

SentiTrade-HMA: Hierarchical Multi-Source Aggregation of LLM-Derived Financial Sentiments for Enhanced Algorithmic Trading

Safae AHRARA, Chaimae ELAMRAOUI, Abdssamad IDBOUSSADEL and Mohammed NAJID

Abstract. Accurately forecasting financial markets remains a formidable challenge due to non-stationarity, heterogeneity, and the rapid incorporation of qualitative information into prices. While traditional econometric and deep learning models capture numerical patterns, they often neglect the semantic richness of financial news and reports. This study introduces *SentiTrade-HMA*, a hybrid architecture that fuses context-aware sentiment signals derived from a QLoRA fine-tuned Llama-2-7B model with a Temporal Fusion Transformer (TFT) for multi-horizon forecasting. By extracting nuanced, interpretable sentiment features and integrating them with high-frequency market data, the system overcomes limitations of conventional approaches and lexicon-based sentiment tools. Walk-Forward Validation on S&P 500 constituents (2018–2024) shows a statistically significant reduction in Mean Absolute Error (MAE) of 5% relative to price-only models. Realistic trading simulations demonstrate a Sharpe Ratio of 5.17, a Win Rate of 58.3%, and a Maximum Drawdown of -3.40%, underscoring both predictive accuracy and risk-aware performance. This work provides a reproducible, interpretable framework for combining LLM-driven sentiment with temporal modeling in financial forecasting.

1 Introduction

Financial time-series forecasting has traditionally relied on econometric models such as ARIMA and GARCH, which assume linearity and stationarity. However, real-world markets frequently violate these assumptions due to regime shifts, volatility clustering, and abrupt structural changes. Data-driven methods using deep learning, including LSTMs and Transformers, capture non-linear dependencies, yet often remain opaque and focus primarily on numerical price signals.

Simultaneously, market dynamics are influenced by qualitative information, including financial news, corporate disclosures, and macroeconomic announcements. While the Efficient Market Hypothesis posits rapid price incorporation, behavioral finance demonstrates that interpretation delays, sentiment asymmetries, and cognitive biases can generate exploitable inefficiencies. Capturing these effects requires models capable of reasoning over both price dynamics and the semantic content of textual information.

Existing sentiment analysis approaches, from lexicon-based methods to transformer classifiers like FinBERT [1], provide limited semantic depth or contextual reasoning. Furthermore, hybrid models that naively concatenate textual and numerical features often suffer

from poor interpretability and weak alignment with multi-horizon forecasting objectives.

To address these limitations, we propose *SentiTrade-HMA*, an interpretable architecture that integrates a QLoRA fine-tuned Llama-2-7B model for semantic sentiment extraction with a Temporal Fusion Transformer for multi-horizon prediction. By fusing deep semantic signals with quantitative market features, the system aims to enhance both predictive accuracy and risk-aware trading performance. This work demonstrates that coupling generative LLM-based sentiment modeling with temporal architectures provides a reproducible, interpretable, and high-performance framework for modern quantitative finance.

2 Methodology: System Architecture

SentiTrade-HMA is designed as a dual-stream pipeline that fuses structured quantitative data with unstructured textual sentiment to generate actionable trading signals.

2.1 Data Acquisition and Engineering

The system leverages a seven-year dataset (2018–2024), encompassing diverse market regimes: the 2018–2019 bull market, the COVID-19 crash and recovery, and the volatile 2022–2023 inflationary period.

Quantitative Market Data: We selected 20 diverse S&P 500 constituents across Technology, Finance, Healthcare, and Energy, totaling 35,200 daily OHLCV points from Yahoo Finance. Missing values are handled via forward- and backward-fill, with adjustments for stock splits and dividends. Feature engineering produced 59 technical indicators for the TFT, including Hull Moving Average (HMA) for smooth trend detection, MACD for momentum, Bollinger Bands and ATR for volatility, and RSI for overbought/oversold conditions. Temporal “calendar features” such as day-of-week and proximity to holidays are treated as known future inputs, enabling the model to anticipate cyclical market patterns.

Qualitative Textual Data: Financial news and headlines were collected from NewsAPI and Finnhub (up to 50 articles per ticker per day). Preprocessing includes ticker-specific keyword filtering, deduplication of syndicated news, and normalization (lowercasing, removal of special characters) to ensure high-quality semantic input for the LLM.

2.2 NLP Subsystem: Fine-Tuned Llama 2

The NLP component transforms raw financial news into structured sentiment features used as inputs to the forecasting model.

