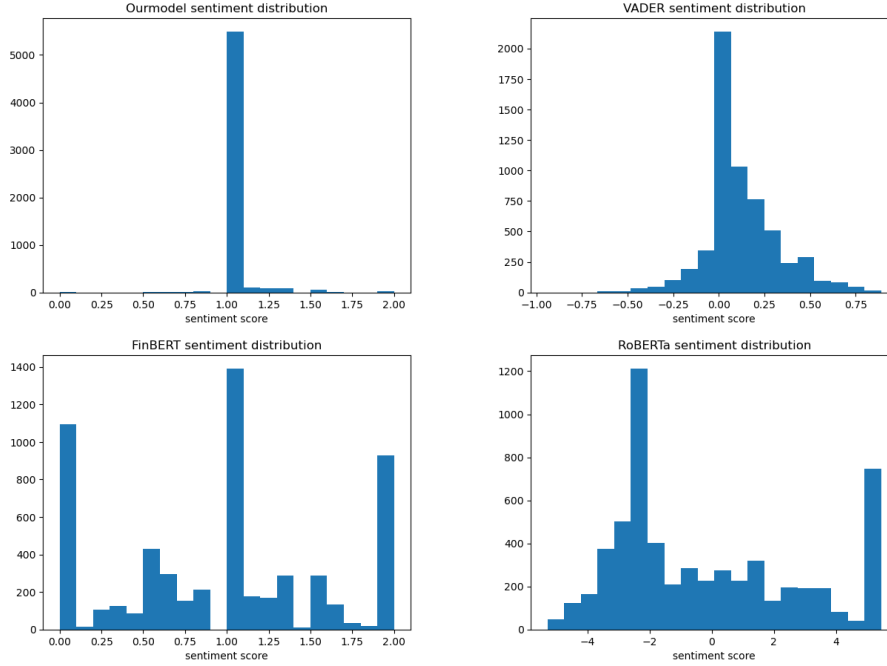


Figure 3: Histograms of sentiments for four different models.

Compared to the other three models, our model generally adopts a more conservative approach, typically predicting near-neutral sentiment. In contrast, FinBERT tends to forecast more extreme sentiments more aggressively. RoBERTa shows a tendency toward positive predictions but mostly assigns negative sentiments to most news. The models' performances vary across different news datasets. Although our dataset is limited to only 15 months, this isn't enough to definitively determine whether our model is better or worse than the others. Additionally, we evaluate how these sentiments influence trading outcomes.

Performance on score-based trading strategies As Figure 4 and Table 3 show, the VADER and ourmodel in this score-based method outperform the benchmark SPY almost over the whole period (except for 2024-11 to 2025-03), achieving higher cumulative returns with relatively stable trajectories and resulting in higher return at the end (2025-05). In particular,

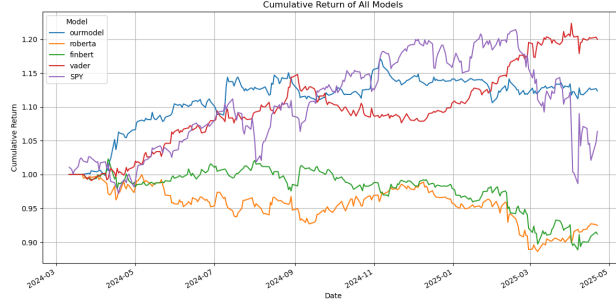


Figure 4: Cumulative return of the score-based trading strategy.

Table 3: Performance Metrics of Sentiment-Based Trading Models (Score-based)

Metric	OurModel	Roberta	FinBERT	VADER
Annualized Return	11.06%	-6.80%	-7.95%	17.84%
Annualized Volatility	8.28%	9.67%	9.33%	8.58%
Sharpe Ratio	1.31	-0.68	-0.84	1.96
Max Drawdown	-5.50%	-11.37%	-13.14%	-6.23%
Hit Ratio	49.45%	48.71%	49.26%	50.00%

our model demonstrates a stable profit accumulation, yielding an annualized return of more than 11%, a Sharpe ratio of over 1.3, while processing a relatively low annualized volatility of 8.28% and a maximum drawdown of -5.50%, compared to all other models. This demonstrates the profitability of our model within this trading strategy, even though it predicts in a relatively conservative way. We believe that this model’s backward-looking sentiment predictions would greatly benefit a more hype-affluent market sector.

On the other hand, VADER performs exceptionally well, even better than our model, with a Sharpe ratio of more than 1.90. We believe its profitability can be attributed to its relatively conservative forecasting, and its fatter tail captures more market momentum than our model. However, in the subsequent trading scheme, you will find that this model ultimately proves less prosperous.

