```
                    mentioned in `Rationale`.
123
124        - ICL Examples:
125          - "`SPX_Close_Slope`_>_0_&&_`VIX_Close_Slope
               `_<_0_We_have_Market_Confidence"
126          - "`GDP_QoQ`_Falling_&&_`PMI`_<_50_We_have_
               an_Economic_Slowdown."
127
128    5. Options Analysis:
129        - Compare `OTM_Skew`, `ATM_Skew`, and `
             ITM_Skew` IV Skews: Assess differences to
             gauge market sentiment and directional
             bias using their `20Day_Moving_Averages`.
130        - Leverage IV spikes to capitalize on
             speculative directional trades.
131        - Example: "Rising_`ATM_Skew_MA`_>_0,_market_
             pricing_up_move,_with_stable_HV_supports_a
             _LONG_position,_as_it_indicates_growing_
             upside_expectations_without_excessive_fear
             ."
132
133    6. News Analysis:
134        - Use `News_Sentiment` and `News_Impact_Score`
             (1-3).
135        - Only strong directional news (score = 3)
             should override other signals.
136        - Medium news (score = 2) supports but does
             not lead.
137        - Always check if news contradicts macro or
             technical trend.
138
139    7. Performance Reflection and Strategic
         Adaptation:
140        - If `Last_Strategy_Used_Data` is available:
141            - Assess the outcome of the previous
                 strategy by examining `last_returns`
                 and the chosen `last_action`.
142            - Determine if the result aligns with
                 the expectations outlined in the
                 previous `Rationale`.
143            - Identify if the direction (LONG or
                 SHORT) led to desirable or
                 undesirable outcomes.
144            - You must NOT reuse or copy the
                 previous `Rationale`. It is only
                 context for reflection.
145            - Summarize in 1-2 sentences whether the
                 previous strategy performed as
                 expected.
146            - Example: "The_previous_LONG_strategy_
                 yielded_positive_returns,_confirming
                 _the_bullish_setup_based_on_RSI_and_
                 moving_averages."
147            - Do NOT include language or phrasing
                 from the previous rationale.
148        - Confidence assignment:
149            - Assign a Likert score (1 to 3) to your
                 `action_confidence`:
150                - 1: Low confidence; contradictory or
                     weak alignment across features.
151                - 2: Moderate confidence; partial
                     alignment with moderate evidence.
152                - 3: High confidence; strong
                     convergence across key features.
153        - Feature Attribution:
154            - Rank the importance of each major
                 feature used in your current rationale
                 using a Likert scale (1 to 3):
155                - 1: Minimal contribution; not
                     required for the decision.
156                - 2: Moderate contribution; relevant
                     but not critical.
157                - 3: High contribution; pivotal to the
                     trading decision.
158
159  Output:
160    action: Str. LONG or SHORT.
161    action_confidence: int. Likert scale (1-3)
           confidence in the proposed `action`, adjusted
           based on prior strategy outcome if `
           Last_Strategy_Used_Data` is available.
162    explanation: >
```

```
163      A concise rationale (max 350 words) justifying
           the proposed `action`.
164      Include:
165        - The top 5 weighted features used in the
             decision, each labeled with its Likert
             importance (1-3).
166          (e.g., "Stock_Data.Price.Close,_Weight_3,_
             Technical_Analysis.RSI.Value,_Weight_1,_
             Options_Data.ATM_Skew,_Weight_2")
167        - A reflective assessment of `
             Last_Strategy_Used_Data`, including
             whether the past `action` was successful
             and was it maintained given prior `
             Rationale`.
168    features_used:
169      - feature: the features used from the prompt's_
           context.
170        direction: LONG, SHORT, or NEUTRAL
171        weight: A Likert score (1 to 3) described in
             Feature Attribution.
```

# APPENDIX B
## ANALYST PROMPT

The Analyst prompt used in Experiment 1 is presented in Listing 2, adapted from [10]. News corpora were anonymized prior to prompting.

Listing 2. Analyst Prompt

```
1  User_Context:
2    Monthly_News_Articles_List: |
3      "{articles_list}"
4
5  System_Context:
6    Persona: Financial Market Analyst
7    Instructions: |
8      Extract the `Top 3` news factors influencing
           stock price movements from the `
           Monthly_News_Articles_List`. Follow these
           steps:
9
10      1. Rank the news by relevance to stock price
             movements:
11          - Prioritize news related to significant
               financial or market impacts (e.g.,
               acquisitions, partnerships, guidance
               revisions).
12          - Weigh industry trends, macroeconomic
               influences, and analyst ratings based on
               their expected effect on the company
               valuation.
13          - News with broad or long-term implications
               ranks higher.
14
15      2. Summarize content into key factors and
             corporate events affecting stock prices,
             using concise language and causal
             relationships.
16
17      3. For each factor, assign:
18          - `Sentiment`: +1 for positive, -1 for
               negative, 0 for neutral or mixed
19          - `Market_Impact_Score`: Likert scale from 1
               to 3, where:
20            - 1 = minimal relevance
21            - 2 = moderate influence
22            - 3 = high impact driver
23
24      Examples of factors influencing stock prices
             include:
25        - Strategic partnerships or competitor
             activity.
26        - Industry trends or macroeconomic influences.
27        - Product launches or market expansions.
28        - Analyst ratings, significant stock price
             moves, or expectations.
29        - Corporate events: guidance revisions,
             acquisitions, contracts, splits,
             repurchases, dividends.
30
31      Example:
```

