

FinLlama: LLM-Based Financial Sentiment Analysis for Algorithmic Trading

Anonymous ACL submission

Abstract

Online sources of financial news have a profound influence on both market movements and trading decisions. Standard sentiment analysis employs a lexicon-based approach to aid financial decisions, but struggles with context sensitivity and word ordering. On the other hand, Large Language Models (LLMs) are powerful, but are not finance-specific and require significant computational resources. To this end, we introduce a finance specific LLM framework, based on the Llama 2 7B foundational model, in order to benefit from its generative nature and comprehensive language manipulation. Such a generator-discriminator scheme, referred to as FinLlama, both classifies sentiment valence and quantifies its strength, offering a nuanced insight into financial news. The FinLlama model is fine-tuned on supervised financial sentiment analysis data, to make it handle the complexities of financial lexicon and context, and is equipped with a neural network-based decision mechanism. The subsequent parameter-efficient fine-tuning optimises trainable parameters, thus minimising computational and memory requirements without sacrificing accuracy. Simulation results demonstrate the ability of FinLlama to increase market returns in portfolio management scenarios, yielding high-return and resilient portfolios, even during volatile periods.

1 Introduction

The ever increasing prominence of algorithmic trading in quantitative finance has necessitated the need for reliable and actionable AI-aided domain knowledge from vast streams of data with multiple modalities. Of particular interest is generative AI, owing to its ability to distill insights from non-numerical sources such as news, earnings calls, financial reports, and other textual sources. In this context, sentiment analysis from text promises to bridge the gap between market movements caused by geopo-

litical and socioeconomic events, human actions, and quantitative trading.

The sentiment contained in on-line textual sources can drive market movements; such information harbours intrinsic advantages and gives a competitive edge to those equipped with the tools to harness it. Sentiment analysis rests upon the quantification of opinions present in unlabeled textual data, and aims to categorize whether the overall perspective is positive, negative, or neutral. When applied to large-scale information sources, this promises to enhance the understanding for the overall direction of macroscopic trends, a task which is both challenging and time-consuming for human analysts.

Despite conceptual benefits, the diverse, nuanced, and vast nature of financial text presents unique challenges when it comes to extracting sentiment in a manner that is both accurate and actionable. For example, the words ‘bull’ and ‘bear’ are neutral in the general vocabulary, but in financial markets, their respective connotations are strictly positive or negative (Mishev et al., 2020). This highlights the need for context-aware sentiment extraction, and underpins the complexities of employing natural language processing (NLP) in financial applications.

To address these issues, we consider the following fundamental questions:

- Can large language models (LLMs), which have already revolutionized manifold areas of NLP, be specifically tailored for sentiment analysis in the finance domain, particularly for enhancing algorithmic trading?
- Can this be achieved in a way which does not require vast computational resources, typically associated with NLP models, thus making the approach accessible to anyone equipped with standard computational resources?

081 .
082 Our proposed solution, termed *FinLlama*, is
083 is obtained by fine-tuning a pre-trained LLM
084 (namely Llama 2 7B (Touvron et al., 2023)) on spe-
085 cialised, labelled and publicly available financial
086 news datasets. The ultimate goal of FinLlama is
087 to enhance the performance of financial sentiment
088 analysis, whilst leveraging on parameter-efficient
089 fine-tuning (PEFT) and 8-bit quantization, through
090 LoRA (Hu et al., 2021), to minimise resource re-
091 quirements.

092 The main contributions of this work are:

- 093 • **Targeted fine-tuning:** Rather than utilising
094 one general LLM for financial tasks, our ap-
095 proach capitalizes on the foundational pre-
096 trained Llama 2 model, whereby fine-tuning
097 is performed specifically for the purpose of
098 sentiment classification through a SoftMax
099 classification layer at its output.
- 100 • **Efficient resource utilization:** Our approach
101 ensures that even standard computational re-
102 sources, with no high-end GPUs, can be em-
103 ployed. By virtue of the pre-trained Llama
104 2 model and through targeted parameter-
105 efficient fine-tuning, computational demands
106 are dramatically reduced compared to the ex-
107 isting methods, thus bridging the gap between
108 academic benchmarks and practical utility.
- 109 • **Benchmarking and real-world application:**
110 The success of fine-tuned LLMs for finance
111 has also highlighted that these have not yet
112 adequately addressed the domain of portfolio
113 construction. To this end, we integrate the
114 extracted sentiment signals by FinLlama into a
115 long-short portfolio, which allows us to obtain
116 finance-specific real-world metrics including
117 cumulative returns and the Sharpe ratio.

