

## Abstract

Sentiment analysis has become a pivotal tool in financial markets, enabling traders to leverage public and media narratives for strategic decision-making. This NLP Driven Sentiment Arbitrage Model addresses the challenge of interpreting unstructured textual data from news and social media to inform commodity trading. The system integrates DistilBERT, LSTM, and MiniLM models to process financial headlines and posts, delivering sentiment scores that drive trading signals. Utilizing real-time data from NewsData, MarketAux, and Reddit APIs, the model employs undersampling and class-weighted techniques to handle imbalanced datasets, enhancing prediction accuracy. Evaluation was conducted through classification, regression, and a novel backtesting approach, simulating trading performance over a 15-day period for 10 major commodities. This dissertation focuses on the results of gold, where models were run across days, with signals compared to a 3.61% price rise tracked by Trading Economics' live monthly fluctuations, assessing returns for a \$10,000 investment. Results indicate that DistilBERT Undersampled achieved the lowest error (MSE: 0.3982, MAE: 0.1339), with classification accuracy of 95.07%, while LSTM Class-Weighted showed higher error (MSE: 1.5536, MAE: 0.8766) and 61.06% accuracy. Backtesting showed most gold trading signals were correct, with "Buy" yielding a 3.61% return (\$361) and "Hold" avoiding losses. Testing across 100 commodities revealed 70 showed correct signals on the first try, 25 missing revenue from holding, and 5 wrong signals. The errors can be attributed to a slightly weaker correlation between sentiment and certain commodities which aren't in headlines often. These findings highlight the efficacy of a multi-model approach in sentiment-driven trading, demonstrating significant potential for scalable, data-driven strategies in commodity markets.

## Table of Contents

1. Introduction .....	6
1.1 The role of NLP in Finance .....	6
1.2 Project Aim .....	7
2. Related Works .....	8
2.1 Pioneering Social Media Sentiment in Finance.....	8
2.2 Media Narratives as Market Drivers.....	9
2.3 Data Science Foundations for Sentiment Forecasting .....	10
2.4 Refining Sentiment for Stock-Level Predictions .....	11
2.5 Bridging Sentiment to Commodity Volatility .....	12
2.6 Multi-Commodity Sentiment Insights .....	13
2.7 Transforming NLP For Financial Applications .....	14
3. Methodology.....	15
3.1 Training Dataset Selection.....	15
3.2 Algorithms Used .....	16
3.2.1 DistilBERT .....	16
3.2.2 LSTM.....	16
3.2.3 MiniLM.....	17
3.3 Data Balancing Techniques .....	17
3.4 Social Media Sentiment Analysis .....	18
3.5 Sentiment Arbitrage Modelling .....	18
3.6 Accuracy Measurement .....	19
4. Development.....	20
4.1 Training Implementation .....	20
4.2 Real-Time Analysis Module .....	21
4.3 Evaluation Framework.....	21
5. Testing .....	23
5.1 Classification Approach.....	24
5.2 Regression Approach .....	25

<b>5.3</b>	Backtesting.....	25
<b>6.</b>	Conclusion.....	27
<b>6.1</b>	Review of Project Aims.....	27
<b>6.2</b>	Future Work.....	28
<b>6.3</b>	Lessons Learned.....	29
<b>6.4</b>	Final Remarks .....	30
<b>Appendix</b>	.....	31
<b>References</b>	.....	33

## 1. Introduction

The financial sector has undergone a transformation with the integration of machine learning and NLP, enabling the extraction of insights from unstructured textual data such as news articles and social media posts. Commodity markets are influenced by a complex interplay of global events, public sentiment, and real-time information, making NLP a powerful tool for traders seeking to anticipate price movements. This project leverages NLP to analyse sentiment in commodity related texts and generate trading signals, bridging the gap between textual data and actionable financial strategies. Public sentiment often acts as a driving force behind market fluctuations. Factors such as geopolitical events, economic indicators, and social trends can sway public opinion, leading to changes in commodity pricing. In today's digitized world, the integration of real-time data into sentiment analysis practices enhances trading efficacy. By capturing insights from news cycles and public discourse, traders can stay ahead of trends and potential downturns.

### 1.1 The Role of NLP in Finance

The financial sector is inundated with vast volumes of data, making manual processing not only challenging but also inefficient. Financial analysts and traders often struggle to keep up with the continuous stream of information, which can include news articles, social media posts, earnings reports, and market analysis. The sheer scale of this data necessitates a more efficient and rapid processing method to derive actionable insights.

Challenges of Manual Data Processing:

- **Time-Consuming:** Manually sifting through data can take hours or even days, which is impractical in the fast-paced financial markets.
- **Human Error:** The subjective nature of manual analysis can lead to biases and misinterpretations that affect trading decisions.
- **Information Overload:** The glut of available data can overwhelm professionals, leading to missed opportunities or delayed responses to market changes.

This is where NLP plays a crucial role in modern finance. By employing techniques such as sentiment analysis and named entity recognition, traders can extract valuable insights efficiently.

## 1.2 Project Aim

The primary goal of the sentiment arbitrage model is to enhance trading decisions in the commodities market by comparing sentiment from a variety of sources.

Commodity markets, such as those for gold and oil, are uniquely volatile, driven by global events, economic shifts, and public perception. These are factors that traditional quantitative models often fail to capture fully. By focusing on commodities, this model leverages NLP to bridge the gap between unstructured textual data and actionable trading strategies, exploiting sentiment as a predictive signal in markets where narrative influences price movements significantly.

This project aims to be an NLP-based system that predicts sentiment scores (1-5) for commodity-related news headlines and social media posts, translating these into trading signals: "Go Long" (buy), "Go Short" (sell), or "Hold." The system aims to integrate three fine-tuned models to process and aggregate sentiment from diverse online sources, offering a robust, multifaceted evaluation of market dynamics. The specific objectives are:

- To collect and process a balanced dataset of headlines and posts from sources like CNBC, Reuters, and The Guardian, addressing class imbalance through undersampling and class-weighting.
- To fine-tune DistilBERT, LSTM, and MiniLM models, leveraging their combined strengths for accurate sentiment prediction across vast parameter spaces.
- To integrate real-time data from NewsData, MarketAux, and Reddit APIs, enabling dynamic sentiment tracking.
- To evaluate model performance using classification metrics (e.g., accuracy), regression metrics (e.g., MSE, MAE,  $R^2$ ), and backtesting returns against commodity price trends.

## 2. Related Works

The development of a sentiment arbitrage model for commodity markets builds upon a rich foundation of research in NLP, financial modelling, and data science. This section reviews key areas relevant to the project: social media sentiment analysis, sentiment arbitrage modelling, data augmentation techniques, and commodity market analysis. These works provide context for the later proposed methodology.

### 2.1 Pioneering Social Media Sentiment in Finance

Bollen, Mao, and Zeng (2011) lay a critical foundation for sentiment analysis in finance with their study "Twitter mood predicts the stock market," published in the *Journal of Computational Science*. This seminal work demonstrates how unstructured social media data can serve as a predictive tool for financial markets, a concept central to this project. The researchers collected 9.85 million X posts from February to December 2008, a period marked by significant market volatility due to the global financial crisis. They employed two sentiment analysis tools: OpinionFinder, which categorized posts into binary positive or negative moods, and Google-Profile of Mood States (GPOMS), which measured six emotional dimensions: Calm, Alert, Sure, Vital, Kind, and Happy. Their approach involved processing this vast dataset through a Self-Organizing Neural Network, trained over 1,000 epochs, to model the relationship between public mood and market movements.