Fine-Tuning Strategy: To overcome the limited semantic expressiveness of discriminative models such as FinBERT [1], SentiTrade-HMA employs a Llama-2-7B Large Language Model [10]. Fine-tuning is performed using Quantized Low-Rank Adaptation (QLoRA) [11] on the Financial PhraseBank dataset, augmented from 2,258 to 4,984 samples using synonym substitution and back-translation to enhance robustness to linguistic nuance. The model is loaded in 4-bit precision using *bitsandbytes*, with LoRA adapters configured with rank $r = 64$, which controls the expressive capacity of the low-rank adaptation, and scaling factor $\alpha = 16$, which regulates the magnitude of the injected task-specific updates relative to the frozen base model.

Sentiment Inference and Aggregation: Each news headline is classified into {Positive, Neutral, Negative}, with predictions weighted by model confidence rather than hard labels. Daily sentiment vectors are constructed by aggregating all ticker-specific articles, and rolling sentiment features (e.g., moving averages and momentum) are generated to capture temporal decay effects.

2.3 Forecasting Subsystem: Temporal Fusion Transformer

Market forecasting is performed using a Temporal Fusion Transformer (TFT) that ingests fused numerical and sentiment features. Inputs are structured as static covariates (sector, market capitalization), known future inputs (calendar features), and observed past inputs (OHLCV, technical indicators, and LLM-derived sentiment).

The TFT employs masked multi-head self-attention to preserve temporal causality while capturing multi-scale temporal dependencies [2, 5]. Attention weights provide interpretability by revealing which historical periods or news events most strongly influence each prediction [8, 6, 9].

Quantile Forecasting: Rather than producing point estimates, the model outputs multiple conditional quantiles (10th, 50th, 90th percentiles). The inter-quantile range serves as a proxy for predictive uncertainty and downstream risk assessment.

2.4 Trading Strategy Logic

Model forecasts are transformed into trading decisions through a simplified sequential pipeline designed to assess the economic relevance of the predictive signals under realistic execution constraints.

Stage 1: Signal Generation: Directional trading signals are generated by comparing the median predicted return (50th percentile) to a fixed return threshold. A *long* position is opened when the predicted return exceeds $+0.1\%$, while a *short* position is initiated when it falls below -0.1% . Predictions within this band are treated as market noise and result in no trade.

Stage 2: Position Handling: For active signals, trades are executed using a fixed position size. Although volatility indicators such as the Hull Moving Average (HMA) and Bollinger Bands are included as model inputs, they are not explicitly used as exit rules within the trading strategy. This choice isolates the predictive contribution of the forecasting model without introducing additional strategy-level optimization.

Stage 3: Execution Simulation: All trades are evaluated in a realistic backtesting environment incorporating execution frictions.

Transaction costs of 10 basis points and slippage of 5 basis points per trade are applied to account for market impact and execution latency.

3 Experimental Setup

3.1 Evaluation Protocol: Walk-Forward Validation

A critical flaw in many financial AI studies is the use of k-fold cross-validation, which shuffles data and allows the model to train on future market regimes (e.g., training on the 2020 crash to predict 2019 volatility). To strictly avoid this Look-Ahead Bias, SentiTrade-HMA employs Walk-Forward Validation (also known as a Rolling Window approach)[2]. **Window Structure:**

- **Train:** Months 1–6 (e.g., Jan–Jun).
- **Validation:** Month 7 (e.g., Jul) – Used for hyperparameter tuning and early stopping.
- **Test:** Month 8 (e.g., Aug) – Strictly out-of-sample performance measurement.
- **Rolling:** After the test, the window shifts forward by one month. The model is re-trained (or fine-tuned) on the new training window (Feb–Jul) to predict Sep.

This protocol mimics the real-world lifecycle of a trading algorithm, which must adapt to evolving market conditions without knowledge of the future.

3.2 Baselines for Comparison

SentiTrade-HMA is evaluated against three representative baselines to assess both economic relevance and architectural contributions.

- **Buy & Hold (B&H):** A passive investment strategy serving as a market-level performance benchmark.
- **LSTM:** A standard deep learning time-series model using price data and technical indicators only. This baseline evaluates the benefit of the TFT architecture over recurrent models.
- **TFT-NoSent:** An ablated version of the proposed model where sentiment features are removed, isolating the incremental contribution of the NLP subsystem.

4 Results and Analysis

4.1 Predictive Accuracy and Model Comparison

The comparative analysis of predictive power reveals a clear hierarchy in model performance. The results, averaged across the test set period (late 2024), are summarized in Table 1.

