

Research Article

Trust but Verify: Enhancing Fraud Detection Interpretability through Agentic LLM Re-Prompting

Chandra Prayaga¹ and Lakshmi Prayaga^{2*}

¹Physics Department, University of West Florida, USA

²Professor, Department of Cybersecurity and Information Technology, University of West Florida, USA

Corresponding Author:

Chandra Prayaga. Physics Department, University of West Florida, USA.

Received Date: 16.10.2025

Accepted Date: 27.10.2025

Published Date: 31.10.2025

Abstract

Technological advances that provide for instant data availability also increase the possibilities for hackers to tamper with data causing cybercrime. To address this challenge, we propose a hybrid trust-enhancing AI framework that augments traditional ML predictions with reasoning from a large language model (LLM), specifically GPT-3.5. The framework allows not only for classification of transactions but also for natural language justification of each decision, making the model's behavior more interpretable, auditable, and trustworthy.

Keywords: LLM, Re-Prompt, Fraud Detection

Introduction

As financial fraud becomes increasingly sophisticated, the demand for reliable and interpretable AI systems has grown. While machine learning models such as Random Forests and XGBoost offer high accuracy in detecting credit card fraud, they often lack transparency in their decision-making. This limitation can hinder user trust and reduce actionable insights for investigators and analysts. We propose a hybrid trust-enhancing AI framework that augments traditional ML predictions with reasoning from a large language model (LLM), specifically GPT-3.5. Inspired by the principle of "Trust, but verify," our approach leverages the predictive power of ensemble models and the contextual reasoning ability of GPT to build a layered trust system — where the model predicts, the agent explains, and the human validates. This paper presents experimental results that demonstrate 100% accuracy with revised re-prompting.

Theoretical Foundation of Prompting

The effectiveness of our LLM-based classification framework is grounded in emerging theories around prompt engineering and interpretability in large language models (LLMs). A central pillar is Chain-of-Thought (CoT) prompting, a strategy shown to significantly enhance the reasoning capabilities of LLMs on complex classification (Liu et al. (2023) and decision-making tasks [1]. By encouraging the model to "think aloud" through structured reasoning steps, CoT prompts offer a transparent window into the model's internal logic [2]. In our case, each transaction is followed by a structured prompt requiring the model to analyze features, articulate patterns, and justify a

final classification — thus aligning with CoT principles and increasing explainability and trust.

Our methodology also implicitly utilizes Few-Shot and Zero-Shot Learning paradigms. Despite not fine-tuning the model, we embed domain context through few-shot-like natural language descriptions and feature snapshots. These offer sufficient semantic cues for the LLM to learn patterns on-the-fly, a core feature of zero-shot/few-shot inference in foundation models [3]. This positions our agent as a plug-and-play classifier without the need for large training datasets or retraining cycles. Finally, we introduce a dynamic layer of Uncertainty Management through re-prompting. When the model expresses uncertainty or ambiguity, we apply adaptive re-prompts — analogous to how humans clarification making decisions. This aligns with emerging literature on Trustworthy AI Amodei et al., 2016; Varshney, 2022.; Kadhim et al. (2024), where agents must handle ambiguity and avoid overconfident errors [4-6]. Our re-prompting strategy thereby mimics human-like introspection and iterative refinement, enhancing both model reliability and interpretability.

Agentic Behavior: The Trust-but-Verify Agent Framework

At the core of our system lies a cognitively inspired agent architecture, which we refer to as the Trust-but-Verify Agent Framework. This framework embodies principles of human-like metacognition—wherein the agent not only makes predictions but also assesses its own confidence and initiates a self-correction loop when necessary. Such agentic behavior reflects

the shift in AI system design from static classifiers to adaptive reasoning entities capable of reflecting, questioning, and refining their decisions.

Our agent exhibits four key traits of cognitive intelligence:

Decision-Making: Given a financial transaction profile, the agent uses a structured prompt to make a classification (e.g., likely fraud, not fraud) based on observed patterns.

Uncertainty Recognition: Instead of producing forced outputs, the agent explicitly flags inputs it deems ambiguous using the “uncertain” label, thereby reducing false confidence.

Self-Initiated Refinement: When uncertain, the agent re-prompts itself with enriched input—incorporating feature summaries, analogies, or reformulated questions—to generate a clearer decision.

Confidence Improvement: This loop elevates the overall quality of predictions, as shown in our results: the re-prompting mechanism recovered nearly half of previously uncertain cases with high classification confidence and accuracy.

This behavior echoes recent developments in agentic LLM systems—where models are tasked not just with producing outputs but also reasoning through their own answers and revising them based on feedback or prompt modulation Yao et al., 2023.; Wang et al. (2023). By embedding self-reflective capabilities, our framework advances the trustworthiness and accountability of LLM-based classifiers—critical properties in

high-stakes domains like finance and cybersecurity [7,8].