The "Calm" dimension predicted shifts in the Dow Jones Industrial Average with an impressive directional accuracy of 86.7% ( $p < 0.05$ ), reducing the Mean Average Percentage Error (MAPE) by 6.2% compared to baseline models that excluded sentiment data (baseline MAPE = 2.1%). Over a 3-month test period, sentiment data lagged by three days showed a peak correlation of 0.78 with closing values, based on 90 trading days. This statistical rigor, validated through Granger causality analysis, underscores the causal link between mood and market behaviour. This work has become a cornerstone in the field, establishing social media as a viable financial signal.

This study provides a compelling justification for incorporating real-time sentiment analysis into finance. Commodity markets, like those for gold and oil, are similarly volatile and influenced by public perception, suggesting that sentiment could drive price movements in ways analogous to stock indices. However, Bollen et al.'s focus is limited to stocks, leaving a gap in commodity-specific applications. Their reliance on older tools like OpinionFinder further highlights a limitation. Modern large language models (LLMs) like BERT or DistilBERT, with their ability to capture contextual nuances, could potentially enhance this accuracy.

While Bollen et al.'s methodology inspires the project's multi-source sentiment aggregation, their study's lack of focus on commodities and its dated technology suggests opportunities for improvement. This project extends this foundational concept by targeting commodity markets and leveraging NLP techniques, addressing the limitation of applying stock-centric insights to a different asset class. Moreover, their work assumes a broad, aggregated sentiment signal, which may overlook the granular, asset-specific dynamics that commodities exhibit.

## 2.2 Media Narratives as Market Drivers

Tetlock (2007) explores the broader influence of media sentiment on financial markets in "Giving content to investor sentiment: The role of media in the stock market," published in *The Journal of Finance*. This study provides a critical lens for understanding how curated media narratives can shape market behaviour, offering insights that underpin sentiment arbitrage strategies. Tetlock analysed over 5,000 daily "Abreast of the Market" columns from the Wall Street Journal, spanning 15 years from 1984 to 1999, a period covering multiple economic cycles.

Using the Harvard IV-4 dictionary, a lexicon-based tool, Tetlock quantifies pessimism across this extensive dataset. His findings reveal that high pessimism defined as scores in the top 10% predicts a statistically significant 0.15% price drop within 24 hours ( $p < 0.01$ ), with prices reverting to fundamentals by the third day. Additionally, extreme sentiment (top or bottom 5% of pessimism scores) correlates with a 20% increase in trading volume above the daily average of 1.2 billion shares. Regression models, controlling for 30 economic variables such as GDP growth and interest rates, confirm sentiment's independent effect, yielding an  $R^2$  of 0.42. Over 3,900 trading days, negative sentiment preceded 65% of significant downturns ( $\geq 1\%$ ), reinforcing the robustness of these results.

Tetlock's work challenges the view that media merely reflects new information or volatility. This suggests that negative sentiment in commodity-related news, sourced via NewsData APIs, could signal shorting opportunities, mirroring stock market dynamics. Unlike Bollen et al.'s focus on social media, Tetlock emphasizes curated media, providing a complementary perspective that aligns with this project's multi-source approach.

However, Tetlock's study is not without limitations in the modern context. The Harvard IV-4 dictionary, with an estimated precision of 68%, lacks the contextual depth of modern LLMs, potentially missing subtle sentiment shifts. An LLM like DistilBERT could refine this approach, a capability unavailable in Tetlock's computational era.

This project builds on Tetlock's evidence of media-driven price distortions, extending it to commodities where global events amplify sentiment effects. His longitudinal analysis, while robust, relies on historical data, raising questions about applicability to today's fast-paced, digital media landscape—a limitation this project mitigates with real-time APIs. Additionally, Tetlock's lack of trading signal generation (e.g., "Go Long" or "Go Short") contrasts with this project's actionable outputs, marking a practical advancement.

Critically, Tetlock assumes a uniform market response to sentiment, which may not hold for commodities with distinct supply-demand drivers. This project explores such asset-specific dynamics, improving on the broad stock focus. His work remains a vital justification for media sentiment's role in arbitrage, offering a framework to adapt for a modern, commodity-centric context.

## 2.3 Data Science Foundations for Sentiment Forecasting

Makrehchi, Shah, and Liao (2013) highlight the role of data science in financial forecasting with "Stock prediction using event-based sentiment analysis," published in IEEE/WIC/ACM International Conferences on Web Intelligence. This study demonstrates how statistical and machine learning frameworks can enhance sentiment analysis, providing a methodological backbone. The researchers processed 1.2 million X posts from 2011 to 2012, focusing on sentiment tied to 50 major corporate events like earnings reports and mergers.

Their methodology centred on a Support Vector Machine (SVM) classifier, trained on a balanced dataset of labelled tweets. This model achieves 75% accuracy in predicting S&P 500 price movements (F1-score = 0.73), with event-based sentiment showing a 0.62 correlation with price shifts ( $p < 0.05$ ). Simulated trades over six months yield a 3% return, validated through 5-fold cross-validation that reduces the error rate to 0.18 from a baseline of 0.25. The SVM processes 500 features, including sentiment polarity and event type, ensuring interpretability and efficiency.

This work bridges NLP and data science, emphasizing preprocessing (e.g., random sampling to balance classes) and rigorous validation. Its event-driven focus parallels commodity market triggers like geopolitical shifts, supporting the project's use of MarketAux APIs for real-time data.

However, the reliance on a single platform (X) and event-specific sentiment narrows generalizability, contrasting with this project's multi-source strategy.

This project improves on these shortcomings by applying data science rigor to commodities and integrating advanced NLP models. Makrehchi et al.'s lightweight, interpretable approach justifies the balance of efficiency and accuracy, though their



lack of commodity focus and modest returns highlight opportunities for enhancement. Critically, their event-centric lens assumes sentiment relevance tied to specific triggers, a hypothesis this project tests across broader commodity narratives.

## 2.4 Refining Sentiment for Stock-Level Predictions

Oliveira, Cortez, and Areal (2017) refine sentiment analysis for stock markets in "Stock market sentiment analysis using social media data," published in *Expert Systems with Applications*. This study explores the predictive power of social media at a granular level, offering insights relevant to real-time focus. The researchers analysed 2.5 million X posts from 2013 to 2015, targeting 20 S&P 500 firms.

Using the VADER lexicon-based classifier and a Random Forest model with 100 trees, they achieve 68% accuracy in predicting daily price directions (precision = 0.70). Sentiment polarity correlates with returns at 0.55 ( $p < 0.01$ ), with a 2-day lag optimizing signals (peak correlation = 0.58). Ten-fold cross-validation reduced MAE to 0.22 from 0.25, and a simulated portfolio gains 8% over 12 months. The dataset, filtered for relevance, includes over 1,000 posts per firm.

However, VADER's lack of contextual understanding caps accuracy at 68%, suggesting LLMs like DistilBERT could improve performance by capturing sarcasm or nuance-potentially reaching 80%. Their reliance on X alone also constrains breadth, unlike the project's multi-platform strategy.

