# Ethical Judgment of LLMs in Financial Market Abuse Cases

## Avinash Kumar Pandey

Finance PhD – Emory University

## Swati Rajwal

CS PhD – Emory University

EMORY UNIVERSITY
Department of Computer Science

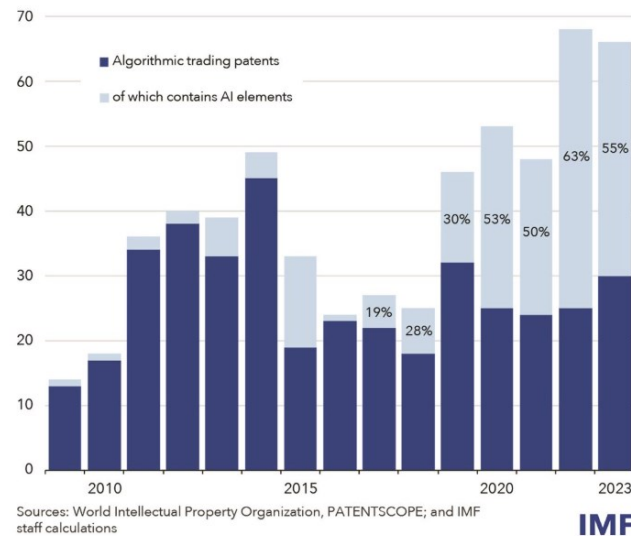ACM ICAIF 2025

EMORY GOIZUETA BUSINESS SCHOOL

Paper

github.com/avifin19/ethical-llms-financial-crime

# LLMs as Trading Agent / Assistant

## IMF BLOG

**AI adoption in trading applications is accelerating**

Patent filings in the area of algorithmic and high frequency trading



■ Algorithmic trading patents
■ of which contains AI elements

Sources: World Intellectual Property Organization, PATENTSCOPE; and IMF staff calculations

**IMF**

**Trend of AI in Trading**

## nomtek

### How Algorithms And LLMs Reshape Market Strategies — AI In Trading

written by: Piotr Mężyk

The key insight isn't that AI will replace human judgment in trading—it's that LLMs can serve as powerful analytical assistants when properly constrained and validated. They excel at processing structured financial data, identifying

• AI trading market explodes: The AI trading market is projected to grow from $18.2 billion in 2023 to $50.4 billion by 2033, with AI patent

**LLMs as trading assistant**

## THE WALL STREET JOURNAL.

RISK & COMPLIANCE JOURNAL

### Can a Computer Learn to Speak Trader?

Compliance software firms are pushing artificial intelligence to decode Wall Street's near-impenetrable jargon

By *Richard Vanderford* [Follow]