118 2 Related Work

119 The potential of sentiment analysis in finance was
120 first recognised in 1970 by Eugene Fama who in-
121 troduced the Efficient Market Hypothesis (EMH)
122 (Fama, 1970), which states that stock prices change
123 in response to unexpected fundamental informa-
124 tion and news. In this context, before the intro-
125 duction of advanced machine learning tools, the
126 financial sector has employed lexicon-driven ap-
127 proaches (Mishev et al., 2020). These methods
128 analyse textual content, sourced from news arti-
129 cles or financial disclosures, based on specific key-

words, which are then linked to established sen-
130 timent ratings (Li et al., 2014; Ke et al., 2019a).
131 However, an exponential increase in the volume
132 and richness of online available information posed
133 significant challenges for lexicon-based analysis,
134 but has opened a fertile ground for machine learn-
135 ing strategies, including techniques such as Naive
136 Bayes and Support Vector Machines (Cristianini
137 and Shawe-Taylor, 2000), as summarised in Figure
138 1.
139

140 In parallel, the advances in deep learning have
141 become instrumental for NLP research and have
142 spurred pioneering works that sought to harness
143 the power of neural networks for NLP tasks. Re-
144 cently, the introduction of the attention mechanism
145 and transformer networks has enabled a significant
146 shift away from recurrent and convolutional meth-
147 ods, traditionally used in deep-learning tasks (Yang
148 et al., 2016). This has led to the development of
149 transformer-based models, such as BERT (Devlin
150 et al., 2019), which owing to its contextual com-
151 prehension of language has been used extensively
152 for sentiment analysis. However, the performance
153 of BERT in the financial domain has encountered
154 limitations, primarily because it is not specifically
155 trained on financial datasets. Moreover, its require-
156 ment for substantial amounts of data for fine-tuning
157 purposes poses a considerable challenge for finan-
158 cial applications, where such data may not be read-
159 ily available.

160 More recently, FinBERT (Araci, 2019), a ver-
161 sion of BERT which is fine-tuned on financial text,
162 has shown promising results for the task of finan-
163 cial sentiment analysis. However, FinBERT still
164 suffers from limitations such as insensitivity to nu-
165 merical values, while due to its relatively small size
166 (110 million parameters) its classification accuracy
167 deteriorates with sentence complexity (Chen et al.,
168 2023). The FinGPT (Liu et al., 2023; Yang et al.,
169 2023) and Instruct-FinGPT (Zhang et al., 2023)
170 aim to enhance their expressive power by using the
171 Llama 7B as their base model. However, FinGPT
172 is not optimized for the task of financial sentiment
173 analysis whilst Instruct-FinGPT only classifies the
174 sentiment valence but is not capable of quantifying
175 the strength of a sentiment class.

176 To the best of our knowledge, BloombergGPT
177 (Wu et al., 2023) is the only pre-trained finance-
178 specific LLM, as Bloomberg was able to train the
179 model using data collected over a span of 40 years.
180 Despite the impressive performance of the model
181 on financial sentiment analysis, the resources re-

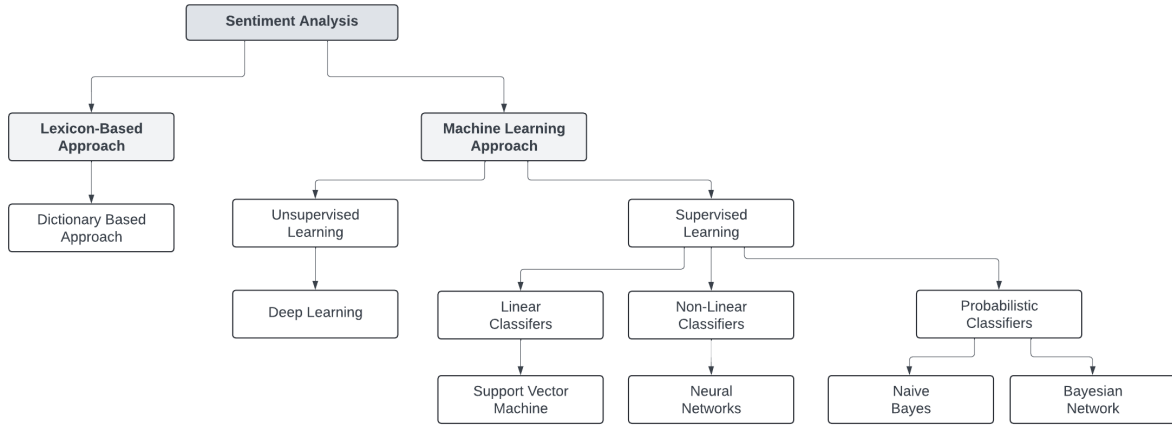


Figure 1: Overview of sentiment analysis methods.

quired to train such a model are substantial (1.3M GPU hours at a cost of \$5M) whilst much of the training data is confidential and not publicly available. This is different from our proposed methodology, which focuses on achieving a high classification accuracy whilst minimizing the training corpus and computational resources, and utilizing publicly available training data. This is achieved by fine-tuning a pre-trained general-purpose LLM on a smaller-scale financial data corpus.

3 Methodology

Our work aims to embark upon the immense expressive power and contextual understanding of general-purpose LLMs in order to make them finance-specific. This is achieved by fine-tuning the state-of-the-art (SOTA) Llama 2 7B model on a finance-specific corpus of online data. The effectiveness of our approach is demonstrated on financial sentiment analysis through a new set of benchmarks that align closely with end portfolio construction – the ultimate goal of financial analysis.