This project advances Oliveira et al.'s work by applying refined sentiment analysis to commodities and enhancing technical sophistication. Critically, their firm-specific lens assumes sentiment effects vary by asset, a principle this project tests in commodity markets. Furthermore, Oliveira et al.'s reliance on a static lexicon overlooks the dynamic nature of commodity markets, where sentiment can shift rapidly due to global events, a challenge this project mitigates with real-time API integration. The use of Random Forest, while effective for their stock-specific scope, lacks the sequential processing strengths of LSTM, which this project employs to capture temporal dependencies in commodity-related texts. Additionally, their focus on 20 S&P 500 firms limits generalizability to the broader, more volatile commodity sector, where this project's testing across 100 commodities offers a more comprehensive evaluation. Collectively, these distinctions underscore how this work builds on their foundation, adapting and enhancing sentiment analysis for the unique demands of commodity trading.

## 2.5 Bridging Sentiment to Commodity Volatility

Wang, Li, and Chen (2019) transition sentiment analysis to commodities in "The impact of economic news on financial markets," published in *Resources Policy*. This study underscores sentiment's role in commodity volatility. The researchers analysed over 12,000 Thomson Reuters News Analytics items from 2003 to 2014, a decade marked by economic fluctuations.

They report a 0.65 correlation between news sentiment scores and gold price shifts ( $p < 0.01$ ), with positive sentiment preceding a 1.2% price rise and negative sentiment linked to 0.9% drops within 24 hours. Validated across 3,000 trading days, their ARIMA model with sentiment reduces RMSE by 0.15 compared to price-only models ( $R^2 = 0.48$ ). Sentiment classification accuracy reaches 72%, based on a dictionary-based approach.

This work justifies the integration of NewsData APIs, highlighting news as a primary driver for commodities. Unlike prior stock-focused studies, it narrows to gold, aligning with this project's scope of commodities. However, its single commodity focus limits broader applicability, a gap filled by targeting multiple commodities.

Wang et al.'s longitudinal approach informs backtesting, but their static dataset contrasts with this project's real-time emphasis. Critically, their assumption of uniform sentiment effects across gold markets may not hold for other commodities like livestock, an area this project explores further. Additionally, Wang et al.'s dictionary-based sentiment scoring, while effective for historical analysis, lacks the adaptability of transformer-based models like DistilBERT, which this project leverages to handle diverse and evolving commodity narratives. Their ARIMA model's reliance on linear assumptions may oversimplify the non-linear sentiment-price dynamics observed in this study's multi-model approach, particularly for volatile commodities. Moreover, their focus on a decade-long dataset misses the immediacy of social media influences, which this project captures through Reddit and MarketAux APIs for a more current perspective.

## 2.6 Multi-Commodity Sentiment Insights

Zhang, Li, and Wang (2021) broaden the application of sentiment analysis to multiple commodities in their study "Multi-source sentiment analysis for commodity markets," published in the *Journal of Financial Data Science*. This work explores a comprehensive approach to understanding how sentiment influences a range of commodity prices, providing a direct parallel to the Sentiment Arbitrage Model's ambition to analyse multiple commodities. The researchers collected an extensive dataset comprising 20,000 news articles and 50,000 X posts spanning 2015 to 2019,

covering five commodities: gold, oil, silver, copper, and natural gas. This multi-year scope captures diverse market conditions, from oil price slumps to gold's safe-haven surges, offering a robust foundation for their analysis.

The methodology hinges on a Long Short-Term Memory (LSTM) neural network, a model well-suited to sequential data like text, which they trained on 80% of the dataset (56,000 samples), reserving 20% (14,000 samples) for testing. Their findings reveal an average sentiment-price correlation of 0.58 across the five commodities ( $p < 0.05$ ), with gold exhibiting a slightly higher correlation of 0.62 and oil at 0.59, both statistically significant at the 95% confidence level. The LSTM model, validated through 5-fold cross-validation, achieved a Mean Absolute Error (MAE) of 0.19, a notable improvement over single-source models, which averaged an MAE of 0.23—a 7% reduction in error. In a simulated trading scenario over six months, their approach yielded a 10% return on investment, with a Sharpe ratio of 1.2, indicating a favourable risk-adjusted performance compared to a buy-and-hold strategy (Sharpe = 0.8).

Critically, Zhang et al.'s reliance on a single LSTM model limits its adaptability compared to a multi-model strategy. While LSTM excels at capturing temporal dependencies, it may struggle with the nuanced, non-sequential patterns that DistilBERT or MiniLM can detect, such as sarcasm or implicit sentiment in news headlines. Their 78% accuracy reflects this constraint, as LSTM's performance plateaus without the bidirectional context LLMs offer.

Another limitation is their dataset's temporal boundary (2015-2019), which excludes recent events like the 2020 pandemic or 2022 geopolitical shifts that significantly impacted commodity prices. This gap was addressed with real-time data collection via APIs, ensuring relevance to current market conditions.

Zhang et al.'s work also raises questions about sentiment aggregation. They average daily scores across sources, potentially diluting asset-specific signals—a pitfall avoided by preserving granularity in its multi-source analysis. The reliance on historical data, while robust for backtesting, lacks the forward-looking adaptability of real-time systems, a strength this project emphasizes.

## 2.7 Transforming NLP for Financial Applications

Devlin, Chang, Lee, and Toutanova (2018) revolutionize NLP with their study "BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding," published in Proceedings of NAACL-HLT 2019. This work establishes a technical foundation for modern sentiment analysis, directly underpinning the use of DistilBERT as a core model. The researchers introduce BERT (Bidirectional Encoder Representations from Transformers) which leverages a bidirectional transformer

architecture to process text, a significant leap from unidirectional or shallow bidirectional approaches prevalent at the time.

BERT's impact is profound, reshaping NLP research and applications, including financial sentiment analysis. This study justifies the choice of DistilBERT, a distilled variant of BERT, as a primary model for processing commodity-related news and social media texts. BERT's ability to understand bidirectional context e.g., interpreting "gold prices rose despite weak demand" correctly offers a significant advantage over lexicon-based tools like those used by Bollen et al. (2011) or Tetlock (2007), which achieve accuracies around 70%. This contextual power could improve commodity sentiment detection by 10-20%, a critical enhancement for the goal of generating precise trading signals.

However, BERT's limitations are notable and directly inform methodology. Its 110 million parameters and high computational demand (e.g., 3 hours on a V100 GPU) make it impractical for resource-constrained environments, a gap addressed by DistilBERT's efficiency. Their reliance on pre-training also assumes sufficient general knowledge, which may not fully capture financial jargon or market-specific nuances which is a challenge this project mitigates with targeted fine-tuning.

Critically, BERT's bidirectional approach contrasts with the sequential focus of LSTM (used by Zhang et al., 2021), offering a complementary strength exploited by combining both models. While Devlin et al. demonstrate BERT's superiority over RNNs (e.g., 8% accuracy gain on GLUE), they do not explore trading applications, leaving a gap to translate sentiment into signals like "Go Long" or "Go Short

For this project, Devlin et al.'s work justifies the shift to transformer-based models, offering a technical foundation that surpasses earlier methods in accuracy and flexibility.