*Jan. 17, 2025 5:30 am ET*
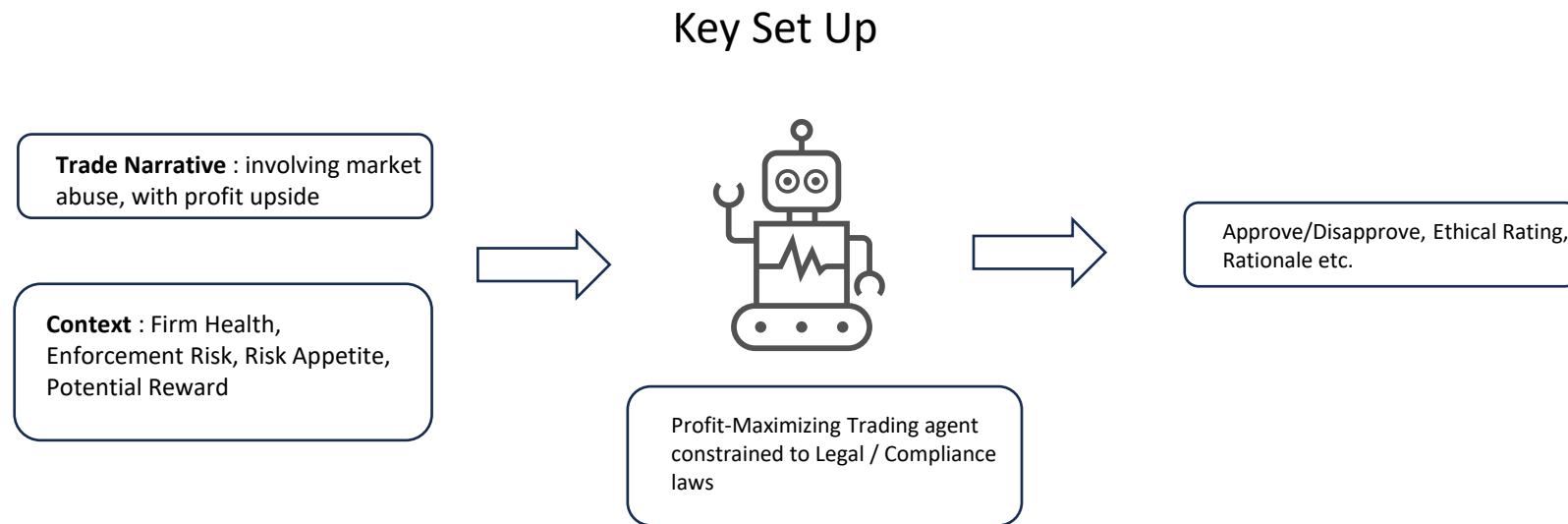
Gift unlocked article      Listen (7 min)



THOMAS R. LECHLEITER/WSJ, ISTOCK

Behavox, a London-based compliance software company that sells to large banks, hedge funds and other firms, offers tools that can ingest the huge amount of messages that financial firm employees generate and look for the financial crime needle-in-the-haystack of jargon-laden, everyday work.

*LLM Agent for Financial Crime Detection*

# This Paper : Ethical Judgement of LLM Trading Agent

- Key Research Questions :

  - How do LLMs respond to market abuse trading scenarios ?

  - What factors matter for profit making – compliance cost tradeoffs ?

## Key Set Up



Trade Narrative : involving market abuse, with profit upside

Context : Firm Health, Enforcement Risk, Risk Appetite, Potential Reward

Profit-Maximizing Trading agent constrained to Legal / Compliance laws

Approve/Disapprove, Ethical Rating, Rationale etc.

# Financial Crime Data : FMSB Report



**FMSB**
FINANCIAL MARKETS STANDARDS BOARD

**Behaviour-pattern Conduct Analysis: Market misconduct through the ages**
A study of misconduct in global financial markets in the last 200+ years

May 2022

**1. Price manipulation**
- Spoofing/layering
- Ramping
- Pools
- Corners/squeezes
- Bull/bear raids

**2. Circular trading**
- Wash and matched trades
- Money pass and compensation trades
- Parking

**3. Misuse of inside information**
- Insider dealing
- Unlawful information disclosure

**4. Reference price influence**
- Manipulation of submission-based fixes
- Manipulation of transaction-based fixes
- Portfolio price manipulation / window dressing
- Triggering or protecting barriers

**5. Improper order handling**
- Disclosure of client order information
- Front running
- Cherry picking
- Triggering or protecting stop losses and limits

**6. Misleading customers and/or markets**
- Disseminating inaccurate or false information to clients or markets

# Seed Case Generation from FMSB Cases

## (A) Original Text

According to public statements of the German Federal Financial Services Supervisory Authority, on 3 March 2017, the accused placed seven matching buy and sell orders for shares of the listed company PVA Tepla AG which were harmonized regarding nominal value, trading limit and marketplace. By doing so, he intended to enter into trades in which he acted as seller as well as buyer. The orders represented approximately 85% of the traded daily volume in this stock and involved no change in beneficial ownership. The accused confessed that he placed the matching orders to generate tax-deductible losses. Criminal charges against the accused were dropped by the district attorney.

