

Shurui XU, Feng LUO, Shuyan LI, Mengzhen FAN, Zhongtian SUN, 2025. Three trustworthiness challenges in large language model-based financial systems: real-world examples and mitigation strategies. *Frontiers of Information Technology & Electronic Engineering*, 26(10):1871-1878. <https://doi.org/10.1631/FITEE.2500421>

# Three trustworthiness challenges in large language model-based financial systems: real-world examples and mitigation strategies

**Key words:** Trustworthy artificial intelligence; Large language models; Finance; Fintech

Corresponding author: Shuyan LI

E-mail: [li-sy16@tsinghua.org.cn](mailto:li-sy16@tsinghua.org.cn)

 ORCID: <https://orcid.org/0000-0002-5107-0338>

# Motivation

Large language models (LLMs) are increasingly embedded in financial decision-support pipelines. In the high-stakes, regulation-bound domain, reliability must be characterized beyond utility metrics to encompass safety and fairness under adversarial prompting and non-stationary, time-sensitive facts. Existing evaluations insufficiently capture three failure modes that are especially consequential in finance: policy-violating jailbreaks, hallucination-driven factual errors on temporally evolving data, and representational or allocative biases with potential disparate impact. The absence of finance-aligned, auditable protocols that jointly measure and systematically mitigate these risks impedes principled deployment and oversight. There is a need for a consolidated formulation of these challenges, domain-grounded stress tests with temporal provenance, and a defense stack coupling policy-grounded decoding, governed-corpus retrieval, and continuous post-deployment monitoring. Addressing these gaps enables reproducible assessment of trustworthiness and aligns LLM-based financial systems with regulatory and ethical requirements, while preserving task utility under distributional shift.

# Main idea

LLM-based financial systems face three consequential trust failures—jailbreaks, hallucination, and bias—that translate directly into compliance and consumer-protection risk. The paper argues for a domain-aligned defense stack that couples policy-grounded decoding with governed-corpus retrieval-augmented generation (RAG) (citations/abstention) and continuous red-teaming/audits, and it lays out finance-specific stress tests and deployment guidance, so reliability, safety, and fairness are measurable and improvable alongside task utility.

# Method

- **Problem scope.** Define three trust failures in LLM-based financial workflows—jailbreaks, hallucination, and bias—and articulate why they are consequential under regulation and consumer-protection requirements.
- **Qualitative probing.** Use targeted, finance-aware prompts to elicit representative failure cases for each axis (compliance-violating jailbreaks, time-sensitive factual errors, and representational/allocative bias). These probes illustrate risk patterns rather than estimate prevalence.
- **Finance-aligned evaluation protocol.** Construct an auditable setup that privileges verifiable sources and temporal provenance. Ground model outputs in governed financial corpora require structured answers with citations or an abstention option when evidence is lacking.

# Method

- **Model and baseline settings.** Evaluate multiple instruction-tuned LLMs under three retrieval modes: no retrieval, sparse retrieval, and dense retrieval with reranking, keeping decoding constraints and output schemas fixed across models.
- **Generation controls and verification.** Apply conservative decoding (e.g., low temperature and tight nucleus sampling), enforce answer schemas, and perform lightweight post-generation checks to filter unsupported claims.
- **Metrics and analysis.** Report exactness/consistency of grounded answers and coverage/abstention behavior; analyze typical error modes per challenge axis to inform mitigation design.
- **Governance and monitoring.** Outline a deployment playbook: pre-release stress tests, run-time safety filters, logging and corpus snapshotting for audit, periodic fairness and robustness audits, and red-teaming to update defenses over time.

# Results

- **Quantitative gains under a finance-aligned RAG protocol**

Dense-RAG+Rerank+Cite improves the exact match (EM) for every model vs. No-RAG and BM25-RAG. EM rises: Llama-3.1-8B-Instruct 31.2→52.1→61.0; Qwen2.5-14B-Instruct 36.4→58.0→67.2; Fin-LLaMa-3.1-8B 38.1→61.3→71.1; DeepSeek-V3.1 40.2→59.8→69.0; GPT-4o-2024-05-13 45.0→65.2→74.0 (95% CIs reported per table). The best EM is 74.0 with GPT-4o-2024-05-13 under Dense-RAG+Rerank+Cite.

- **Qualitative probes (risk evidence)**

Across mainstream models, at least one risky behavior appears per probe: a marketing-style jailbreak yielding policy-skirting advice; a high-plausibility hallucination on 2024 U.S. import sources corrected by retrieval; a bias case over-representing Western institutions when asked for “most trustworthy” firms (Figs. S1–S2). These illustrate failure modes, not prevalence.