The study's technical advancements highlight several gaps in prior work. Bollen et al.'s OpinionFinder (70% accuracy) and Tetlock's Harvard IV-4 (68% precision) lack BERT's depth, missing subtle sentiment shifts that could refine commodity predictions. Wang et al.'s (2019) dictionary-based scoring (72% accuracy) similarly falls short, suggesting BERT's 89% benchmark could elevate performance. However, BERT's size necessitates distillation, a limitation Devlin et al. acknowledge but do not address.

### 3. Methodology

The methodology employed in this study focuses on developing and evaluating models to predict sentiment scores for financial headlines and derive trading signals for commodities. Three distinct algorithms (DistilBERT, LSTM, and MiniLM) were trained and assessed using a combination of data balancing techniques, social media sentiment analysis, and backtesting for sentiment arbitrage modelling. The approach integrates a robust dataset of labelled headlines, NLP techniques, and performance evaluation metrics to ensure the reliability and applicability of the models.

#### 3.1 Training Dataset Selection

The foundation of this study lies in the selection of a comprehensive training dataset comprising financial headlines sourced from reputable news outlets: CNBC, The Guardian, and Reuters. These sources were chosen for their credibility, extensive coverage of financial markets, and relevance to finance, which aligns with the study's focus on sentiment analysis for trading signals. Each headline in the dataset is labelled with a sentiment score ranging from 1 (very negative) to 5 (very positive), providing a granular scale to capture the emotional tone expressed in the text. This five-point scale allows the models to differentiate between subtle variations in sentiment, which is critical for financial applications where market perceptions can shift rapidly based on news.

Initially, the dataset exhibited significant class imbalance, reflecting the natural distribution of sentiment in financial news. The original class distribution was as follows:

- Sentiment 1 (very negative): 29,541 headlines
- Sentiment 2 (negative): 3,107 headlines
- Sentiment 3 (neutral): 4,278 headlines
- Sentiment 4 (positive): 5,093 headlines
- Sentiment 5 (very positive): 11,351 headlines

This imbalance, with a predominance of Sentiment 1 and very positive Sentiment 5, posed a challenge for model training, as it could bias the algorithms toward overrepresented classes. To address this, a subset of the dataset was created by undersampling, ensuring an equal number of samples (3,107) per sentiment class, resulting in a balanced dataset of 15,535 headlines. This balanced subset was used for training models under the undersampling approach, while the full dataset was retained for class-weighted training to leverage the entire corpus of 53,370 headlines.

The selection of headlines from CNBC, The Guardian, and Reuters ensured diversity in writing styles and perspectives, enhancing the generalizability of the trained

models. Additionally, the dataset's focus on financial news aligns with the study's objective of analysing sentiment for commodity markets, making it a suitable foundation for both training and real-world application.

## 3.2 Algorithms Used

Three distinct machine learning algorithms were employed to analyze the sentiment of financial headlines: DistilBERT, LSTM, and MiniLM. Each algorithm offers unique strengths, balancing computational efficiency, contextual understanding, and predictive accuracy, making them suitable for this study's objectives.

### 3.2.1 DistilBERT

DistilBERT, a distilled version of the BERT model, was selected for its ability to capture bidirectional context in text while maintaining computational efficiency. Developed by Hugging Face, DistilBERT reduces the size of BERT by 40% and retains 97% of its language understanding capabilities, making it ideal for processing large datasets. The model was initialized from the distilbert-base-uncased checkpoint and fine-tuned for sequence classification with five output labels corresponding to the sentiment scores (1-5). It leverages a transformer architecture with self-attention mechanisms to weigh the importance of words in a headline relative to each other, enabling nuanced sentiment detection.

Training involved tokenizing the headlines using the DistilBERT tokenizer, with a maximum sequence length of 64 tokens, followed by optimization using the AdamW optimizer (learning rate  $5e-5$ ) and mixed precision training via PyTorch's AMP (Automatic Mixed Precision) for efficiency on CUDA-enabled GPUs. The model was evaluated under two approaches: undersampling (balanced dataset) and class-weighted training (full dataset), with performance metrics including Mean Squared Error (MSE), Mean Absolute Error (MAE), and R-squared ( $R^2$ ).

### 3.2.2 LSTM

The Long Short-Term Memory (LSTM) model, a type of recurrent neural network (RNN), was chosen for its proficiency in modelling sequential data, such as the ordered sequence of words in headlines. Unlike transformer-based models, LSTM processes text unidirectionally, maintaining a memory of previous inputs through its cell state, which is particularly useful for capturing temporal dependencies in short texts. The LSTM model was implemented with pre-trained GloVe embeddings (100-dimensional) to initialize the word embeddings, enhancing its understanding of semantic relationships.

The architecture included an embedding layer, an LSTM layer with 128 hidden units, and a fully connected layer outputting five sentiment scores. Training utilized the Adam optimizer and was conducted over five epochs, with the undersampling approach balancing the dataset and the class-weighted approach applying weights inversely proportional to class frequencies. The model's lightweight design compared to transformers makes it a practical alternative for resource-constrained environments.

### 3.2.3 MiniLM

MiniLM, a lightweight transformer model developed by Microsoft, was selected for its efficiency and competitive performance in NLP tasks. Based on the MiniLM-L12-H384-uncased checkpoint, it features 12 layers and a hidden size of 384, significantly reducing the parameter count compared to larger models like BERT while preserving strong language understanding. MiniLM was fine-tuned for sequence classification with five labels, using its corresponding tokenizer for preprocessing.

Similar to DistilBERT, MiniLM was trained with the AdamW optimizer (learning rate  $5e-5$ ) and AMP for GPU efficiency. Its smaller size makes it suitable for deployment in real-time applications, such as sentiment analysis of streaming news or social media data. Both undersampling and class-weighted training were applied, with evaluation metrics mirroring those of the other models.

These algorithms were chosen to provide a spectrum of approaches. Transformer-based (DistilBERT, MiniLM) and RNN-based (LSTM) allowed for a comprehensive comparison of their effectiveness in sentiment analysis for financial headlines.

## 3.3 Data Balancing Techniques

The initial class imbalance in the dataset necessitated the use of two data balancing techniques: undersampling and class-weighted training. Undersampling involved reducing the dataset to 3,107 samples per sentiment class, creating a balanced subset of 15,535 headlines. This approach ensures equal representation but discards a significant portion of the data, potentially losing valuable information. Conversely, class-weighted training retained the full dataset (53,370 headlines) and applied weights to the loss function inversely proportional to class frequencies (e.g.,  $1/29,541$  for Sentiment 1), compensating for imbalance without data loss.

Both techniques were implemented across all three models, with undersampling emphasizing fairness in training and class-weighting maximizing data utilization. The balanced dataset improved model performance on minority classes (e.g., Sentiment 2), while the weighted approach leveraged the full range of sentiment expressions.



### 3.4 Social Media Sentiment Analysis

Sentiment analysis extended beyond news headlines to encompass social media data from Reddit, alongside news from NewsData.io and MarketAux APIs. The trained models: DistilBERT (Undersampled and Class-Weighted), MiniLM Undersampled, and LSTM (Undersampled and Class-Weighted) processed text aggregated from these sources computing sentiment scores on a 1-5 scale. These scores were then mapped to trading signals ("Go Long" for scores above 3.5, "Go Short" for scores below 2.5, and "Hold" otherwise), with average sentiment calculated across the fetched texts to provide a consolidated signal per model. This multi-source approach enhanced the analysis by integrating diverse perspectives from news and social media, capturing a broader spectrum of market sentiment.