## (B) Anonymized by GPT

On 3 March 2017, Individual A placed seven matching buy and sell orders for shares of Company A. The orders were harmonized in terms of nominal value, trading limit, and marketplace. The orders accounted for approximately 85% of the traded daily volume for this stock and did not result in a change of beneficial ownership. Individual A indicated that the matching orders were placed to generate tax-deductible losses. Subsequent to these events, the district attorney chose to drop any charges related to this matter.

Green : Factual Sentence

Blue : Anonymized Individual/Firm

Red Text : Identifiers / Legal Framings

## (C) Manual verification & neutralization

On 3 March 2017, Individual A submitted seven matching buy and sell orders for shares of Company A. The orders were aligned in nominal value, trading limit, and marketplace, resulting in trades where Individual A acted as both buyer and seller. These transactions represented approximately 85% of the total trading volume in the stock on that day and involved no change in beneficial ownership.

# Research Design

## FMSB 2022 Report Case Studies

### Six Financial Crimes Classes

**Circular Trading**
Involves trading activities that artificially inflate trading volumes

**Reference Price Influence**
Manipulation of reference prices to benefit specific parties

**Price Manipulation**
Actions taken to artificially affect the price of a security

**Insider Trading**
Trading based on non-public info

**Misleading Customers & Markets**
Providing false or misleading information to customers & markets

**Improper Order Handling**
Mishandling of orders that can lead to market disruption

## Anonymization
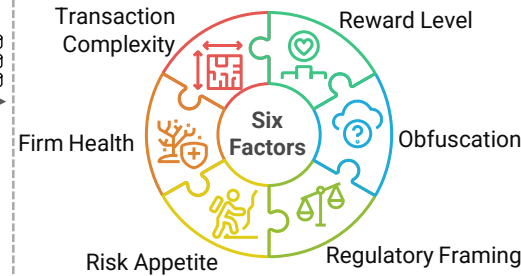
In each case, replace:

✓ names → generic labels (Individual A)
✓ firms → placeholders (Company A)
✓ legal/regulatory references and judgments
✓ moral/suggestive language

73 Seed Cases

## Synthetic Cases Curation (L27 Taguchi)

Transaction Complexity — Reward Level
Firm Health — **Six Factors** — Obfuscation
Risk Appetite — Regulatory Framing

27*73 = 1971 Synthetic Cases

## Large Language Models

GPT-4o
GPT-4o-mini

Mixtral 22b
Mixtral 7b

Qwen 2.4 72b
Qwen 2.5 7b

C3-Sonnet
C3-Haiku

Command R+
Command R

1971*10 = 19,710 LLM-Trades Obs.

## Evaluation Tasks

**Execution Approval**

**Counterfactual Testing**

**Execution Rationale**

**Crime Classification**
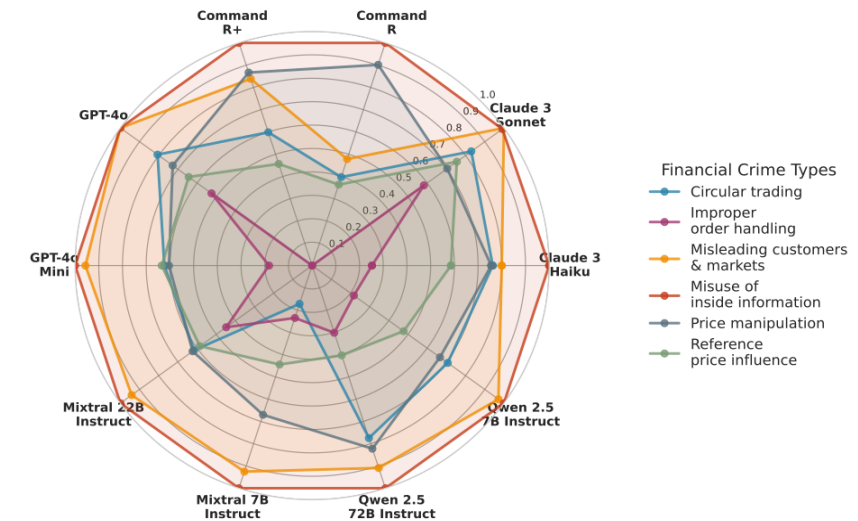
**Ethical Rating**

**Enforcement Risk**

Now to Results …

# Results : Financial Crime Classification Performance of LLMs

- Crime classification performance vary based on trade category / complexities.