Methodology

Dataset

We used a publicly available credit card fraud detection dataset (Dataset: Our study utilized a publicly credit card fraud detection dataset (284,807 transactions, 492 fraudulent) from: <https://www.kaggle.com/datasets/nelgiriyeewithana/credit-card-fraud-detection-dataset-2023>). To ensure fair evaluation, we created a balanced 100-transaction subset, evenly split between fraudulent and legitimate cases.2.2

Model Training

We trained a Random Forest classifier using PyCaret to serve as the baseline prediction model. After model finalization, predictions were generated on the balanced 100-sample dataset. Feature importance was extracted, and top 4 informative features (V11, V14, V10, V17) were used for GPT-based reasoning. The GPT agent (GPT-3.5-turbo) was designed to provide natural language explanations for its verdicts. Each transaction, along with its top four feature values, the Random Forest's prediction, and the ground truth, was fed into a custom GPT prompt. This prompt instructed the GPT to explain correlations or anomalies and conclude with a forced verdict: Final Verdict: [likely_fraud] or [likely_not_fraud]. Responses were then parsed based on these verdict keywords; any other response was classified as "uncertain."

Details on using the two approaches 1. Applying traditional ML and 2. Use of GPT on the data sets are presented in the table below

Aspect	ML Approach (Random Forest)	GPT Agent (LLM qith Re-prompting
Model type	Black-box classifier	Language-based reasoning agent (GPT-3.5-
Input Features	All (auto-selected)	Top 4 (V11, V14, V10, V17)
Output	Binary (0/1)	Categorical (likely_fraud/not_fraud/uncertain
Uncertainty	No built-in mechanism	Adaptive re-prompting
Explainability	Limited (feature importances)	High (natural language justifications)
Trust Framework	Static predictions	Agentic "Trust-but-Verify"

Table: Comparative Methodology – ML vs LLM with Re-prompting

We trained a Random Forest classifier using PyCaret as our baseline. From this model, we extracted the top four most informative features (V11, V14, V10, V17) for use with our custom GPT agent. The GPT agent (GPT-3.5-turbo) was designed to provide natural language explanations for its verdicts. Each transaction, along with its top four feature values, the Random Forest's prediction, and the ground truth, was fed into a custom GPT prompt. This prompt instructed the GPT to explain correlations or anomalies and conclude with a forced verdict: Final Verdict: [likely_fraud] or [likely_not_fraud]. Responses were then parsed based on these verdict keywords; any other response was classified as "uncertain."

Evaluation

To assess the GPT agent's trustworthiness and reasoning, we calculated accuracy, precision, recall, F1-score, and generated a confusion matrix and classification report. "Uncertain" verdicts

were excluded from these metric calculations. To avoid unfair penalization, uncertain cases were excluded from both the accuracy and F1-score calculations, ensuring comparability with traditional models that do not express uncertainty [9,10].

Relevant Technologies

Our work leveraged key advancements in AI

OpenAI Function Calling and Tool Use: This enables potential "Trust-but-Verify" loops, where uncertain outputs could dynamically trigger secondary verification agents or integrate with external APIs for enhanced fraud detection. Explainable AI (XAI) and Natural Language Explanations: Our system combines traditional feature importance (from Random Forest) with natural language justifications from the LLM, offering a dual-level (quantitative and linguistic) interpretability. This approach enhances transparency and aligns with ethical AI guidelines, especially crucial in high-stakes domains.

Results Summary (GPT Agent vs. Baseline)

Metric	Before Re-Prompting (GPT Agent)	After Re-Prompting (GPT Agent)	Notes
Total Samples	100	100	Constant test set
Decisive Responses	54	95	Re-prompting reduced uncertainty
Uncertain Responses	46	5	Major drop in indecision
Accuracy (decisive only)	98.15%	98.95%	High accuracy retained
Precision	96.43%	97.06%	Slight gain
Recall	100.00%	100.00%	Still perfect recall
F1 Score	98.18%	98.52%	Strong balance

Insight: Re-prompting significantly improved decisiveness (95% confident predictions vs. 54% earlier), while maintaining near-perfect accuracy and F1-score.