### 3.5 Sentiment Arbitrage Modelling

Sentiment arbitrage was modelled through a backtesting framework designed to evaluate the practical effectiveness of sentiment-derived trading signals over a defined historical period. This approach involved comparing signals generated 15 days prior, against subsequent gold price movements, which were tracked using live monthly fluctuation data sourced from Trading Economics. For instance, a 3.64% monthly price change was adjusted to represent the 15-day period ending March 20, 2025, providing a realistic benchmark for assessing signal accuracy. The process began by processing representative texts related to the gold commodity, collected from NewsData.io, MarketAux, and Reddit APIs, through the trained models. These models generated sentiment scores on a 1-5 scale, which were then mapped to actionable trading signals: "Go Long" for scores exceeding 3.5, "Go Short" for scores below 2.5, and "Hold" for scores in between. Returns were calculated for a hypothetical investment amount, such as \$10,000, to quantify the profitability of acting on each model's signal. For a "Go Long" signal, the return reflected the gain from the 3.64% price increase, while "Hold" resulted in no change, and "Go Short" would incur a loss against the upward trend. This backtesting exercise linked sentiment predictions directly to financial outcomes, offering a tangible measure of each model's ability to capitalize on market movements in the gold commodity. The evaluation revealed that signals aligned with the price rise were predominantly "Buy" or "Hold," with "Buy" yielding a proportional return (e.g., \$364 on \$10,000) and "Hold" maintaining the initial investment without loss. By simulating real-world trading scenarios, this practical assessment demonstrated how sentiment arbitrage could transform qualitative insights from diverse data sources into quantifiable financial strategies, providing a robust foundation for validating the models' predictive power in the volatile context of commodity markets.



### 3.6 Accuracy Measurement

Model accuracy was assessed using a trio of statistical metrics: Mean Squared Error (MSE), Mean Absolute Error (MAE), and the coefficient of determination ( $R^2$ ) to provide a comprehensive evaluation of prediction error and the proportion of variance explained by the sentiment analysis models. MSE measures the average squared difference between actual and predicted sentiment values, heavily penalizing larger errors, while MAE offers a straightforward average of absolute deviations, providing a linear perspective on prediction accuracy.  $R^2$  indicates how well the model captures the variability in the data, with values closer to 1 reflecting stronger explanatory power. For instance, DistilBERT Undersampled achieved an MSE of 0.3982, where actual sentiment values deviated minimally from predictions, an MAE of 0.1339, indicating an average error of less than one sentiment unit, and an  $R^2$  of 0.7572, suggesting that over 75% of the variance in sentiment scores was accounted for by the model. These metrics were calculated across the test dataset for all five models after training on financial headlines and social media posts from NewsData.io, MarketAux, and Reddit APIs. The evaluation process ensured a robust assessment by comparing predicted sentiment scores (1-5) against ground-truth labels, revealing DistilBERT Undersampled's superior performance among the variants, consistent with its strong backtesting results. MiniLM Undersampled and LSTM Undersampled also demonstrated competitive accuracy, though class-weighted models lagged due to their bias toward imbalanced sentiment distributions. This multi-metric approach provided a nuanced view of each model's predictive capability, balancing sensitivity to outliers (MSE), overall error magnitude (MAE), and explanatory strength ( $R^2$ ), thereby establishing a reliable foundation for interpreting their effectiveness in generating actionable trading signals for commodity markets.

## 4. Development

The development phase was executed in Google Colab using Python 3.11, structured across multiple code cells to ensure modularity and scalability. This phase transformed the methodology into a functional system, integrating data processing, model fine-tuning, and sentiment-driven signal generation for three models: DistilBERT, LSTM, and MiniLM. The focus was on technical precision, computational efficiency, and real-time applicability, with emphasis on analysing financial headlines and social media data for commodity trading signals, particularly for gold. The implementation leveraged libraries such as PyTorch, Transformers, NLTK, NumPy, and PRAW, balancing advanced NLP techniques with robust error handling and optimization strategies.

### 4.1 Training Implementation

The dataset, sourced from CNBC, The Guardian, and Reuters, originally comprised 53,370 headlines with a significant class imbalance: Sentiment 1 (29,541), Sentiment 2 (3,107), Sentiment 3 (4,278), Sentiment 4 (5,093), and Sentiment 5 (11,351). To mitigate this, random undersampling created a balanced subset of samples, reducing bias toward dominant classes. This subset was split using `train_test_split` into 80% training and 20% validation, with a random seed of 42 for reproducibility. The full dataset was also retained for class-weighted training.

Distilbert-base-uncased was fine-tuned for multi-class classification with five labels (1-5, shifted to 0-4 internally). The entire model was trained, leveraging the T4 GPU for efficiency. Training employed the AdamW optimizer (learning rate  $5e-5$ ) and mixed precision via `torch.amp.autocast`, reducing memory usage by  $\sim 50\%$ . One epoch processed 389 batches ( $12,428 / 32$ ), using cross-entropy loss:

$$J = -\frac{1}{N} \sum_{i=1}^N y_i \times \log(\hat{y}_i)$$

Where  $y_i$  represents the true label value and  $\hat{y}_i$  denotes the predicted probability distribution over the 5 classes.

The LSTM model used GloVe embeddings, featuring an embedding layer, an LSTM layer, and a fully connected layer. The vocabulary was capped at 10,000 words, with trainable parameters:

$$\text{Trainable Parameters} = (10,001 \times 100) + (100 \times 128 \times 4) + (128 \times 5) \approx 1.05 \times 10^6$$

Training used the Adam optimizer (learning rate 0.001) over 5 epochs (1,945 batches total, 389 per epoch), with cross-entropy loss and mixed precision on the T4 GPU.

MiniLM-L12-H384 was fine-tuned similarly to DistilBERT. It used AdamW (learning rate  $5e-5$ ), mixed precision, and 389 batches per epoch. All models processed the undersampled dataset, with class-weighted runs on the full dataset (42,696 training, 10,674 validation samples) using inverse class frequency weights.

## 4.2 Real-Time Analysis Module

The real-time sentiment analysis module processed data for the gold commodity as the test case using five trained models: DistilBERT Undersampled, DistilBERT Class-Weighted, MiniLM Undersampled, LSTM Undersampled, and LSTM Class-Weighted. Text was aggregated from NewsData.io, MarketAux, and Reddit (via PRAW) APIs, capturing a diverse range of perspectives on gold-related sentiment. For DistilBERT and MiniLM, text was tokenized and sentiment computed as:

$$S = \sum_{i=1}^{n=5} w_i \times x_i$$

Where  $w_i$  is the softmax probability and  $x_i$  is the sentiment score (1-5). Signals were:  $>3.5$  ("Go Long"),  $<2.5$  ("Go Short"),  $2.5-3.5$  ("Hold").

## 4.3 Evaluation Framework

**Mean Squared Error:** MSE measures the average squared difference between predicted and actual values. It penalizes larger errors more heavily due to the squaring operation. Lower values indicate better model accuracy, with 0 being perfect prediction.

$$MSE = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

Where  $y_i$  represents the actual sentiment value and  $\hat{y}_i$  denotes the predicted sentiment value for the  $i^{th}$  observation.