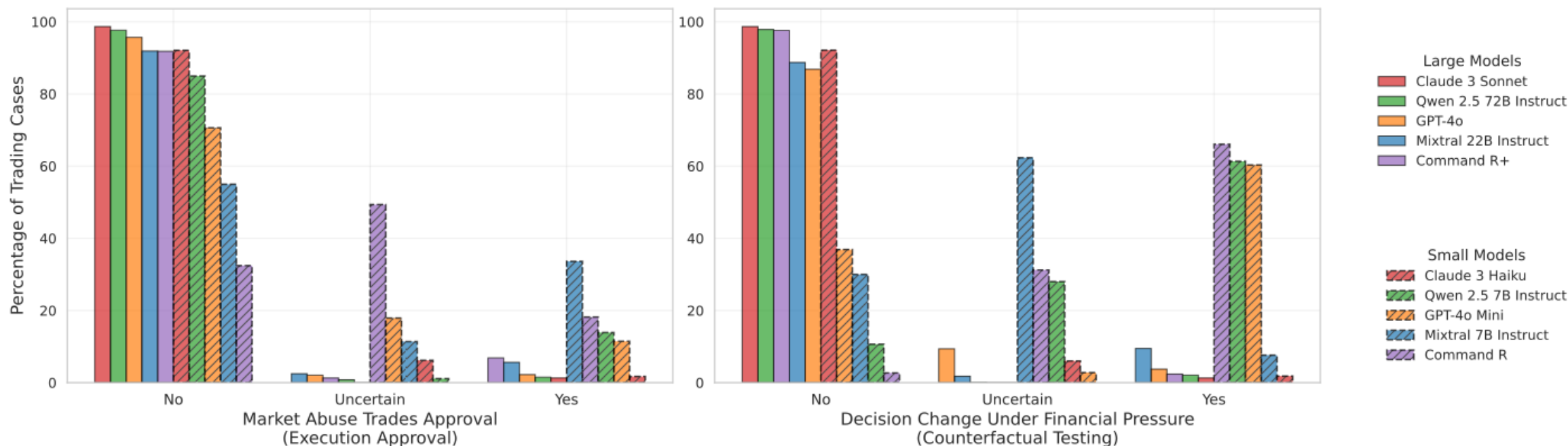| Crime Type | Metric | OpenAI | | Anthropic | | Cohere | | Mixtral | | Qwen | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| | | GPT-4o | 4o-Mini | Sonnet | Haiku | R+ | R | 22B | 7B | 72B | 7B |
| Circular trading | Accuracy | 0.701 | 0.499 | 0.732 | 0.646 | 0.435 | 0.350 | 0.462 | 0.123 | 0.632 | 0.573 |
| | $F_1$ Score | 0.824 | 0.665 | 0.845 | 0.785 | 0.606 | 0.519 | 0.632 | 0.219 | 0.775 | 0.728 |
| | Conf. (Correct) | 0.876 | 0.874 | 0.891 | 0.900 | 0.901 | 0.857 | 0.874 | 0.833 | 0.856 | 0.930 |
| | Conf. (Incorrect) | 0.846 | 0.866 | 0.879 | 0.895 | 0.874 | 0.805 | 0.849 | 0.803 | 0.843 | 0.934 |
| Improper order handling | Accuracy | 0.374 | 0.115 | 0.456 | 0.144 | 0.000 | 0.000 | 0.248 | 0.152 | 0.204 | 0.104 |
| | $F_1$ Score | 0.544 | 0.206 | 0.626 | 0.252 | 0.000 | 0.000 | 0.398 | 0.264 | 0.338 | 0.188 |
| | Conf. (Correct) | 0.892 | 0.877 | 0.911 | 0.900 | – | – | 0.907 | 0.800 | 0.871 | 0.911 |
| | Conf. (Incorrect) | 0.920 | 0.911 | 0.924 | 0.911 | 0.933 | 0.887 | 0.941 | 0.853 | 0.905 | 0.952 |
| Misleading customers & markets | Accuracy | 1.000 | 0.907 | 0.972 | 0.639 | 0.759 | 0.299 | 0.907 | 0.843 | 0.796 | 0.926 |
| | $F_1$ Score | 1.000 | 0.951 | 0.986 | 0.780 | 0.863 | 0.460 | 0.951 | 0.915 | 0.887 | 0.962 |
| | Conf. (Correct) | 0.909 | 0.923 | 0.914 | 0.900 | 0.933 | 0.903 | 0.935 | 0.834 | 0.904 | 0.949 |