3.1 Fine-tuning the Llama 2 model

Even though pre-trained LLMs offer a range of capabilities such as reasoning, translation, summarising and text generation, they often struggle when applied to a specific task of interest, such as sentiment analysis. This limitation becomes even more critical in the finance domain, where the nuanced language, media hype and extensive length of financial news articles pose significant challenges.

To tackle these challenges, our work revisits the first principles of LLMs in order to align them to the task of financial sentiment analysis. This

is achieved by using four labelled financial text datasets as training data to fine-tune the Llama 2 model. Such finance-specific training equips the model with the ability to understand the linguistic nuances present in the financial domain. Furthermore, a three-class SoftMax classification layer is employed at the output of the foundational model. This made it possible to alter the primary function of the LLM from text generation to sentiment classification. In this way, the proposed fine-tuned FinLlama model acts as a generator-discriminator and produces sentiment decision outputs for three labels: positive, negative or neutral.

3.1.1 Training datasets

The training data was a combination of four labelled publicly available financial news datasets, namely the Financial PhraseBank (FPB) dataset (Malo et al., 2014), FiQA dataset (Maia et al., 2018), Twitter Financial News dataset (Wang, 2023) and GPT-labelled Financial News dataset (Magic, 2022). This resulted in a comprehensive collection of 34,180 labelled samples, as outlined below.

- **Financial PhraseBank (FPB) Dataset.** This dataset, accessible via HuggingFace, consists of 4,840 samples which are randomly extracted from financial news articles. In order to ensure high quality annotation, the samples were annotated by 16 experts with backgrounds in finance and business. Each sample was annotated with one of the three labels: positive, negative, and neutral.
- **FiQA Dataset.** This dataset is also accessible via HuggingFace and consists of 1,210

labelled sentences. Each sentence was annotated with one of the three labels: positive, negative, and neutral.

- Twitter Financial News Sentiment.** This dataset, accessible via HuggingFace, includes 11,930 tweets with content from the financial domain. Each tweet was annotated as positive, negative, and neutral.
- GPT-labelled Financial News.** This dataset, accessible via HuggingFace, consists of 16,200 financial news articles labelled by GPT-3.5. Each article was annotated with one of the five labels: strongly negative, mildly negative, neutral, mildly positive, and strongly positive. To align this dataset with the three-class output of our FinLlama model, the strongly and mildly negative classes were combined into a single negative class, and similarly, the strongly and mildly positive classes were combined into a single positive class.

3.1.2 Model Training

The proposed FinLlama model was first initialised with the Llama 2 7B model, followed by fine-tuning over 5 epochs. The training process utilised the AdamW optimizer (Loshchilov and Hutter, 2017), as it effectively decouples the weight decay from the optimization steps, leading to more effective training. The initial learning rate was deliberately kept small as the Llama 2 7B model is already pre-trained on a large corpus of data, whilst the warm-up ratio and weight decay served as key regularisation techniques to prevent overfitting, a crucial aspect given the limited size of our fine-tuning dataset.

Moreover, the LoRA implementation was employed in the fine-tuning process with a rank, $r = 8$, a scaling factor, $\alpha = 16$, and a dropout of 0.05, in order to minimize the number of trainable parameters whilst achieving high and robust end performance. Through the LoRA implementation, the number of trainable parameters was set to 4.2M, amounting to just 0.0638% of the total number of parameters in the Llama 2 7B model. This made it possible for our fine-tuning process to be **implemented on a single A100 (40 GB) GPU**, thus avoiding the need for excessive computational resources. A summary of the most important training parameters used in the fine-tuning process is given in Table 1.

3.2 Proposed Framework

After establishing the proposed fine-tuned Llama 2 model, we followed the framework shown in Figure 2, with the aim of assessing the performance of our FinLlama model against other established sentiment analysis methods, using finance-specific real-world metrics.

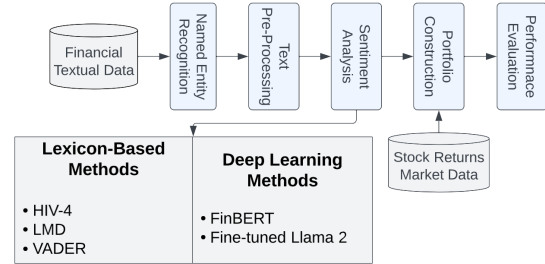


Figure 2: Framework for sentiment analysis.

Data Collection and Processing. Both textual and market data were analysed in order to construct appropriate long-short (L/S) portfolios. Regarding the textual data, 204,017 articles dating between 2015 to 2021 were collected from online sources such as Reuters, The Motley Fool and MarketWatch. These sources were selected due to their reliability, reputation, lack of bias and focus on major corporations. Financial market data were collected for the same time period from Yahoo Finance. These market data contained daily stock returns for the 500 companies in our Investable Universe (S&P 500), resulting in 1,672 days of stock returns data for each company. Data processing in the form of Named Entity Recognition (NER) and text pre-processing was then applied to the textual data, to remove irrelevant articles and ensure the compatibility of the articles with our sentiment methods.