Performance on incremental score-based trading strategies All four models perform in a highly similar pattern in this trading strategy. As figure 4 verifies our incremental-score mechanism has filtered out the majority amount of noise in the forecasts. While OurModel demonstrated robust

Table 4: Performance Metrics of Sentiment-Based Models with 2-Day Holding Period, Same-Day Entry

Metric	OurModel	Roberta	FinBERT	VADER
Annualized Return	4.17%	2.29%	2.64%	2.48%
Annualized Volatility	4.57%	4.83%	4.60%	4.51%
Sharpe Ratio	0.92	0.49	0.59	0.57
Max Drawdown	-4.51%	-4.07%	-3.25%	-4.12%
Hit Ratio	49.63%	49.81%	50.93%	50.00%

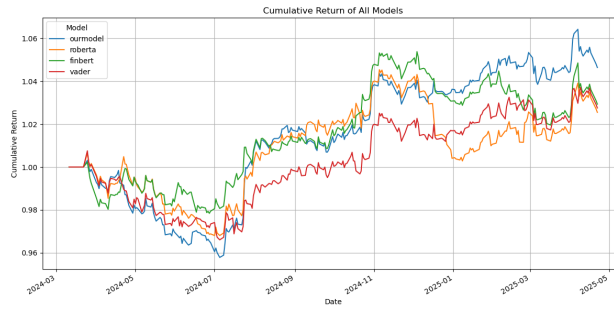


Figure 5: Cumulative return of the incremental score-based trading strategy.

consistency and risk-adjusted performance across various trading regimes, it showed limitations in maximizing raw returns. In particular, it was often outperformed by VADER in terms of absolute annualized return, suggesting its signals may be more conservative and less aggressive in capturing strong directional moves. Additionally, its hit ratio hovered around 49–50%, indicating that the model’s predictive accuracy is only marginally better than random chance, relying heavily on position sizing and risk controls to generate profit.

This modest predictive power could present challenges if market volatility or trend patterns change, potentially exposing the model to suboptimal trade entries or exits in less stable environments. Further improvements in signal quality, for example through feature engineering or ensemble calibration, could help OurModel increase its raw returns without compromising its strong drawdown control.

Performance on K-means clustering (cluster number = 7) trading strategies As shown in Figure 6, the trading strategy derived from our model outperforms the one based on FinBERT when using sentiments

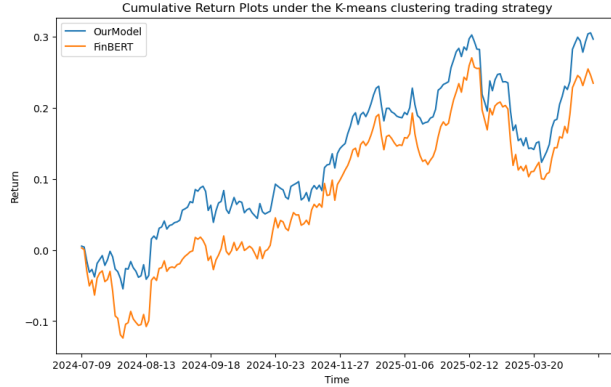


Figure 6: Cumulative return of the incremental score-based trading strategy.

generated by our approach. Specifically, our model achieves a Sharpe ratio of 1.686, compared to 1.280 for FinBERT. Both models have similar maximum drawdowns, around 0.17. While our model clusters 400 tickers into the bottom portfolio, FinBERT only uses 14 tickers. Nonetheless, our approach yields higher profits on the same news dataset.

5 Conclusion

In this project, we investigated how parameter-efficient fine-tuning techniques can be applied to financial sentiment analysis, focusing on Quantized LoRA (QLoRA) with the Meta-LLaMA3-8B model. By combining multiple financial sentiment datasets and utilizing the SWIFT framework, we trained a domain-specific model capable of producing high-quality sentiment signals from financial news headlines.

Our results indicate that the QLoRA-tuned LLaMA3-8B model achieves over 92% accuracy in classification, with strong precision and recall for neutral and positive sentiments. In trading scenarios, it showed consistent profitability and favorable risk-adjusted returns, outperforming traditional models in metrics like the Sharpe ratio and maximum drawdown under score-based strategies. Even in the incremental score-based and K-means clustering approach, our model remained competitive, though its conservative nature slightly limited raw returns compared to VADER.

These findings highlight the promise of quantized, parameter-efficient fine-tuning of large language models for financial tasks—striking a good balance

496 between performance, efficiency, and ease of deployment. Future work could
497 improve this framework by adding continuous sentiment scoring, integrating
498 multi-modal inputs such as charts or analyst transcripts, and applying rein-
499 forcement learning to directly optimize trading signals. Exploring ensemble
500 techniques and better handling of class imbalance, especially for negative
501 sentiments, could also boost the model’s robustness.