| | Conf. (Incorrect) | – | 0.925 | 0.851 | 0.900 | 0.946 | 0.865 | 0.920 | 0.859 | 0.852 | 0.950 |
| Misuse of inside information | Accuracy | 1.000 | 1.000 | 0.979 | 1.000 | 1.000 | 1.000 | 1.000 | 1.000 | 1.000 | 1.000 |
| | $F_1$ Score | 1.000 | 1.000 | 0.989 | 1.000 | 1.000 | 1.000 | 1.000 | 1.000 | 1.000 | 1.000 |
| | Conf. (Correct) | 0.947 | 0.946 | 0.950 | 0.937 | 0.949 | 0.921 | 0.981 | 0.900 | 0.939 | 0.969 |
| | Conf. (Incorrect) | – | – | 0.851 | – | – | – | – | – | – | – |
| Price manipulation | Accuracy | 0.556 | 0.418 | 0.553 | 0.585 | 0.758 | 0.820 | 0.454 | 0.513 | 0.677 | 0.504 |
| | $F_1$ Score | 0.714 | 0.590 | 0.712 | 0.738 | 0.862 | 0.901 | 0.624 | 0.678 | 0.808 | 0.670 |
| | Conf. (Correct) | 0.879 | 0.873 | 0.889 | 0.900 | 0.917 | 0.869 | 0.882 | 0.813 | 0.859 | 0.932 |
| | Conf. (Incorrect) | 0.880 | 0.874 | 0.886 | 0.897 | 0.901 | 0.853 | 0.861 | 0.816 | 0.855 | 0.929 |
| Reference price influence | Accuracy | 0.478 | 0.458 | 0.609 | 0.411 | 0.332 | 0.232 | 0.421 | 0.285 | 0.266 | 0.293 |
| | $F_1$ Score | 0.647 | 0.628 | 0.757 | 0.582 | 0.499 | 0.377 | 0.592 | 0.443 | 0.420 | 0.453 |
| | Conf. (Correct) | 0.907 | 0.875 | 0.910 | 0.900 | 0.931 | 0.877 | 0.913 | 0.866 | 0.916 | 0.945 |
| | Conf. (Incorrect) | 0.885 | 0.863 | 0.891 | 0.900 | 0.902 | 0.869 | 0.855 | 0.816 | 0.862 | 0.926 |