Sentiment Analysis. In total, five sentiment analysis methods were applied. For the lexicon-based approaches (see Appendix A.1), LMD (Loughran and McDonald, 2011) and HIV-4 (Stone et al., 1966) were implemented using the pysentiment2 Python library, while VADER (Hutto and Gilbert, 2015) was implemented using the NLTK library. Regarding the deep learning methods (see Appendix A.2), both the FinBERT model and our FinLlama model were obtained through HuggingFace, and were utilised via the Transformers library.

The considered methods were evaluated on every article within each corpus for a given company. In

Parameter	Definition	Value
Learning rate	Determines the step size at each iteration of gradient descent	0.0003
Weight decay	Regularization technique to prevent overfitting by penalizing large weights	0.01
Batch size	Number of training samples used in one iteration of gradient descent	128
Training epochs	A full training pass over the entire training set	5
LR scheduler	Framework that adjusts the learning rate between iterations	Cosine Annealing
Warmup ratio	Increases the learning rate gradually over a certain number of epochs	0.1
GPUs	Number of GPUs used	1 A100 (40GB)
LoRA rank	Defines the dimensions of low-rank matrices	8
LoRA alpha	Scaling factor for the weight matrices within LoRA	16
LoRA dropout	Proportion of randomly deactivated neurons during training	0.05

Table 1: Training parameters used in the fine-tuning process of the proposed FinLlama.

cases where multiple articles were published on the same day for a given company, the average sentiment for that day was calculated as

$$S_t = \frac{1}{N_t} \sum_{i=1}^{N_t} S_{it} \quad (1)$$

Here, S_t represents the average sentiment for the t -th day, N_t denotes the number of news articles published on that same t -th day for a given company, while S_{it} designates the sentiment strength of the i -th news article on a particular t -th day. The daily sentiment outputs for each company were merged to arrive at the final sentiment data that were utilised as a parameter in the portfolio construction stage.

Portfolio Construction. Once the sentiment for each method was defined for every company, the long-short portfolio was constructed. We used the sentiment as a parameter to determine which companies should be in a long or a short position, aiming to maximise returns from both positions. The long-short portfolio was constructed using the following procedure:

- *Define the Investable Universe:* Even though the S&P 500 comprises 500 companies, the financial textual data collected did not contain articles associated to some of the companies for the test period of February 2015 to June 2021. Consequently, 417 companies were considered.
- *Define the long and short position:* The sentiment signal obtained from each of the five methods was used to construct five distinct portfolios. For each method, companies were ranked daily according to their sentiment. Companies that did not have sentiment data on a particular day were omitted from the ranking. As the daily sentiment score for each

company ranges between -1 and 1, those with the highest positive sentiment were placed in a long position, whilst those with the strongest negative sentiment were placed in a short position.

- *Allocation:* An equally-weighted portfolio strategy was considered in our portfolio construction as this strategy is mostly utilised by hedge funds (Ke et al., 2019b). The percentage of companies in a long and short position was fixed at 35%. Consequently, the top 35% of companies in terms of performance were allocated to long positions, while the bottom 35% were allocated to short positions.
- *Determine daily returns:* The daily return for each company that was held in a long or short position was obtained by the market data on that particular day. The average daily return of companies that were held in a long position, r_{Long} , was defined as

$$r_{Long} = \frac{1}{N_{Long}} \sum_{i=1}^{N_{Long}} r_{Long}(i) \quad (2)$$

Similarly, the average daily return of companies that were held in a short position, r_{Short} , was defined as

$$r_{Short} = \frac{1}{N_{Short}} \sum_{i=1}^{N_{Short}} r_{Short}(i) \quad (3)$$

For each particular day, the number of companies that were held in either a long position (N_{Long}) or a short position (N_{Short}) were equal. Consequently, the total portfolio return on a particular day was the difference between the daily long return, $r_{Long}(i)$, and daily short return, $r_{Short}(i)$, and is given by

$$r_{daily}(i) = r_{Long}(i) - r_{Short}(i) \quad (4)$$

Portfolio Evaluation. The performance of the portfolio constructed using our fine-tuned model was assessed against the portfolios constructed using other SOTA sentiment methods. To this end, the employed real-world financial metrics were: cumulative returns, r_{cum} , annualized return, R_p , annualized volatility, σ_p , and the Sharpe ratio, S_a (Berk and DeMarzo, 2019), defined as

$$r_{cum} = \sum_{i=1}^N r_{daily}(i) \quad (5)$$

$$R_p = \frac{1}{N} \sum_{i=1}^N r_{log}(i) \times 252 \quad (6)$$

$$\sigma_p = \sqrt{\frac{\sum_{i=1}^N (r_{log}(i) - \bar{r})^2}{N - 1}} \times \sqrt{252} \quad (7)$$

$$S_a = \frac{R_p - R_f}{\sigma_p} \quad (8)$$

where N is the total number of investing days, totaling 1,672, $r_{log}(i)$ represents the logarithmic daily return, \bar{r} denotes the average daily logarithmic return, R_f designates the annualized risk-free rate of return, and 252 is the number of business days in a year. The risk-free return, R_f , typically represents the yield of the 10-Year Treasury Note; however, due to its prolonged low yield (Yahoo Finance, 2023) during the analysed period, a 0% rate is commonly used and was adopted in our analysis.