---

**Algorithm 1:** Expert Trade Heuristic

---
**Data:** Time-indexed price series
**Result:** Trade action: LONG (1) or SHORT (0)
**1 foreach** *date t in dataset* **do**
**2**     $P_t \leftarrow \text{Close}(t)$;
**3**     $r^{(10)} \leftarrow \frac{P_{t+10}}{P_t} - 1, \ r^{(20)} \leftarrow \frac{P_{t+20}}{P_t} - 1$;
**4**     $r^{\text{weighted}} \leftarrow 0.4 \cdot r^{(10)} + 0.6 \cdot r^{(20)}$;
**5**     **if** $r^{\text{weighted}} >= 0$ **then**
**6**        **Action** $\leftarrow$ LONG     (Trade_Action = 1);
**7**     **else**
**8**        **Action** $\leftarrow$ SHORT     (Trade_Action = 0);

---

```
32            'A_major_tech_company_partners_with_a_leading_
                  automotive_firm_for_EV_battery_innovation.
                  _Analysts_predict_this_could_boost_
                  revenues_significantly.'
33
34        Ranked Factors:
35           1. factor: Strategic partnership in EV
                     battery technology expected to increase
                     revenue.
36                 sentiment: +1
37                 market_impact: 3
38           2. factor: Positive sentiment driven by
                     projected long-term gains.
39                 sentiment: +1
40                 market_impact: 2
41           3. factor: Growing demand for EV technology
                     anticipated to support future earnings.
42                 sentiment: +1
43                 market_impact: 2
44
45     Output:
46        factors:
47           - factor: str. Summary of the news item. Max 70
                     words.
48           - sentiment: int. One of Positive +1, Negative
                     -1, or Neutral 0
49           - market_impact: int. Likert scale 1 to 3
```

## APPENDIX C
## ALGORITHMS

The labeling algorithm emulates expert trading behavior by deliberately leveraging future return information to assign proxy trade actions in hindsight. This approach offers a cost-effective and scalable addition to manual annotation, capturing the general direction an informed trader might take. These synthetic labels are then provided to the LLM, along with a smaller set of HITL annotated examples.

## APPENDIX D
## DATASET

*Market Data*

This market data ($\mathcal{S}_{\text{mk}}$) included OHLCV price series as well as macro-level indicators and forward-looking sentiment signals. Specifically, it comprised:

- Daily returns of the S&P 500 Index (SPX) and NASDAQ-100 Index (NDX). These are market and sector indices,
- Implied Volatility (IV) and Historical Volatility (HV) metrics, derived from the stock's derivatives,
- The CBOE Volatility Index (VIX) as a proxy for market fear and option market expectations,
- *Weekly Past Returns*, which record the percentage change over the past four weekly intervals. The four-week span was selected empirically to align with the model's monthly strategy generation frequency.

These features help in modeling short-term market dynamics.

*Fundamental Data*

Fundamental data ($\mathcal{S}_{\text{fund}}$) has firm-level fundamentals and macroeconomic indicators. Macroeconomic variables provided contextual narrative for interpreting observed signals, and supporting regime identification [8], [9]. This set covered:

- **Liquidity ratios:** Current Ratio, Quick Ratio;
- **Leverage and coverage:** Debt-to-Equity, Interest Coverage;
- **Profitability metrics:** Gross Margin, Operating Margin, Return on Equity (ROE), Return on Assets (ROA);
- **Valuation:** Price-to-Earnings (P/E), Price-to-Book (P/B), Enterprise Value (EV), and Earnings Before Interest, Taxes, Depreciation, and Amortization (EBITDA).
- **Growth:** Revenue and Earnings Growth;
- **Macroeconomic indicators:** Gross Domestic Product (GDP), Purchasing Managers' Index (PMI), Producer Price Index (PPI), Consumer Confidence Index (CCI), U.S. 10-Year Treasury Yield, and the 10Y–2Y yield curve slope.

To enhance temporal abstraction, all variables were computed as quarter-over-quarter (QoQ) or year-over-year (YoY) percentage changes. It is critical to take first-order dynamics as LLMs can recall absolute numbers for economic details, allowing look-ahead bias in the backtests [20].