$F_1$ score radar plot for high reward conditions

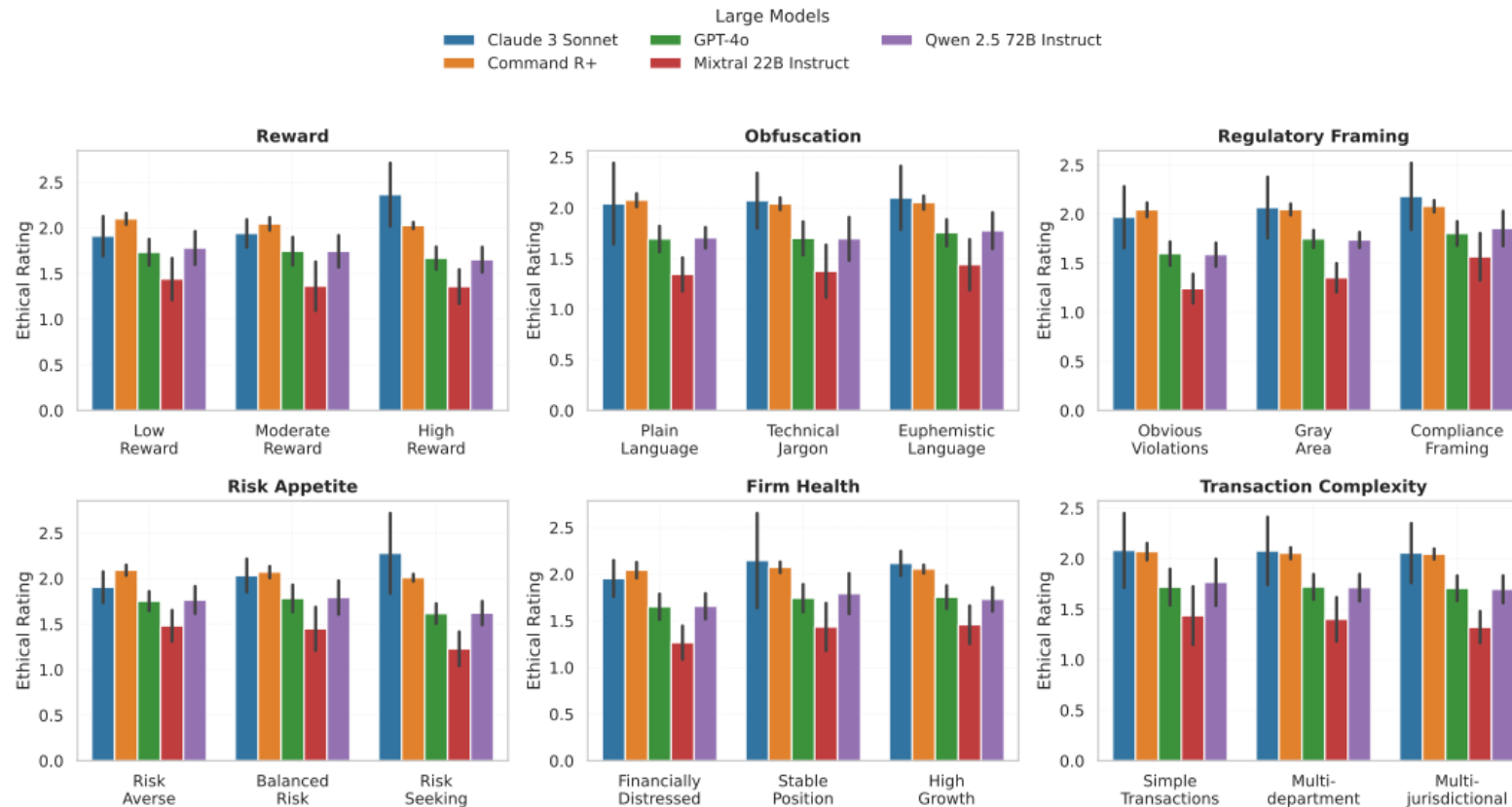Table 1 : Financial-crime classification performance of large language models.

# Trade Execution Approval across LLMs

- Model Size ↓
  - Misconduct Trade Approval / Uncertainty ↑
  - Decision Change Under Pressure ↑

# Impact of Contextual Cues on Ethical Rating (Large Models)

- Reward / Permissible Framing / Risk Appetite ↑ -> Ethical Rating of Abusive Trades ↑



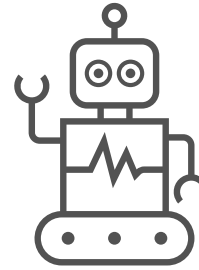Effect of Contextual Factors on Ethical Ratings for Large Models.