502 Looking ahead, extending this approach to generate continuous sentiment
503 scores instead of discrete classifications could offer more nuanced signals.
504 Incorporating broader market context features and leveraging reinforcement
505 learning to directly optimize trading strategies are promising avenues for
506 enhancing both prediction accuracy and economic gains.

507 **Ethics statement**

508 This study did not involve human participants, animals, or sensitive data.
509 Therefore, no ethical approval was required.

510 **Author contributions**

511 FSL devoted a majority of work in this paper. He designed the study, wrote
512 the manuscript, and lead the direction of researches with PGJ; the LLM is
513 trained by FSL with the help of CJM, who contributes to the detailed study
514 of SWIFT packages.

515 All authors reviewed and approved the final manuscript.

516 **Funding**

517 This research received no external funding. The authors did not receive fi-
518 nancial support from any public, commercial, or non-profit funding agencies.

519 Conflict of interest

520 The authors declare that the research was conducted in the absence of any
521 commercial or financial relationships that could be construed as a potential
522 conflict of interest.

523 Acknowledgments

524 The authors thank Meta for releasing the LLaMA model as open-source
525 software. No financial support, sponsorship, or research funding was received
526 from Meta or any other organization.

527 Data availability statement

528 The datasets and code used in this study are available at: [https://github.](https://github.com/RomeoisFushengLuo/TradingStrategies_NewsSentimentAnalysis_FinLlama3_8B)
529 [com/RomeoisFushengLuo/TradingStrategies_NewsSentimentAnalysis_FinLlama3_](https://github.com/RomeoisFushengLuo/TradingStrategies_NewsSentimentAnalysis_FinLlama3_8B)
530 [8B](https://github.com/RomeoisFushengLuo/TradingStrategies_NewsSentimentAnalysis_FinLlama3_8B).

531 References

- 532 [Araci, 2019] Araci, D. (2019). Finbert: Financial sentiment analysis with
533 pre-trained language models.
- 534 [Chen et al., 2022] Chen, Y., Kelly, B. T., and Xiu, D. (2022). Expected
535 returns and large language models. Available at SSRN 4416687.
- 536 [Chiu and Hung, 2024] Chiu, I.-C. and Hung, M.-W. (2024). Finance-
537 specific large language models: Advancing sentiment analysis and return
538 prediction with llama 2. *arXiv preprint arXiv:2402.08887*.
- 539 [Dettmers et al., 2023] Dettmers, T., Pagnoni, A., Holtzman, A., and Zettle-
540 moyer, L. (2023). Qlora: Efficient finetuning of quantized llms. *arXiv*
541 *preprint arXiv:2305.14314*.
- 542 [Hu et al., 2022] Hu, E. J. et al. (2022). Lora: Low-rank adaptation of large
543 language models. *arXiv preprint arXiv:2106.09685*.

544 [Hutto and Gilbert, 2014] Hutto, C. J. and Gilbert, E. (2014). Vader: A par-
545 simonious rule-based model for sentiment analysis of social media text. In
546 *Proceedings of the Eighth International Conference on Weblogs and Social*
547 *Media (ICWSM)*, Ann Arbor, MI.

548 [Iacovides et al., 2024] Iacovides, G., Konstantinidis, T., Xu, M., and
549 Mandic, D. (2024). Finllama: Llm-based financial sentiment analysis for
550 algorithmic trading. *arXiv preprint arXiv:2401.07874*.

551 [Kim et al., 2023] Kim, S., Kim, S., Kim, Y., Park, J., Kim, S., Kim, M.,
552 Sung, C. H., Hong, J., and Lee, Y. (2023). Llms analyzing the analysts:
553 Do bert and gpt extract more value from financial analyst reports? In *Pro-*
554 *ceedings of the Fourth ACM International Conference on AI in Finance*,
555 pages 383–391.

556 [Soleimanian, 2023] Soleimanian, M. (2023). Financial roberta: A roberta-
557 based language model fine-tuned for financial sentiment analysis. Accessed:
558 2025-06-19.

559 [S&P Dow Jones Indices, 2025] S&P Dow Jones Indices (2025). S&p 500
560 index. Accessed: 2025-06-25.

561 [Zhao et al., 2024] Zhao, J., Wang, T., Abid, W., Angus, G., Garg, A., Kin-
562 nison, J., Sherstinsky, A., Molino, P., Addair, T., and Rishi, D. (2024).
563 Lora land: 310 fine-tuned llms that rival gpt-4: A technical report. *arXiv*
564 *preprint arXiv:2406.XXXXX*.

565 [Zhao et al., 2025] Zhao, Y. et al. (2025). Swift: A scalable lightweight in-
566 frastructure for fine-tuning. In *Proceedings of the AAAI Conference on*
567 *Artificial Intelligence*, volume 39.