*Analytics*

Technical indicators ($\mathcal{S}_{\text{an}}$) were computed over rolling 20-day windows using the open-source TA-Lib[7] library. These features include:

- Simple Moving Averages (SMA) over 20, 50, 100, 200 trading-day horizons,
- Relative Strength Index (RSI),
- Average True Range (ATR) for volatility,
- Moving Average Convergence Divergence (MACD) with its signal line and derived strength,
- Volume-Weighted Average Price (VWAP) as a reference anchor for intraday valuations.

Each indicator was extended with slope and z-score to assist the LLM in capturing directional shifts and the statistical significance of deviations. These technical indicators are widely used in trading practice and academic research [18].

*Alternative Data*

Structured representations of financial news headlines ($\mathcal{S}_{\text{alt}}$) were extracted using a large language model (LLM), which anonymized and synthesized the content into latent factors. Following the LLMFactor methodology [10], each news item was distilled into 2–5 interpretable factors, capturing macroeconomic and firm-specific signals.

---

[7]https://ta-lib.org/

| Instrument | Paper SR | SR (±$\sigma$) [$p$-value] | MDD (±$\sigma$) |
|---|---|---|---|
| AB InBev | 0.187 | **1.21 (0.30) [0.00]** | 0.18 (0.08) |
| Alibaba | 0.021 | **0.06 (0.02) [0.00]** | 0.09 (0.01) |
| Amazon | 0.419 | 0.39 (0.45) [0.85] | 0.30 (0.09) |
| Apple | 1.424 | 1.19 (0.55) [0.22] | 0.29 (0.09) |
| Baidu | 0.080 | **0.20 (0.17) [0.00]** | 0.36 (0.09) |
| CCB | 0.202 | **0.33 (0.25) [0.04]** | 0.24 (0.14) |
| Coca Cola | 1.068 | 1.07 (0.53) [0.50] | 0.25 (0.04) |
| Dow Jones | 0.684 | 0.70 (0.30) [0.91] | 0.25 (0.05) |
| ExxonMobil | 0.098 | 0.10 (0.35) [0.91] | 0.34 (0.08) |
| FTSE 100 | 0.103 | **0.50 (0.23) [0.00]** | 0.31 (0.08) |
| Google | 0.227 | **-0.54 (0.59) [0.00]** | 0.43 (0.13) |
| HSBC | 0.011 | **0.38 (0.17) [0.00]** | 0.29 (0.05) |
| JPMorgan Chase | 0.722 | 0.72 (0.31) [0.98] | 0.26 (0.06) |
| Kirin | 0.852 | 0.85 (0.42) [0.99] | 0.39 (0.07) |
| Meta | 0.151 | **0.63 (0.61) [0.01]** | 0.45 (0.27) |
| Microsoft | 0.987 | 0.70 (1.00) [0.38] | 0.28 (0.16) |
| NASDAQ 100 | 0.845 | 0.85 (0.35) [1.00] | 0.16 (0.05) |
| Nikkei 225 | 0.019 | **0.26 (0.29) [0.02]** | 0.29 (0.07) |
| Nokia | -0.094 | **0.07 (0.24) [0.00]** | 0.57 (0.15) |
| PetroChina | 0.156 | 0.22 (0.29) [0.29] | 0.67 (0.00) |
| Philips | 0.675 | **1.40 (0.50) [0.00]** | 0.25 (0.03) |
| S&P 500 | 0.834 | 0.83 (0.25) [1.00] | 0.14 (0.04) |
| Shell | 0.425 | 0.42 (0.37) [0.95] | 0.51 (0.05) |
| Siemens | 0.426 | 0.39 (0.23) [0.43] | 0.26 (0.12) |
| Sony | 0.424 | 0.42 (0.36) [0.97] | 0.16 (0.04) |
| Tesla | 0.621 | 0.48 (0.41) [0.29] | 0.52 (0.09) |
| Tencent | -0.198 | -0.19 (0.33) [0.98] | 0.10 (0.09) |
| Toyota | 0.304 | 0.36 (0.27) [0.37] | 0.45 (0.10) |
| Volkswagen | 0.216 | **0.45 (0.18) [0.00]** | 0.48 (0.09) |

TABLE XI
REPLICATION METRICS FOR [3]

To mitigate memorization and data leakage risks, named entities and dates were anonymized (e.g., "Tesla" becomes "the Company").

APPENDIX E
REPLICATED BENCHMARK METRICS

We report the replicated benchmark metrics in Appendix E for the assets used in [3]. We include the mean SR and MDD, each averaged across 25 runs with standard deviation $\sigma$.

For the SR, we conduct a two-sided one-sample $t$-test to assess whether the metric is significantly different from the published value. The null hypothesis $H_0$ assumes equivalence, i.e. $H_0 : \mu_{SR} = SR_{paper}$.

Since this is a replication test, failing to reject $H_0$ indicates successful replication. $p$-values are computed only for SR; other metrics are reported without significance testing.

All assets have been successfully replicated within acceptable bounds, with exceptions highlighted in bold. Notably, GOOGL, one of the stocks included in our test environment, exhibited a statistically significant deviation from the original benchmark, with a $p$-value below 0.05.