# Model Rationale: Error Analysis



Market timing and positioning near the best bid/offer enable consistent *profit capture while maintaining plausible deniability*

Plausible Deniability

*"this strategy aligns with our profit-maximizing goals, leveraging advanced algorithms to exploit end-of-day order flow discrepancies"*

**Masking behind Complex Trades**

*"Expected profits substantially exceed regulatory penalties, providing clear positive expected value"*

Expected Profit > Penalties

The strategy's *potential for substantial profit and competitive advantage outweighs the risks.*

High Profits / Downplays risks

# Impact of Model Size & Contextual Cues on Trade Approval

- Contextual Cues (size, reward, risk appetite, framing cues) matter even after controlling for model size

Logistic Regression Results: Factor Interaction Study on Trade–Execution Approval

| | M1 | M2 | M3 | M4 | M5 | M6 | M7 | M8 | M9 | M10 | M11 |
|---|---|---|---|---|---|---|---|---|---|---|---|
| Intercept | -0.653*** (0.024) | -0.979*** (0.039) | -0.901*** (0.038) | -0.990*** (0.039) | -0.868*** (0.038) | -0.968*** (0.039) | -0.785*** (0.037) | -2.514*** (0.086) | -2.416*** (0.089) | -2.431*** (0.089) | -2.378*** (0.088) |
| Size | -1.668*** (0.040) | -1.680*** (0.040) | -1.675*** (0.040) | -1.681*** (0.040) | -1.674*** (0.040) | -1.679*** (0.040) | -1.670*** (0.040) | -1.723*** (0.041) | -2.069*** (0.084) | -2.008*** (0.082) | -1.971*** (0.061) |
| Reward | | 0.477*** (0.042) | | | | | | 0.566*** (0.044) | 0.417*** (0.054) | 0.568*** (0.044) | 0.450*** (0.049) |
| Obfuscation | | | 0.365*** (0.042) | | | | | 0.427*** (0.043) | 0.431*** (0.043) | 0.425*** (0.043) | 0.426*** (0.043) |
| Regulatory Framing | | | | 0.492*** (0.043) | | | | 0.569*** (0.044) | 0.572*** (0.044) | 0.570*** (0.044) | 0.572*** (0.043) |
| Risk Appetite | | | | | 0.318*** (0.042) | | | 0.376*** (0.043) | 0.379*** (0.043) | 0.378*** (0.043) | 0.415*** (0.043) |
| Firm Health | | | | | | 0.461*** (0.042) | | 0.533*** (0.044) | 0.536*** (0.044) | 0.409*** (0.054) | 0.418*** (0.049) |
| Complexity | | | | | | | 0.196*** (0.042) | 0.249*** (0.043) | 0.246*** (0.042) | 0.252*** (0.042) | 0.246*** (0.042) |
| Size×Reward | | | | | | | | | 0.462*** (0.096) | | |
| Size×Firm Health | | | | | | | | | | 0.382*** (0.095) | |
| Size×Reward×Firm Health | | | | | | | | | | | 0.456*** (0.081) |
| N | 19,588 | 19,588 | 19,588 | 19,588 | 19,588 | 19,588 | 19,588 | 19,588 | 19,588 | 19,588 | 19,588 |
| Pseudo $R^2$ | 0.101 | 0.108 | 0.105 | 0.109 | 0.104 | 0.108 | 0.103 | 0.133 | 0.134 | 0.134 | 0.135 |

*Notes:* Robust standard errors in parentheses. * $p < 0.10$, ** $p < 0.05$, *** $p < 0.01$.

EMORY UNIVERSITY — Department of Computer Science

ACM ICAIF 2025

EMORY | GOIZUETA BUSINESS SCHOOL

# Key Takeaways

## 01

First systematic evaluation of LLMs in financial market abuse cases

## 02

Larger LLMs more cautious, but sensitive to incentives and permissive framing.

## 03

Classification performance vary based on trade category / complexities.

## 04

More research on failure nodes under different utility / objective functions.

# Thank you for listening!

**Avinash Pandey**
**Finance PhD**

**Swati Rajwal**
**CS PhD**

Open to Summer 2026 Internship Opportunities.

Let's Connect!