4 Experimental Results

The performances of the five portfolios which were constructed as described in Section 3 are illustrated in Figure 3. Notice that the deep learning approaches outperformed the lexicon-based approaches in terms of cumulative returns, particularly those relying on general-purpose dictionaries (HIV-4 and VADER). This was to be expected, given that lexicon-based approaches often fail to capture the contextual meaning of sentences, whilst the nuanced nature of financial text significantly reduces the accuracy of general-purpose dictionaries.

Moreover, observe from the top-right panel of Figure 3 and Table 2 that the difference in cumulative returns between our model and the best performing method among the considered ones increased over time. The significant advantage of our FinLlama from 2019 onwards can be explained by a significant rise in the daily average number of companies traded, as a result of an increasingly

Date	Daily Companies Traded	Return Difference	Best existing method
1/1/2016	14.7	-8.1	LMD
1/1/2017	19.0	40.1	FinBERT
1/1/2018	20.0	59.3	FinBERT
1/1/2019	20.0	54.7	FinBERT
1/1/2020	28.0	73.2	FinBERT
1/1/2021	49.2	98.5	FinBERT

Table 2: Difference in cumulative returns between our FinLlama model and the best-performing existing method (among LMD, HIV-4, VADER, and FinBERT) on the first day of each year, along with the daily average number of companies traded during the previous year. A negative difference in returns indicates that the cumulative returns of our model are lower than those of the best existing method at that date.

more diverse set of articles in our news corpus over the years. Indeed, this difference in returns exhibits a positive correlation of 0.81 with the daily average number of companies invested, with a P-value of 0.048, indicating the statistical significance of the trend (significant if P-value < 0.05). The summary of the difference in cumulative returns between our model and the best performing existing method on the first day of each year, along with the daily average number of companies traded during the previous year, is shown in Table 2.

It is important to note that the increase in the daily average number of companies traded coincides with a rise in the number of articles used to calculate the daily sentiment of each company from 2018 onwards. This behaviour is attributed to Reuters first starting to produce digital content in 2018, followed by a dramatic increase from 2020 onwards, when MarketWatch began producing AI-generated articles on stock price updates, as shown in Figure 4. Additionally, there has been a natural increase in the amount of digital articles produced by all three sources since 2019.

The increased returns resulting from more informed trading decisions, along with the growing gap between the returns of our model and those of the best existing method, highlight the superior ability of our model to achieve accurate financial sentiment valence and strength quantification, compared to existing methods. This is because, the accuracy of sentiment parameters becomes increasingly important with the rise in the number of companies traded and the volume of articles used to make trading decisions. Such trend has been observed over time due to the expanding corpus of financial news articles used during the trading stage.