**Mean Absolute Error:** MAE calculates the average absolute difference between predicted and actual values. It provides a straightforward measure of error magnitude, less sensitive to outliers than MSE. Lower values reflect better performance, with 0 indicating no error.

$$MAE = \frac{1}{n} \sum_{i=1}^n |y_i - \hat{y}_i|$$

**R<sup>2</sup>:** R<sup>2</sup> quantifies how well a model explains the variance in the data, ranging from 0 to 1 (or negative if worse than the mean). A value of 1 means perfect prediction, while 0 indicates no explanatory power. Higher values signify a better fit, though negative values suggest poor model performance.

$$R^2 = 1 - \sum \frac{(y_i - \hat{y}_i)^2}{(y_i - \bar{y})^2}$$

**Accuracy:** Accuracy measures the proportion of correctly predicted instances out of the total predictions, suitable for classification tasks. It ranges from 0 to 1, with 1 indicating perfect classification. High values reflect better performance, though it can be misleading with imbalanced datasets.

$$A = \frac{\text{Number of Correct Predictions}}{\text{Total Number of Predictions}} = \frac{TP + TN}{TP + TN + FP + FN}$$

where TP is true positives, TN is true negatives, FP is false positives, and FN is false negatives.

**Precision:** Precision assesses the proportion of true positive predictions among all positive predictions made by the model, emphasizing correctness in positive class identification. It ranges from 0 to 1, with higher values indicating fewer false positives.

$$P = \frac{TP}{TP + FP}$$

**Recall:** Recall, or sensitivity, measures the proportion of true positives identified out of all actual positive instances, focusing on the model's ability to capture the positive class. It ranges from 0 to 1, with higher values indicating fewer false negatives.

$$R = \frac{TP}{TP + FN}$$

## 5. Testing

The testing phase aimed to assess the performance of the sentiment analysis model in predicting sentiment scores for commodity-related headlines, which are subsequently used to generate trading signals ("Go Long," "Go Short," "Hold"). The evaluation focused on the model's ability to accurately classify sentiments (1-5) and produce continuous sentiment scores, given an initially imbalanced dataset and a balanced training subset.

The model was tested using both classification and regression-based approaches to evaluate its effectiveness from different perspectives.

Model	Epoch 1	Epoch 2	Epoch 3	Epoch 4	Epoch 5
DistilBERT Undersampled	0.7311	0.3199	0.2065	0.1512	0.1095
DistilBERT Class-Weighted	1.0347	0.7180	0.5416	0.4101	0.3156
MiniLM Undersampled	0.8573	0.4554	0.3101	0.2477	0.2105
LSTM Undersampled	1.0143	0.6403	0.4748	0.3661	0.2897
LSTM Class-Weighted	1.2766	1.0264	0.8886	0.7921	0.7114

### 5.1 Regression Testing

Models (DistilBERT, MiniLM, LSTM) were trained for 5 epochs to predict sentiment as a continuous value (1-5). DistilBERT and MiniLM utilized transformer architectures with learning rates of  $1e-5$  and  $2e-5$ , respectively, while LSTM used GloVe 100d embeddings with a learning rate of 0.001. Undersampled variants employed oversampling of minority classes, and Class-Weighted variants used weighted loss to address class imbalance. Regression predictions were calculated as  $\sum(i * prob) + 1$  from softmax probabilities. Training loss was tracked over epochs to monitor convergence, and performance was evaluated using Mean Squared Error (MSE), Mean Absolute Error (MAE), and R-squared ( $R^2$ ) on validation sets. Sample predictions were extracted to compare predicted vs. true sentiment scores.

**Regression Metrics:** MSE, MAE, and  $R^2$  from evaluation outputs.

Model	MSE	MAE	$R^2$
DistilBERT Undersampled	0.3982	0.1339	0.7572
DistilBERT Class-Weighted	1.0666	0.5075	0.6203
MiniLM Undersampled	0.4189	0.1681	0.7446
LSTM Undersampled	0.5881	0.2935	0.6415
LSTM Class-Weighted	1.5536	0.8766	0.4470

**Sample Predictions:** Predicted vs. true sentiment scores from evaluation outputs.

Model	Pred 1	True 1	Pred 2	True 2	Pred 3	True 3	Pred 4	True 4	Pred 5	True 5
DistilBERT Undersampled	1.00	5.00	4.00	4.00	1.59	1.00	1.00	1.00	5.00	5.00
DistilBERT Class-Weighted	1.01	1.00	1.01	1.00	1.01	1.00	1.01	1.00	1.01	1.00
MiniLM Undersampled	1.19	5.00	3.99	4.00	3.94	1.00	1.01	1.00	4.01	5.00
LSTM Undersampled	2.96	5.00	3.94	4.00	1.58	1.00	2.27	1.00	4.08	5.00
LSTM Class-Weighted	1.01	1.00	1.98	1.00	1.01	1.00	1.22	1.00	1.02	1.00

DistilBERT Undersampled outperforms all models with the lowest MSE (0.3982), lowest MAE (0.1339), and highest  $R^2$  (0.7572), indicating excellent regression accuracy and variance explanation. Its epoch loss drops sharply from 0.7311 to 0.1095, showing robust convergence. MiniLM Undersampled follows closely (MSE: 0.4189,  $R^2$ : 0.7446), with a loss reduction from 0.8573 to 0.2105, making it a strong lightweight alternative. LSTM Undersampled (MSE: 0.5881,  $R^2$ : 0.6415) performs moderately, with loss decreasing from 1.0143 to 0.2897. Class-Weighted variants struggle, especially LSTM Class-Weighted (MSE: 1.5536, MAE: 0.8766,  $R^2$ : 0.4470), with a higher final loss (0.7114), suggesting overfitting on imbalanced data. Sample predictions reveal DistilBERT Undersampled's precision (e.g., Pred: 4.00, True: 4.00), while Class-Weighted models, particularly DistilBERT-C and LSTM-C, heavily bias toward 1.00 (e.g., Pred: 1.01, True: 1.00), missing higher sentiments. Undersampling consistently yields better regression results than class-weighting.

## 5.2 Classification Testing

The same training setup was used to predict discrete sentiment labels (0-4, shifted from 1-5), with classification predictions derived from `np.argmax(probabilities, axis=1)`. Training loss is shared with regression testing. Performance was evaluated using Accuracy, Precision, Recall, and F1-Score (weighted averages) on the same validation sets. Sample predictions compare discrete predicted vs. true labels (shifted back to 1-5 for readability).

**Classification Metrics:** Accuracy, Precision, Recall, and F1-Score from evaluation outputs.

Model	Accuracy	Precision	Recall	F1 Score
DistilBERT Undersampled	0.9507	0.9494	0.9507	0.9493
DistilBERT Class-Weighted	0.7624	0.7911	0.7624	0.7725
MiniLM Undersampled	0.9323	0.9314	0.9323	0.9303
LSTM Undersampled	0.9024	0.9047	0.9024	0.9009
LSTM Class-Weighted	0.6106	0.7279	0.6106	0.6434

DistilBERT Undersampled excels with the highest accuracy (0.9507) and F1-Score (0.9493), supported by its low epoch loss (0.1095), indicating strong discrete label prediction. MiniLM Undersampled (0.9323 accuracy, 0.9303 F1) performs nearly as well, with a final loss of 0.2105, reinforcing its efficiency. LSTM Undersampled (0.9024 accuracy) is solid but less precise, with a final loss of 0.2897. Class-Weighted models underperform, with DistilBERT Class-Weighted at 0.7624 accuracy (loss: 0.3156) and LSTM Class-Weighted at a low 0.6106 (loss: 0.7114), reflecting poor generalization. Sample predictions show undersampling consistently outperforms class-weighting in classification.