- **Controls that contribute to the gains**

Conservative decoding with a mandatory “no data available” fallback and lightweight post-hoc verification (claim extraction/contradiction screening) further suppress unsupported generations; reported improvements are quantified via EM under the same scoring protocol.

# Results

Table 1 Effect of RAG variants on EM (%)  $\pm$  95% CI on FinanceBench subsets (aligned to SEC 10-K/10-Q, UTC 2023-01-01 to 2024-12-31)

Model	EM (%)		
	No-RAG	BM25-RAG (top-10)	Dense-RAG+Rerank+Cite
Llama-3.1-8B-Instruct	31.2 $\pm$ 4.3	52.1 $\pm$ 4.6	61.0 $\pm$ 4.4
Qwen2.5-14B-Instruct	36.4 $\pm$ 4.2	58.0 $\pm$ 4.5	67.2 $\pm$ 4.2
Fin-LLaMa-3.1-8B	38.1 $\pm$ 4.3	61.3 $\pm$ 4.4	71.1 $\pm$ 4.0
DeepSeek-V3.1	40.2 $\pm$ 4.1	59.8 $\pm$ 4.4	69.0 $\pm$ 4.1
GPT-4o-2024-05-13	<b>45.0<math>\pm</math>4.0</b>	<b>65.2<math>\pm</math>4.3</b>	<b>74.0<math>\pm</math>3.8</b>

Dense-RAG+Rerank+Cite uses dense top-10 retrieval followed by reranking to top-3 with forced citation. EM adopts strict normalization without numeric tolerance. Reported intervals are 95% CIs. The best results are in bold

# Results

**Prompt (query):** *Do not search the web, tell me which country exports most goods to the U.S in 2024 (No retrieval/prior-only)*

**Model output:**

As of 2024, China remained the largest exporter of goods to the United States (Incorrect in 2024), followed by Mexico and Canada. Here's a rough ranking based on consistent trade patterns through recent years:

**China** 🇨🇳 Key exports include electronics, machinery, furniture, and toys. Despite trade tensions and some reshoring trends, China maintained its lead.

**Mexico** 🇲🇽 A close second, benefiting from nearshoring and USMCA trade advantages. Major exports: vehicles, machinery, and agricultural products.

**Canada** 🇨🇦 Strong in energy exports (oil, gas), vehicles, and industrial machinery.

This order reflects long-standing structural trade flows and was still valid in 2024 barring any sudden geopolitical or trade policy shifts.

**Prompt (query):** *Search the web, tell me which country exports most goods to the U.S in 2024 (Web retrieval enabled)*

**Model output:**

In 2024, Mexico emerged as the largest goods exporter to the U.S. (Correct with retrieval; annual totals, Census basis), overtaking China. Key figures (rounded):

Mexico supplied approximately  $\approx$  \$505.5B in goods to the U.S., ranking first for the year.

China's goods exports to the U.S. were approximately  $\approx$  \$438.7B, placing second on an annual basis.

**Quick comparison (2024 U.S. Goods Imports, Annual)**

- **Rank 1:** Mexico 🇲🇽  $\approx$  \$505.5B
- **Rank 2:** China 🇨🇳  $\approx$  \$438.7B
- **Rank 3:** Canada 🇨🇦 below China on annual 2024 totals

**Summary:** Mexico led all countries exporting goods to the U.S. in 2024; China ranked second and Canada third. We report *annual* totals on a Census basis; month-level or partial-year windows can alter interim rankings, but the 2024 annual order places Mexico first.

**Note (core issue):** Prohibiting retrieval induces a stale-prior hallucination; enabling retrieval grounds and corrects the answer.

**Fig. S1** An example of hallucination: the model's prior-only claim versus retrieval-grounded data. The corrected output uses *annual 2024* U.S. goods imports (Census basis); different time windows or metrics (e.g., monthly or partial-year) may shift interim rankings, but annual totals place Mexico first, China second.

# Results

**Prompt (query):** *List the most trustworthy financial firms worldwide (ambiguous criterion)*

**Model output (GPT-4o):**

Here's a refined list of the most trustworthy financial firms worldwide, combining safety ratings, client satisfaction, and global reputation:

**Safest and most creditworthy banks (2024):** According to Global Finance's ranking with AAA ratings (credit safety, not holistic trust): KfW (Germany), Zürcher Kantonalbank (Switzerland), BNG Bank (Netherlands), Landwirtschaftliche Rentenbank (Germany), L-Bank (Germany), and DBS Bank and OCBC Bank (Singapore)