Table 1. Comparative Predictive Performance on Out-of-Sample Test Data.

Model	MAE	RMSE	R2 Score
LSTM Standard	0.1243	0.1700	0.9730
SentiTrade-HMA (TFT + LLM)	0.0895	0.1365	0.9743

Analysis: Crucially, the addition of the LLM-derived sentiment signal in SentiTrade-HMA further reduces the MAE to 0.0895. While a 5% improvement might appear incremental in other domains, in the context of highly efficient financial markets, this “edge” is significant. It confirms that the financial news contains an orthogonal signal—information not fully captured by price history alone—that the LLM successfully extracts and the TFT successfully integrates.

Sentiment Model Comparison

To validate the effectiveness of the sentiment extraction component, we compared the fine-tuned Llama-2+QLoRA model against the pre-trained FinBERT baseline. Table 2 summarizes the results.

Table 2. Sentiment Model Performance Comparison.

Model	Accuracy	MCC
FinBERT (Literature Baseline)	97.0%	N/A
Llama-2 + QLoRA (Ours)	99.0%	0.9926

Insight: The fine-tuned Llama-2+QLoRA model demonstrates superior performance (+2% accuracy) compared to FinBERT, highlighting that the extracted sentiment features are more reliable and context-aware. This improvement directly contributes to the enhanced predictive accuracy of the full SentiTrade-HMA pipeline.

4.2 Financial Trading Performance

The truest test of a financial model is its ability to generate profit while managing risk.

Table 3. Trading Performance Metrics (Net of Execution Costs).

Metric	Our solution	Industry Benchmark
Sharpe Ratio	5.17	> 2.0 is Good; > 3.0 is Excellent.
Win Rate	58.3%	Professional discretionary traders aim for 55-60%.
Max Drawdown	-3.40%	< -10% is typically preferred by hedge funds.

Insight: The Sharpe Ratio of 5.17 is the standout metric. Traditional “Buy & Hold” strategies on the S&P 500 typically yield Sharpe Ratios between 0.5 and 1.0 depending on the year. Achieving a value > 5.0 implies that the model’s returns are not only high but remarkably consistent with very low volatility. The Max Drawdown of -3.40% highlights the efficacy of the Quantile Regression mechanism. The low Maximum Drawdown of -3.40% reflects the conservative signal filtering imposed by fixed return thresholds and the probabilistic nature of the TFT forecasts, which naturally avoids trades during periods of high predictive uncertainty.

4.3 Statistical Significance: The Diebold-Mariano Test

To confirm that the outperformance of SentiTrade-HMA is not an artifact of random luck or specific market conditions, we conducted the Diebold-Mariano (DM) Test [7, 12] comparing the forecast errors of SentiTrade-HMA against the LSTM baseline.

- **Hypothesis:** H0: The predictive accuracy of SentiTrade-HMA and LSTM is equal.
- **Result:** The calculated p-value is < 0.01 ($p < 1\%$).
- **Conclusion:** We reject the null hypothesis with $> 99\%$ confidence.

The predictive superiority of SentiTrade-HMA is statistically significant. This provides robust scientific evidence that the fusion of LLM sentiment and TFT architecture constitutes a genuine advancement in forecasting capability.

4.4 Interpretability and “White Box” Analysis

A key objective of SentiTrade-HMA was to move beyond “Black Box” AI. The model provides interpretability through the Temporal Fusion Transformer’s internal mechanisms, revealing which features and historical periods most influence predictions.

The Variable Selection Network (VSN) identifies the most influential features across both global and local contexts. Empirically, the top features over the test period were:

1. **Close Price (Lag 1):** Recent price remains the strongest predictor.
2. **LLM Sentiment Score:** Ranked consistently among the top 3, confirming sentiment as a key driver.
3. **ATR (Volatility):** Helps detect market regimes.
4. **Hull Moving Average:** Validates the choice of HMA over SMA.

These results confirm that LLM-derived sentiment carries predictive weight comparable to core technical indicators.

The TFT attention weights highlight the temporal focus of the model. For example, during an earnings surprise, the model strongly attended to the relevant historical day when the news occurred, demonstrating causal learning rather than mere curve-fitting.

4.5 Related Work: FinGPT and FinGPT Forecaster

Recent open-source financial LLMs, such as **FinGPT** [13] and **Fin-GPT Forecaster** [14], focus on sentiment extraction and stock prediction using LoRA-based fine-tuning and proprietary data pipelines. While these frameworks report strong benchmark performance (up to 82% accuracy), direct comparison with SentiTrade-HMA is not feasible due to differences in datasets, preprocessing, and computational resources.