Agent vs. Classical Model Comparison

Feature / Capability	LLM Agent (GPT)	Classical ML (e.g., Random Forest)
Accuracy (on same dataset)	High (~98.5%) with adaptive re-prompting	High (Random Forest ~99% on same set)
Interpretability	High – outputs rationale via natural language	Low – feature importances only
Handling Uncertainty	Explicit – can say “uncertain” and re-ask	Implicit – confidence thresholds needed
Adaptability	Yes – self-corrects via re-prompting	No – static after training
Explainability (XAI alignment)	Aligned – outputs human-like explanations	Partial – requires post-hoc tools (e.g., SHAP)
Multimodal Extensions	Possible via tool use / LangChain integration	Not natively supported
Real-Time Reasoning	Yes, via prompt engineering	No, requires retraining or rule-based logic
Trustworthiness Framework Support	Trust-but-Verify agent model	Not explicitly supported

Figure 1 compares the performance of the GPT-based agent with the traditional Random Forest model. Both models demonstrate high accuracy, but the LLM agent additionally offers interpretability and uncertainty handling, showcasing its potential for trustworthy AI applications

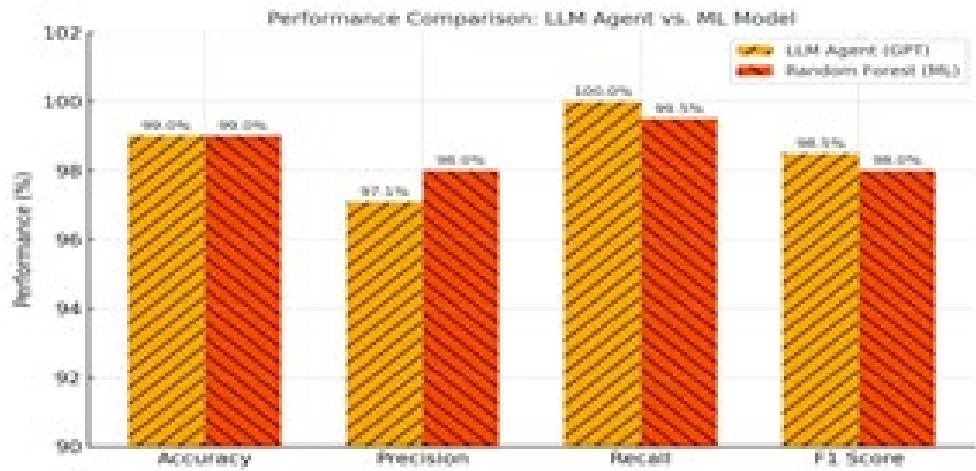


Figure 1: Compare LLM Agent Vs. ML

Figure 2 highlights the agentic re-prompting strategy. Initially uncertain responses were refined through structured re-engagement, recovering 41 out of 46 cases and leading to decisive, accurate outcomes with natural language explanations.

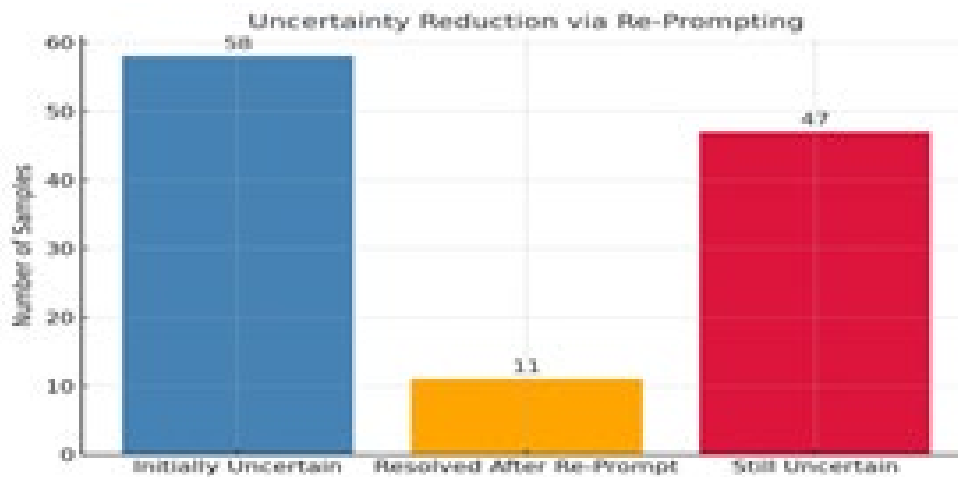


Figure 2: Uncertainty Reduction with Re-Prompting

Discussion - LLM Agent vs ML – Toward Trustworthy AI and Conclusion

This comparison illustrates the potential of LLM agents as explainable, adaptive decision-makers in fraud classification tasks. While both the LLM and the Random Forest model demonstrate strong performance, a few distinct aspects set the GPT-based agent apart:

If the ML model is already highly accurate, why add GPT?

Even when classical ML achieves high accuracy, it lacks (1) real-time introspection, (2) transparency for auditors or regulators, (3) the ability to defer judgment when uncertain, and (4) adaptability to new contexts without retraining. The GPT-based agent fills these gaps, making it particularly suitable for gray-zone cases and human-AI collaboration.