The improved sentiment classification accuracy

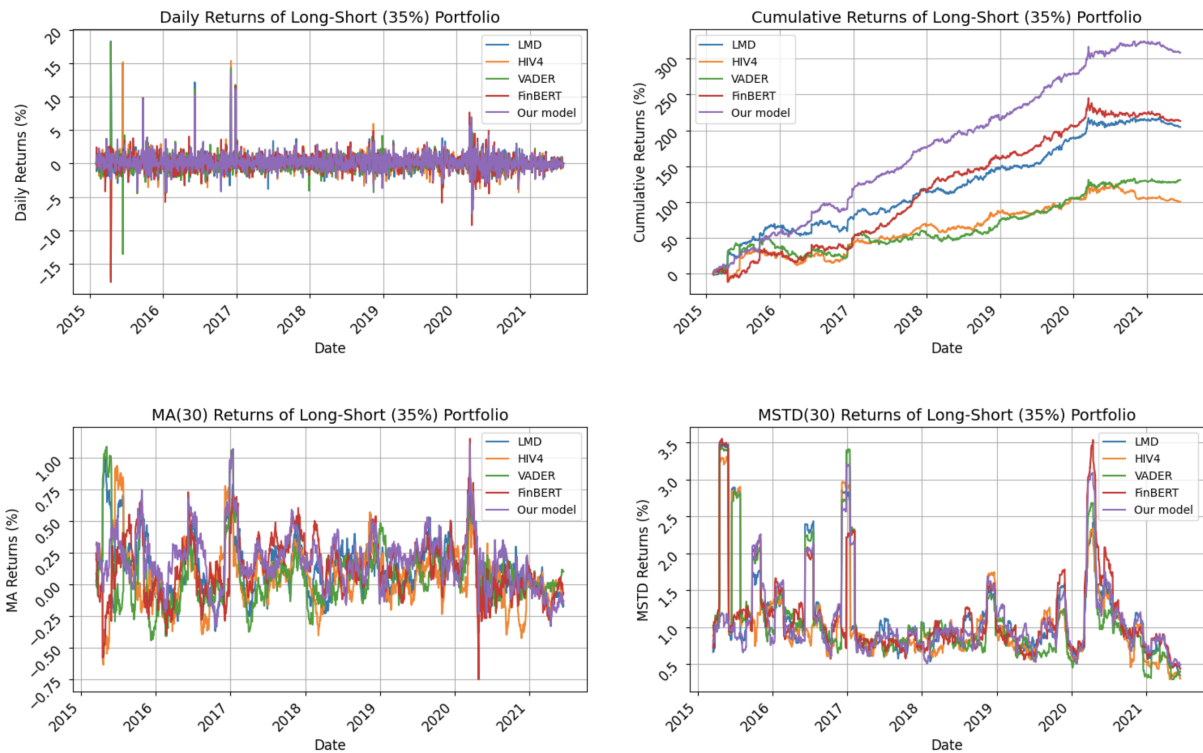


Figure 3: Comparison of the performance of the 35% long-short portfolios which were constructed using the five considered sentiment analysis methods, for the time period of February 2015 to June 2021. The MA(30) and MSTD(30) represent, respectively, the moving average and the moving standard deviation of the returns calculated over a 30-day rolling window.

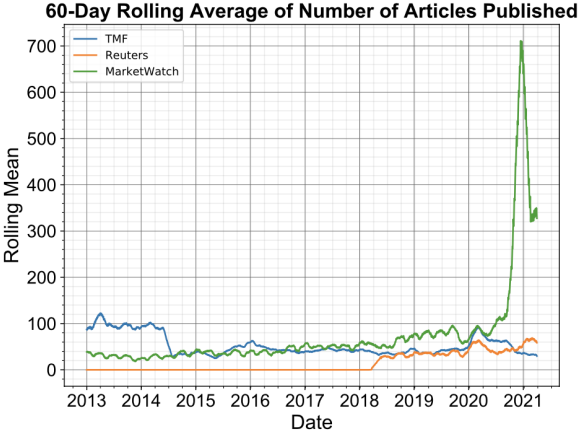


Figure 4: The 60-day rolling average of total number of articles published on each of The Motley Fool, Reuters and MarketWatch from 01/01/2013 to 31/05/2021

exhibited by our model also leads to more robust trading decisions, as indicated in the bottom two panels of Figure 3. In particular, a comparison of our FinLlama model with FinBERT, the current best performing model in the literature, shows that during turbulent economic periods caused by unexpected events or economic changes, the standard deviation of our model was lower than that of Fin-

BERT, while achieving similar or higher returns. The enhanced robustness of FinLlama is evident across a range of socio-economic and geo-political events that caused significant movements in the S&P 500, identified through the business information database Factiva, most notably:

- New trading regulations in China, renewed worries about the Greek economy running out of money, and tepid US corporate earnings in April 2015.
- Concerns about the Federal Reserve increasing interest rates, uncertainty about Greece defaulting on their debt, and geopolitical events and tensions, including the Saint-Quentin-Fallavier attack in June 2015.
- Apprehension about the economic impact of the 2016 US elections, including potential changes in trade policies, tax reforms, regulatory adjustments, and shifts in domestic and international economic relations in January 2017.
- Significant fears about the economic effects of the COVID-19 pandemic, including concerns

499
500
501
502
503
504
505
506
507
508
509
510
511
512
513
514
515
516
517
518
519
520
521

491
492
493
494
495
496
497
498

	LMD	HIV-4	VADER	FinBERT	FinLlama (Ours)	S&P 500
Cumulative Returns (%)	204.6	100.4	130.6	213.0	308.2	83.1
Annualized Return (%)	29.1	13.5	17.9	30.3	45.0	11.3
Sharpe Ratio	1.5	0.7	0.9	1.5	2.4	0.62
Annualized Volatility (%)	19.5	18.9	19.6	20.3	18.6	18.5

Table 3: Statistical comparison between the performances of the five considered sentiment analysis methods using a 35% long-short portfolio. For Cumulative Returns, Annualized Return and Sharpe Ratio, higher is better. For Annualized Volatility, lower is better.

about a severe economic downturn, increased unemployment rates, corporate bankruptcies, and a dramatic decline in consumer spending and business investments in March 2020.

The quantitative results, displayed in Table 3, support the qualitative observations mentioned above and suggest that the 35% long-short portfolio, constructed using our fine-tuned Llama-2 model, was the most successful.

Overall, our FinLlama model successfully generated significantly higher returns for investors compared to all other considered methods, and most importantly FinBERT, whilst simultaneously reducing portfolio risk and being more robust to turbulent economic periods, as indicated by the higher Sharpe ratio and lower annualized volatility.

5 Conclusion and Future Work

We have introduced an innovative approach to financial sentiment analysis which rests upon the fine-tuning of a general-purpose LLM. The proposed method has capitalised on the extensive knowledge base and generative nature of LLMs, combining their inherent text generation with the classification ability. In addition, such an approach has enabled the LLMs to become more attuned to the nuanced language of the finance sector, whilst minimising their resource utilisation and computational demands.