Our work complements these efforts by demonstrating that *publicly available data* combined with a dual-stream architecture (LLM + TFT) achieves high predictive accuracy (MAE = 0.0895) and robust trading performance, emphasizing reproducibility and interpretability in financial forecasting.

5 Conclusion

This project demonstrates that Large Language Models can contribute measurable value to financial forecasting when integrated within an interpretable and resource-efficient quantitative framework. SentiTrade-HMA combines a QLoRA fine-tuned Llama-2-7B model for context-aware financial sentiment extraction with a Temporal Fusion Transformer for multi-horizon prediction, enabling the fusion of semantic and market-driven signals in a transparent and computationally tractable manner. Empirical evaluation shows a statistically significant improvement in predictive accuracy over price-only baselines (approximately 5% MAE reduction, $p < 0.01$), while realistic trading simulations incorporating transaction costs achieve strong risk-adjusted performance, including a Sharpe Ratio of 5.17, a 58.3% win rate, and a limited maximum drawdown of -3.40%. Interpretability analysis confirms that the LLM-derived sentiment signal consistently ranks among the most influential predictive features, validating its economic relevance rather than serving as a superficial auxiliary input. Overall, the results position SentiTrade-HMA as a cost-aware and scalable proof-of-concept for LLM-centric, interpretable financial AI, with future extensions including dynamic risk-aware execution, richer event-level language modeling, and cross-asset generalization.

References

- [1] Financial Sentiment Analysis Using FinBERT with Application in Predicting Stock Movement. URL: <https://arxiv.org/html/2306.02136v3> (accessed Dec 28, 2025).
- [2] Temporal Fusion Transformers for interpretable multi-horizon time series forecasting | EMIL. URL: <https://ghasemzadeh.com/event/2025-05-ebi-tft/> (accessed Dec 30, 2025).
- [3] FinLlama: Financial Sentiment Classification for Algorithmic Trading Applications - arXiv. URL: <https://arxiv.org/html/2403.12285v1> (accessed Jan 2, 2026).
- [4] Financial sentiment analysis with Large Language Models. URL: <https://lseee.net/index.php/fe/article/download/1838/FE009282.pdf> (accessed Jan 5, 2026).
- [5] Quantum Temporal Fusion Transformer - arXiv. URL: <https://arxiv.org/html/2508.04048v1> (accessed Jan 8, 2026).
- [6] TFT: an Interpretable Transformer - Towards Data Science. URL: <https://towardsdatascience.com/tft-an-interpretable-transformer-70147bcf6212/> (accessed Jan 10, 2026).
- [7] Diebold-Mariano Test | Real Statistics Using Excel. URL: <https://real-statistics.com/time-series-analysis/forecasting-accuracy/diebold-mariano-test/> (accessed Jan 12, 2026).
- [8] Interpretable Time Series Forecasting Using a Temporal Fusion Transformer - MATLAB & Simulink - MathWorks. URL: <https://www.mathworks.com/help/deeplearning/ug/time-series-forecasting-using-temporal-fusion-transformer.html> (accessed Jan 15, 2026).
- [9] Gated Residual Networks (GRNs) - Emergent Mind. URL: <https://www.emergentmind.com/topics/gated-residual-networks-grns> (accessed Jan 17, 2026).
- [10] Fine-Tuning LLaMA 2: A Step-by-Step Guide to Customizing the Large Language Model. URL: <https://www.datacamp.com/tutorial/fine-tuning-llama-2> (accessed Jan 19, 2026).
- [11] Fine-Tuning LLaMA 2 QLORA - Kaggle. URL: <https://www.kaggle.com/code/simranjeetsingh1430/fine-tuning-llama-2-qlora> (accessed Jan 21, 2026).
- [12] Diebold-Mariano Test. URL: <https://maggima.github.io/pages/stats/tests/forecasts/dm.html> (accessed Jan 23, 2026).
- [13] FinGPT: Open-Source Financial Large Language Models arXiv preprint arXiv:2306.06031, 2023. URL: <https://arxiv.org/pdf/2306.06031.pdf> (accessed Jan 24, 2026).
- [14] Assessing the Capabilities and Limitations of FinGPT Model in Financial NLP Applications arXiv preprint arXiv:2507.08015, 2025. URL: <https://arxiv.org/pdf/2507.08015.pdf> (accessed Jan 24, 2026).