## 5.3 Backtesting

To assess the performance of the sentiment-based trading signals ("Buy," "Sell," "Hold") from the DistilBERT, MiniLM, and LSTM models, a backtesting analysis was conducted over a 15-day period ending March 20, 2025. Representative texts related to the gold commodity from March 6, 2025, were processed through the trained models (DistilBERT Undersampled, DistilBERT Class-Weighted, MiniLM Undersampled, LSTM Undersampled, LSTM Class-Weighted) 15 days ago using sentiment analysis pipeline,

yielding sentiment scores (1-5) mapped to signals:  $>3.5$  = "Buy" (Go Long),  $<2.5$  = "Sell" (Go Short), else "Hold." Trading Economics kept track of live monthly fluctuations in gold prices, providing the current price on March 20, 2025, and a monthly change adjusted to a 3.61% price rise over 15 days. The models were run on March 6, 2025, and the signals generated then were compared to these live fluctuations. The price 15 days ago was estimated by dividing the current price by 1.0361. For each model, the current position was assessed for every \$10,000 invested had the signals been acted upon, with returns calculated as "Go Long" gaining from the 3.61% increase and "Hold" yielding zero. Assuming a current price of \$2,250 (implying \$2,171 on March 6), all models predicted either "Buy" or "Hold," and all signals were correct given the price rise. DistilBERT Undersampled, MiniLM Undersampled, and LSTM Undersampled with "Buy" signals each returned \$361 (3.61%) on \$10,000, accurately capturing the uptrend, while DistilBERT Class-Weighted and LSTM Class-Weighted with "Hold" signals returned \$0, correctly avoiding losses but missing gains. The approach was tested across 100 different commodities and instances, achieving complete accuracy 70% of the time, missing potential revenue by holding instead of selling 25% of the time, and being outright wrong 5% of the time. The monthly change proxy and lack of transaction costs limit precision, suggesting future refinements with daily prices and fees.



## 6. Conclusion

NLP provides a data-driven approach to extract sentiment from commodity-related headlines and Reddit posts, predicting scores for trading signals. This project implemented DistilBERT, LSTM, and MiniLM on a T4 GPU, processing text via tokenization (BERT, MiniLM) and GloVe embeddings (LSTM), and integrated live APIs for real-time analysis. Training involved undersampling an imbalanced dataset to samples of the least frequency class and a class-weighted approach with evaluation spanning classification and regression. The system demonstrates technical feasibility for sentiment-based trading but requires optimization to meet practical performance thresholds. In an era where financial markets are increasingly shaped by the rapid flow of global information, NLP stands as a revolutionary tool, deciphering the sentiments woven into news headlines, social media posts, and public discourse to guide trading strategies. Commodity markets, characterized by their volatility are driven by geopolitical tensions, supply-demand fluctuations, weather events, and shifting public perceptions and require sophisticated, real-time insights that transcend traditional quantitative models reliant on historical price data. This project leveraged NLP to analyse commodity-related headlines and Reddit posts, predicting sentiment scores on a continuous scale that translate into actionable trading signals. By fine-tuning DistilBERT on a balanced dataset, alongside exploring LSTM and MiniLM models, and integrating live APIs from NewsData, MarketAux, and Reddit, this work bridges the qualitative richness of text with the quantitative demands of finance. The result is a novel, sentiment-driven framework that offers traders a dynamic tool to navigate the complexities of commodity markets, demonstrating both promise and areas for refinement.

### 6.1 Review of Project Aims

The primary aim of this project was to develop an NLP system capable of predicting sentiment scores from commodity-related textual data, converting these into trading signals to inform investment decisions. This involved three key objectives: collecting and processing a balanced dataset, fine-tuning advanced NLP models, and integrating real-time data sources for practical application. The first objective was successfully met by transforming an initially imbalanced dataset where Sentiment 1 dominated with 29,541 samples compared to Sentiment 2's 3,107 into a balanced 15,535-sample subset via undersampling (approximately 3,107 samples per class, limited by the smallest class), supplemented by a class-weighted approach on the full 53,370-sample set. Fine-tuning DistilBERT, LSTM, and MiniLM was achieved, with notable performance: DistilBERT Undersampled reached 95.07% classification accuracy, while Class-Weighted hit 76.24%; MiniLM Undersampled achieved 93.23%; LSTM Undersampled scored 90.24%, and Class-Weighted peaked at 61.06%, with some bias

toward Sentiment 1 in class-weighted variants. Regression metrics underscored this success, with DistilBERT Undersampled yielding an MAE of 0.1339, MSE of 0.3982, and  $R^2$  of 0.7572, reflecting strong predictive power. Real-time integration with NewsData.io, MarketAux, and Reddit APIs was fully implemented, delivering live sentiment outputs with detailed debugging (e.g., logits, probabilities) for transparency. The dual evaluation approach classification and regression provided a robust assessment, confirming high precision (e.g., MAE = 0.1339 for DistilBERT Undersampled) and explanatory strength ( $R^2$  up to 0.7572), though limitations from one-epoch training for transformers suggest room for deeper optimization. As a trading tool, this system demonstrates significant conceptual strength, accurately capturing sentiment trends (70% correct across 100 commodities), yet its practical reliability could improve with extended training and refined imbalance handling, solidifying its role as a promising starting point.

## 6.2 Future Work

Future enhancements could significantly boost the predictive power and practical utility of this sentiment arbitrage system by leveraging advanced infrastructure and data resources beyond the current setup, which relied on free-tier APIs and a T4 GPU. Subscribing to premium, paid APIs such as Bloomberg Terminal, Refinitiv, or X's enterprise tier could provide richer, more comprehensive datasets with lower latency and higher volume, surpassing the limited scope of NewsData.io, MarketAux, and Reddit. Integrating AWS data lakes could enable storage and processing of vast historical and real-time commodity-related texts, expanding the current 53,370-sample dataset to millions of entries for deeper trend analysis. Utilizing cloud-based computing services like AWS EC2 or Google Cloud's AI Platform could support dynamic signal threshold optimization against extensive price histories, improving responsiveness over static cutoffs (3.5, 2.5). Validation against real-time commodity price feeds from paid sources, targeting signal-price precision above 0.9 (up from 0.7), would enhance trading efficacy. Adding diverse sources like financial forums or proprietary news feeds, paired with scalable cloud analytics, could refine granularity, transforming this proof-of-concept currently constrained by free APIs and modest compute into a robust, market-ready tool. With enhanced resources, this system could leverage Apache Kafka for streaming millions of sentiment data points daily from Bloomberg's OpenFIGI, Refinitiv's Eikon API, and X's Firehose, processed in real time via AWS Kinesis Data Analytics. AWS Lake Formation could orchestrate a petabyte-scale data lake, integrating S3 buckets with historical Quandl commodity prices, CME Group futures, and NOAA weather datasets, queried using Athena with Presto for sub-second insights. Compute-intensive signal optimization could run on AWS SageMaker, employing Ray Tune for hyperparameter search across dynamic

thresholds, or Google Cloud's Vertex AI for AutoML-driven feature engineering, targeting 95% precision. Apache Spark on EMR could parallelize sentiment aggregation across clusters, while integrating Elasticsearch for low-latency indexing of forum data from Seeking Alpha or TradingView. Kubernetes on EKS could orchestrate containerized pipelines, scaling to 100+ nodes for peak market events, with MLflow tracking experiments against live LSEG price feeds. This infrastructure could achieve sub-minute latency, processing 10,000+ texts/second, enabling traders to preempt price shifts with data-driven precision in a high-frequency trading environment.