Our fine-tuned Llama2 7B model, termed FinLlama, has been used to construct a long-short portfolio, yielding results that have surpassed those of the existing methods in the field. The FinLlama has achieved cumulative returns which have outperformed the currently leading FinBERT model by 44.7%, while achieving a significantly higher Sharpe ratio and lower annualized volatility. This demonstrates that fine-tuning an LLM can yield superior results, even with a small amount of task-specific data. In addition, the present work has set a new benchmark in the field, transcending traditional measures such as the accuracy and F1-score,

which are commonly used in the literature. It is our hope that such an approach is a step towards narrowing down the divide between academic research and practical applications within quantitative finance.

Our future research will aim to enhance both the sentiment classification accuracy and efficiency of fine-tuned LLM models by incorporating additional techniques to produce a tractable and interpretable platform to facilitate the application of artificial intelligence (AI) in the finance sector.

Disclaimer: Nothing herein is financial advice, and NOT a recommendation to trade real money. Please use common sense and always first consult a professional before trading or investing.

6 Limitations

While the proposed FinLlama has successfully achieved its objectives of improving sentiment classification accuracy, it occasionally misclassifies articles, resulting in losses on a small minority of trading days. These misclassifications exemplify the limitations in handling certain nuances of financial language and context. Future work will involve the analysis of the causes of such misclassifications, followed by rigorous performance bounds and risk analysis. In addition, the current fine-tuning process would benefit from incorporating Reinforcement Learning from Human Feedback (RLHF) (Ouyang et al., 2022), which could further enhance the accuracy and robustness of FinLlama in understanding complex financial language.

In terms of portfolio construction, our study does not integrate additional technical indicators and trading costs, in combination with sentiment strength, which could enhance our portfolio strategy. Moreover, our current work has been limited to equities within the S&P 500. In future work, we aim to investigate the performance of FinLlama in trading other financial instruments, such as bonds and derivatives, as well as its effectiveness in dif-

604	ferent markets.		
605	References		
606	Dogu Araci. 2019. FinBERT: Financial sentiment analysis with pre-trained language models. <i>arXiv preprint arXiv:1908.10063</i> .		
607			
608			
609	J.B. Berk and P. M. DeMarzo. 2019. Corporate finance. volume 5.		
610			
611	Ziwei Chen, Sandro Gössi, Wonseong Kim, Bernhard Bermeitinger, and Siegfried Handschuh. 2023. FinBERT-FOMC: Fine-tuned FinBERT Model with sentiment focus method for enhancing sentiment analysis of FOMC minutes. pages 357–364. Proceedings of the 4th ACM International Conference on AI in Finance.		
612			
613			
614			
615			
616			
617			
618	Nello Cristianini and John Shawe-Taylor. 2000. <i>An introduction to support vector machines and other kernel-based learning methods</i> . Cambridge University Press.		
619			
620			
621			
622	Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. BERT: Pre-training of deep bidirectional transformers for language understanding. In <i>North American Chapter of the Association for Computational Linguistics</i> .		
623			
624			
625			
626			
627	Eugene F. Fama. 1970. Efficient capital markets: A review of theory and empirical work. <i>The Journal of Finance</i> , 25(2):383–417.		
628			
629			
630	Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. 2021. LoRA: Low-rank adaptation of large language models. <i>arXiv preprint arXiv:2106.09685</i> .		
631			
632			
633			
634			
635	C.J. Hutto and Eric Gilbert. 2015. VADER: A parsimonious rule-based model for sentiment analysis of social media text. volume 08, pages 216–225. Proceedings of the 8th International Conference on Weblogs and Social Media, ICWSM 2014.		
636			
637			
638			
639			
640	Zheng Tracy Ke, Bryan T Kelly, and Dacheng Xiu. 2019a. Predicting returns with text data. Technical report, National Bureau of Economic Research.		
641			
642			
643	Zheng Tracy Ke, Bryan T. Kelly, and Dacheng Xiu. 2019b. Predicting returns with text data. NBER Working Papers 26186, National Bureau of Economic Research, Inc.		
644			
645			
646			
647	Xiaodong Li, Haoran Xie, Li Chen, Jianping Wang, and Xiaotie Deng. 2014. News impact on stock price return via sentiment analysis. <i>Knowledge-Based Systems</i> , 69:14–23.		
648			
649			
650			
651	Xiao-Yang Liu, Guoxuan Wang, and Daochen Zha. 2023. FinGPT: Democratizing internet-scale data for financial large language models. <i>arXiv preprint arXiv:2307.10485</i> .		
652			
653			
654			
	Ilya Loshchilov and Frank Hutter. 2017. Fixing weight decay regularization in Adam. <i>ArXiv</i> , abs/1711.05101.	655	656
			657
	Tim Loughran and Bill McDonald. 2011. When is a liability not a liability? Textual analysis, dictionaries, and 10-Ks. <i>The Journal of Finance</i> , 66:35 – 65.	658	659
			660
	Neural Magic. 2022. Twitter financial news sentiment.		661
	Macedo Maia, Siegfried Handschuh, André Freitas, Brian Davis, Ross McDermott, Manel Zarrouk, and Alexandra Balahur. 2018. WWW’18 Open Challenge: Financial Opinion Mining and Question Answering. <i>Companion Proceedings of the The Web Conference 2018</i> .	662	663
		664	665
		666	667
	Pekka Malo, Ankur Sinha, Pekka Korhonen, Jyrki Wallenius, and Pyry Takala. 2014. Good debt or bad debt: Detecting semantic orientations in economic texts. <i>Journal of the Association for Information Science and Technology</i> , 65(4):782–796.	668	669
		670	671
		672	
	Kostadin Mishev, Ana Gjorgjevikj, Irena Vodenska, Lubomir T. Chitkushev, and Dimitar Trajanov. 2020. Evaluation of sentiment analysis in finance: From lexicons to transformers. <i>IEEE Access</i> , 8:131662–131682.	673	674
		675	676
		677	
	Long Ouyang, Jeff Wu, Xu Jiang, Diogo Almeida, Carroll L. Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, John Schulman, Jacob Hilton, Fraser Kelton, Luke E. Miller, Maddie Simens, Amanda Askell, Peter Welinder, Paul Francis Christiano, Jan Leike, and Ryan J. Lowe. 2022. Training language models to follow instructions with human feedback. <i>ArXiv</i> , abs/2203.02155.	678	679
		680	681
		682	683
		684	685
		686	
	P. J. Stone, D. C. Dunphy, M. S. Smith, and D. M. Ogilvie. 1966. <i>The General Inquirer: A Computer Approach to Content Analysis</i> . MIT Press.	687	688
		689	
	Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. 2023. Llama 2: Open foundation and fine-tuned chat models. <i>arXiv preprint arXiv:2307.09288</i> .	690	691
		692	693
		694	695
	Oliver Wang. 2023. News with GPT instructions.		696
	Shijie Wu, Ozan Irsoy, Steven Lu, Vadim Dabravolski, Mark Dredze, Sebastian Gehrmann, Prabhajan Kam-badur, David Rosenberg, and Gideon Mann. 2023. BloombergGPT: A large language model for finance. <i>ArXiv</i> , abs/2303.17564.	697	698
		699	700
		701	
	Yahoo Finance. 2023. Treasury yield 10 years historical data.	702	703
	Hongyang Yang, Xiao-Yang Liu, and Christina Dan Wang. 2023. FinGPT: Open-source financial large language models. <i>arXiv preprint arXiv:2306.06031</i> .	704	705
		706	