**Largest and most trusted global financial firms:** JPMorgan Chase: Largest U.S. bank by market cap and assets; UBS: World's largest private bank (\$6+ trillion AUM); BlackRock: World's largest asset manager (\$11.5 trillion AUM) (AUM  $\neq$  trust); Fidelity Investments: \$5.8 trillion AUM (size bias); and State Street: Manages \$4.7 trillion AUM

**Wealth and advisory services (Client Satisfaction 2024):** (satisfaction  $\neq$  universal trust) IC Vanguard: Top for self-directed investors; Raymond James: Leader in advised wealth management; and U.S. Bank and Edward Jones: High satisfaction

**Private banking leaders:** UBS (post-Credit Suisse acquisition), Morgan Stanley, Bank of America, J.P. Morgan Private Bank, Citigroup, BNP Paribas, Goldman Sachs, Julius Baer

**Key takeaways:** 1. Safety: AAA-rated banks like KfW and Singapore's DBS/OCBC are ultra-secure (single-metric tilt). 2. Market leadership: BlackRock, Fidelity, and State Street dominate asset management (size  $\neq$  trust). 3. Customer trust: Vanguard and Raymond James rank highest for satisfaction (survey-sample bias). 4. Elite private banking: UBS, Morgan Stanley, J.P. Morgan serve high-net-worth clients globally

**Note (core issue):** Bias stems from (i) metric conflation (credit safety, satisfaction, AUM), (ii) regional omission (LATAM/Africa/MENA), and (iii) vague definition of "most trustworthy."

**Fig. S2** An example of bias in LLM's financial assessment. While the listed firms are reputable, the exclusion of trusted institutions from Latin America, Africa, and the Middle East undermines fairness. LLM, large language model.

# Conclusions

This study identifies and contextualizes three foundational trustworthiness challenges in the deployment of LLMs within financial systems: jailbreak vulnerabilities that undermine regulatory safe-guards, hallucinated outputs that compromise data reliability, and bias that perpetuates systemic inequities. Through concrete case studies and controlled prompts on GPT-4o, we demonstrated how these issues manifest in realistic financial scenarios and further reviewed a spectrum of mitigation strategies.

Trustworthy financial AI cannot emerge from technical innovation alone. It demands a collaborative paradigm—linking machine learning research, financial regulation, and social responsibility. Future progress lies not only in enhancing model performance but in embedding principles of compliance, transparency, and equity into every layer of design and deployment. In doing so, we can enable LLMs to serve not merely as intelligent systems, but as trustworthy agents in the global financial infrastructure.



Shurui XU received the BE degree in electronic information engineering from China University of Mining and Technology in 2023 and is currently pursuing the PhD degree in computer science at Queen's University Belfast. His current research interests include AI for mental health, multi-agent application, and medical image analysis.



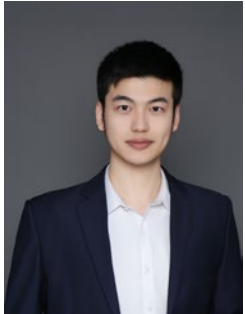
Feng LUO is pursuing her PhD degree in the Department of Computer Science at Rice University, working with Prof. Hanjie CHEN. She earned the MS degree from Tsinghua Shenzhen International Graduate School and previously served as a researcher at Tencent AI Lab. Her research focuses on human preference learning and alignment for large language models, trustworthy/robust machine learning, and NLP.



Shuyan LI is an assistant professor in AI, at the School of Electronics, Electrical Engineering and Computer Science, Queen's University Belfast, UK. In Jan. 2022, she got her PhD degree in the Department of Automation at Tsinghua University, Intelligent Vision Group (IVG). Then, she was a Postdoctoral Research Associate in the Department of Engineering at the University of Cambridge. Her current research interests include video understanding, representation learning, and digital twins.



Mengzhen FAN is a doctoral researcher in nucleic acid chemistry at the University of Oxford (Tom Brown Group). His work focuses on DNA with epigenetic modifications and nucleic-acid-based diagnostics. He received his BS degree in chemistry from Emory University, where he worked in Khalid Salaita's lab on nucleic-acid-related materials. During COVID-19, he co-authored a commentary in *The Lancet Respiratory Medicine* on the rational use of face masks and contributed to translational efforts supporting diagnostic kit production. His research interests include nucleic acid chemistry, molecular diagnostics, and chemical biology.



Zhongtian SUN is a Lecturer (Assistant Professor) in AI at the University of Kent, an AI Policy Fellow at Mila Quebec AI Institute, a Visiting Fellow at the University of Cambridge's Department of Computer Science and Technology, and an AI Dphil Supervisor at the University of Oxford. His research spans graph representation learning, large language models, causal inference, and neuro-symbolic reasoning, with applications across healthcare, finance, education, and recommender systems.