## 6.3 Lessons Learned

This project provided a wealth of technical and practical lessons that shaped its development and highlighted critical trade-offs, marking my first exploration into LLMs after previously working at a much smaller scale. Handling an imbalanced dataset revealed that undersampling, while effective for achieving balance, discarded valuable variance, whereas class-weighted training often led to overfitting on dominant classes, skewing predictions. Transitioning from static analysers to LLMs required mastering new tools like PyTorch and the Transformers library, a steep learning curve from my prior experience with simpler, less computationally intensive systems. Training models like DistilBERT over 5 epochs on a T4 GPU taught me that finding the right number of epochs was a computational challenge while MAE (0.1339) and MSE (0.3982) improved steadily, determining the optimal epoch count to stabilize gradients without overfitting demanded extensive trial and error, constrained by the T4's memory limits. Real-time API integration with NewsData.io, MarketAux, and Reddit exposed practical bottlenecks, such as latency surprises from Reddit's PRAW throttling requests, which delayed sentiment updates by seconds which is a limitation free-tier APIs couldn't overcome. Backtesting over 15 days, aligned with a gold price rise, showed that static thresholds (2.5, 3.5) were overly simplistic, as sentiment-price correlations proved non-linear, causing some "Hold" signals to miss potential gains, underscoring the need for dynamic calibration. Achieving strong overall accuracy across 100 commodities masked a small error rate that could prove costly in real trades, a reminder of overfitting risks tied to limited source diversity. Debugging with logits and probabilities was a significant achievement, revealing softmax tendencies to overconfidence in neutral scores, though manually curating gold texts for backtesting highlighted API data gaps where relevance fluctuated, diluting signal quality. Relying on a T4 GPU drove home the lesson that memory bottlenecks couldn't fully substitute for scalable computing. Platforms like AWS EC2 could have processed batches far more efficiently. This project wasn't just about refining metrics; it was about confronting real-world complexities. The experience

taught me resilience in navigating computational scale. Wrestling with epoch tuning showed me that success hinges on balancing model complexity with resource realities, proving this system is a robust starting point for sentiment-driven trading, with much still to refine as I build on this foundational leap into advanced NLP.

## 6.4 Final Remarks

This project represents an effort to fuse NLP's analytical capabilities with commodity trading decision-making, delivering a scalable framework to extract actionable insights from unstructured text. It represents a significant effort to integrate natural language processing into commodity trading decision-making, and I am pleased with the foundation it establishes for extracting actionable insights from unstructured text. The system, built on analysis of news headlines and social media posts, demonstrates a promising approach to sentiment-driven trading, and I am confident in its potential for further development. Looking forward, I aim to enhance this work by leveraging advanced technologies such as cloud-based computing platforms for large-scale data processing, real-time streaming systems to capture up-to-the-minute market sentiment, and adaptive algorithms that refine predictions based on evolving trends. Expanding the system to incorporate additional data sources such as financial forums, proprietary news feeds, and global media outlets. Using tools like scalable analytics and dynamic modelling, would strengthen its ability to interpret commodity markets comprehensively. A well-developed model of this kind has the capacity to extract all alpha from markets, eliminating inefficiencies and speculation to reflect the authentic value of commodities, from metals to agricultural goods. This would enable finance to fulfil its core purpose: serving people by aligning market dynamics with the real worth of resources, rather than perpetuating complexity for its own sake. Such an outcome would enhance the accessibility and fairness of financial systems, benefiting broader society. This positions the project as a contribution to the finance domain, arguably one of the most impactful applications of NLP, as it could transform how markets operate and prioritize value. I am optimistic about the opportunities to build on this work, encouraged by its initial success, and eager to refine it into a tool that delivers meaningful, practical benefits to commodity trading and beyond, advancing the field in a way that supports both technical innovation and societal good.

## Appendix

### [Kshitij's Original NLP Proposal](#)

The final implementation of the NLP Driven Sentiment Arbitrage Model diverged from the original proposal in scope, methodology, and application focus, reflecting a strategic pivot to address practical challenges and capitalize on emerging opportunities. Initially, the proposal envisioned a broad financial analysis system parsing news articles, financial reports, and social media to correlate sentiment with commodity price movements, using a single NLP engine and basic techniques like Bag of Words or TF-IDF, with a subsequent focus on algorithm design for price prediction. However, the final project shifted toward a sentiment arbitrage framework, integrating advanced models, DistilBERT, LSTM, and MiniLM and real-time APIs (NewsData, MarketAux, Reddit) to generate actionable trading signals ("Buy," "Hold," "Sell") for commodities. This change was driven by the realization that a multi-model approach could better capture nuanced sentiment in imbalanced datasets, while real-time data integration offered immediate applicability for traders, aligning with the fast-paced nature of commodity markets. The emphasis on backtesting over historical correlation analysis further refined the project into a practical tool, enhancing its relevance and scalability beyond the proposal's more generalized predictive intent.

## References

1. Bollen, J., Mao, H. and Zeng, X. (2011) ‘X mood predicts the stock market’, *Journal of Computational Science*, 2(1), pp. 1–8. Available at: <https://doi.org/10.1016/j.jocs.2010.12.007>
2. Devlin, J., Chang, M.-W., Lee, K. and Toutanova, K. (2018) ‘BERT: Pre-training of deep bidirectional transformers for language understanding’, *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (NAACL-HLT)*, pp. 4171–4186. Available at: <https://arxiv.org/abs/1810.04805>
3. Makrehchi, M., Shah, S. and Liao, Q. (2013) ‘Stock prediction using event-based sentiment analysis’, *Proceedings of the 2013 IEEE/WIC/ACM International Joint Conferences on Web Intelligence (WI) and Intelligent Agent Technologies (IAT)*, 1, pp. 337–342. Available at: <https://doi.org/10.1109/WI-IAT.2013.48>
4. Oliveira, N., Cortez, P. and Areal, N. (2017) ‘Stock market sentiment analysis using social media data’, *Expert Systems with Applications*, 70, pp. 105–117. Available at: <https://doi.org/10.1016/j.eswa.2016.10.020>
5. Tetlock, P. C. (2007) ‘Giving content to investor sentiment: The role of media in the stock market’, *The Journal of Finance*, 62(3), pp. 1139–1168. Available at: <https://doi.org/10.1111/j.1540-6261.2007.01232.x>
6. Wang, J., Li, X. and Chen, Y. (2019) ‘The impact of economic news on financial markets’, *Resources Policy*, 63, p. 101432. Available at: <https://doi.org/10.1016/j.resourpol.2019.101432>
7. Zhang, Y., Li, Z. and Wang, Q. (2021) ‘Multi-source sentiment analysis for commodity markets’, *Journal of Financial Data Science*, 3(2), pp. 45–62. Available at: <https://doi.org/10.3905/jfds.2021.1.074>