Zichao Yang, Diyi Yang, Chris Dyer, Xiaodong He, Alex Smola, and Eduard Hovy. 2016. [Hierarchical attention networks for document classification](#). pages 1480–1489.

Boyuan Zhang, Hongyang Yang, and Xiao-Yang Liu. 2023. [Instruct-FinGPT: Financial sentiment analysis by instruction tuning of general-purpose large language models](#). *ArXiv*, abs/2306.12659.

A Existing Sentiment Analysis Methods

A.1 Lexicon-Based Approaches

A.1.1 Harvard IV-4 Psychological Dictionary (HIV-4)

The HIV-4 is one of the oldest manually constructed lexicons, and is used for objectively identifying specified characteristics of messages in areas involving social science, political science, and psychology. The latest version of the HIV-4 dictionary contains over 11,000 words which are classified into one or more of 183 categories. In this work, we focus on the 1,045 words labelled as positive and the 1,160 words labelled as negative.

A.1.2 Loughran and McDonald (LMD) Dictionary

Loughran and McDonald evaluated standard dictionaries and found that these frequently misclassify terms within financial texts. This insight led to the development of the LMD dictionary, which is specifically tailored for the financial sector. The dictionary categorizes words into six distinct sentiment categories: negative, positive, uncertainty, litigious, strong modal, and weak modal. It was constructed using data from 50,115 10-K filings from 8,341 firms listed on the New York Stock Exchange (NYSE) and the National Association of Securities Dealers Automated Quotations (NASDAQ), covering the period from 1994 to 2008. Overall, the LMD dictionary contains 2,355 negative financial words and 353 positive financial words.

A.1.3 Valence Aware Dictionary for Sentiment Reasoning (VADER)

The VADER dictionary combines lexical features, derived from micro-blog contexts, with the grammatical and syntactical conventions that humans typically employ to express or emphasize sentiment intensity. This enables VADER to accurately quantify the sentiment strength of text. The model

contains approximately 9,000 token features, which are each assigned a sentiment score ranging from -4 (indicating extremely negative sentiment) to +4 (indicating extremely positive sentiment). The overall polarity score for a text is calculated by summing the sentiment scores of each word present in the lexicon, with the final score normalized to fall within the range of -1 to +1.

A.2 Deep Learning Approaches

A.2.1 FinBERT

FinBERT leverages the BERT model architecture, and is specifically tailored for financial contexts. It was pre-trained on a substantial financial text corpus consisting of 1.8M news articles sourced from the Thomson Reuters Text Research Collection (TRC2) dataset, spanning the years between 2008 to 2010. Further refinement was achieved through fine-tuning on the Financial Phrasebank (FPB) dataset, thus enhancing its capabilities in financial sentiment classification. FinBERT generates SoftMax outputs for three labels: positive, negative, and neutral.