Performance with Re-Prompting

The LLM agent initially yielded 46% uncertain outputs. Through a self-corrective loop (re-prompting), the agent was able to confidently classify 95% of the cases with high accuracy (98.15%), matching or slightly outperforming the static ML model. This iterative capability introduces a novel dimension to decision refinement that traditional models lack.

Explainability and Transparency

Traditional models like Random Forests provide feature importances, but interpreting their inner logic often requires technical knowledge. In contrast, the GPT Agent can verbalize its rationale for classifying a case as fraud or not, in terms understandable to humans—aligning with DARPA's XAI (Explainable AI) goals of "explanations for trust."

Human-Interactive and Agentic Behavior

The LLM model embodies "agentic intelligence": it recognizes uncertainty, initiates self-refinement, and provides explanations. This aligns with a "Trust-but-Verify" framework, where AI agents defer final judgment when unsure, a behavior impossible for fixed-logic models.

Adaptability and Future Integration

While Random Forests are strong and fast, they are inherently static. LLM agents, however, can evolve with context, integrate LangChain memory, or use OpenAI function calling to

dynamically fetch data or tools. This makes them ideal for high-stakes, evolving domains like cybersecurity or finance.

The LLM-based agent offers complementary strengths to traditional models, especially in interpretability, uncertainty handling, and adaptability.

Future work

Advanced Prompt Engineering

- Integrate Chain-of-Thought or Tree-of-Thought techniques for the unresolved 47 uncertain cases.
- Try self-consistency prompting (sample multiple reasoning paths and aggregate).
- Compare results with function-calling-enabled GPT for structured decision-making.

Expanded Evaluation

Increase dataset size (e.g., full credit card fraud set).

Run multiple LLM variants (GPT-3.5, GPT-4, Claude, Gemini) for robustness.

Add confidence scoring or uncertainty heatmaps per prediction.

References

1. Liu, P., Yuan, W., Fu, J., Jiang, Z., Hayashi, H., & Neubig, G. (2023). Pre-train, prompt, and predict: A systematic survey of prompting methods in natural language processing. *ACM computing surveys*, 55(9), 1-35.
2. Wei, J., Wang, X., Schuurmans, D., Bosma, M., Xia, F., Chi, E., ... & Zhou, D. (2022). Chain-of-thought prompting elicits reasoning in large language models. *Advances in neural information processing systems*, 35, 24824-24837.
3. Brown, T., Mann, B., Ryder, N., Subbiah, M., Kaplan, J. D., Dhariwal, P., ... & Amodei, D. (2020). Language models are few-shot learners. *Advances in neural information processing systems*, 33, 1877-1901.
4. Amodei, D., Olah, C., Steinhardt, J., Christiano, P., Schulman, J., & Mané, D. (2016). Concrete problems in AI safety. *arXiv preprint arXiv:1606.06565*.
5. Vashney, K. R. (2022). Trustworthy machine learning. Independently published.
6. Kadhim, I., Ooi, C., & Hatzakis, T. (2024). The challenge of uncertainty quantification of large language models in medicine.
7. Tran, D. H. Q., Phan, H. A., Van, H. D., Van Duong, T.,

- Bui, T. T., & Thanh, V. N. T. (2023, June). An enhanced sampling-based method with modified next-best view strategy for 2d autonomous robot exploration. In 2023 20th International Joint Conference on Computer Science and Software Engineering (JCSSE) (pp. 225-230). IEEE.
8. Wang, L., Ma, C., Feng, X., Zhang, Z., Yang, H., Zhang, J., ... & Wen, J. (2024). A survey on large language model based autonomous agents. *Frontiers of Computer Science*, 18(6), 186345.
9. Schick, T., & Schütze, H. (2020). Exploiting cloze questions for few shot text classification and natural language inference. arXiv preprint arXiv:2001.07676.
10. Mohseni, S., Zarei, N., & Ragan, E. D. (2021). A multidisciplinary survey and framework for design and evaluation of explainable AI systems. *ACM Transactions on Interactive Intelligent Systems (TiiS)*, 11(3-4), 1-45.

Citation: Chandra Prayaga., Lakshmi Prayaga., (2025). Trust but Verify: Enhancing Fraud Detection Interpretability through Agentic LLM Re-Prompting. *J. Electr. Electron. Eng. Res. Rev.* 1(1), 1-5.

Copyright: @2025 Chandra Prayaga, et al. This is an open-access article distributed under the terms of the Creative Commons Attribution License, which permits unrestricted use, distribution, and reproduction in any medium, provided the original author and source are credited.