---

| **RESEARCH ARTICLE**

# Human-in-the-Loop AI in Financial Services: Data Engineering That Enables Judgment at Scale

**Rahul Joshi**

*IIT Kharagpur, India*

**Corresponding Author:** Rahul Joshi, **E-mail**: reachrahuljoshi@gmail.com

---

| **ABSTRACT**

A person involved in the process of AI systems combines human knowledge with machine learning skills to produce hybrid decision-making platforms that meet the intricate needs of risk management, regulatory compliance, and customer protection. This is a major technical advancement in the financial services industry. Financial organisations can benefit from the efficiency of automated processing while maintaining the contextual knowledge and regulatory accountability that human analysts provide, thanks to these systems, which purposefully incorporate human judgement at critical decision points. The architectural underpinnings of these hybrid systems necessitate advanced data engineering solutions that facilitate real-time streaming for prompt decision-making as well as extensive historical analysis functions. Infrastructure for data traceability and explainability is crucial for ensuring regulatory compliance, necessitating unalterable audit trails that document both automated decisions and human inputs, justifications, and alterations to outcomes. The incorporation of feedback systems facilitates ongoing learning processes, allowing human knowledge to improve algorithm performance via organized decision collection and model adjustment pathways. Scalability issues require smart workload management solutions that can address unpredictable processing needs while ensuring peak performance among distributed analyst teams, necessitating advanced caching techniques and pre-computation abilities that facilitate quick case resolution without sacrificing decision quality.

| **KEYWORDS**

Human-in-the-loop systems, financial services AI, data engineering architecture, explainable artificial intelligence, continuous learning pipelines

| **ARTICLE INFORMATION**

---

**Introduction**

Financial services organizations encounter a core difficulty: attaining scale while preserving the human insight necessary for regulatory adherence, risk control, and safeguarding customers. The changing environment of financial fraud introduces unparalleled challenges, with check fraud standing out as a notably enduring risk even amid the digital evolution of payment methods. Recent industry analysis indicates that companies frequently undervalue the intricate characteristics of contemporary check fraud schemes, which have progressed from basic forgery to elaborate operations utilizing advanced counterfeiting methods and social engineering strategies [1]. Human-in-the-loop AI systems embody a strategic method that merges machine learning effectiveness with human supervision, which is especially important in crucial situations like fraud detection, credit assessments, and anti-money laundering inquiries.

In contrast to completely automated systems, HITL architectures intentionally include decision points that allow human analysts to approve, modify, or enhance machine-generated results. The incorporation of artificial intelligence in anti-money laundering monitoring has shown the vital role of human skill in understanding algorithmic results within regulatory environments.

Contemporary AML systems produce thousands of alerts each day, necessitating human analysts to examine intricate transaction trends, evaluate suspicious activity reports, and make decisions that have substantial legal and compliance consequences [2]. This method recognizes that although AI is proficient in identifying patterns and analyzing large data sets, human knowledge is crucial for contextual insights, addressing edge cases, and ensuring regulatory compliance.

The complexity of modern financial offenses requires a refined method that integrates technological tools with human insights. Check fraud schemes now encompass various layers of deceit, featuring seemingly legitimate business documents, advanced printing technologies, and synchronized attacks targeting several financial institutions. These operations' complexity frequently necessitates human analysts to assemble seemingly disconnected information, recognizing links that automated systems may overlook because of their dependence on set patterns and guidelines [1]. Likewise, investigations into money laundering often require examining intricate ownership frameworks, grasping cultural and regional business norms, and deciphering regulatory obligations that differ greatly among jurisdictions.

The effectiveness of these hybrid systems relies significantly on the foundational data engineering framework that enables smooth cooperation between algorithmic processing and human understanding. Contemporary AML monitoring systems are required to handle vast amounts of transaction data while providing the adaptability needed for human investigators to delve into individual cases, cross-check various data sources, and create thorough narratives for compliance reporting [2]. The infrastructure should support the rapid processing demands of machine learning algorithms alongside the intricate, contextual information requirements of human analysts who need to justify their decisions under regulatory examination.

Data engineering issues in HITL systems go beyond conventional processing methods to include the combination of structured transactional data with unstructured information sources, real-time alert creation, and detailed audit trails. The system should assist researchers who might require tracking transaction flows among various institutions, examining patterns over long durations, and linking seemingly unrelated events to uncover complex money laundering operations. This necessitates data systems capable of preserving detailed transaction information while offering summarized perspectives that facilitate pattern detection and trend analysis across extensive datasets [2].

## Architectural Foundations for Human-AI Collaboration

### Real-Time Data Ingestion and Processing

Successful HITL systems need data architectures that facilitate batch processing for historical insights and real-time streaming for prompt decision-making support. The advancement of data integration platforms has greatly improved the ability of financial institutions to handle intricate data ecosystems, with contemporary solutions offering improved connectivity to various financial data sources such as core banking systems, payment processors, and regulatory reporting systems [3]. Transaction logs, customer profiles, external risk databases, and regulatory feeds are just a few of the data sources that the ingestion layer must handle. This often calls for sophisticated transformation capabilities to standardise data formats and resolve schema conflicts between antiquated systems and contemporary applications.

Modern financial organizations encounter the difficulty of consolidating information from numerous distinct systems, all functioning with varying protocols, data formats, and update intervals. The intricacy of this integration is evident in real-time fraud detection situations, where transaction approval decisions need to occur in milliseconds while also enhancing the decision context with past customer behavior, merchant risk assessments, and external threat intelligence sources [3]. Stream processing frameworks facilitate the prompt enhancement of incoming data with contextual details, ensuring that analysts obtain thorough case perspectives shortly after event triggers. However, this necessitates advanced orchestration to handle the timing alignment of data streams that may exhibit different latency traits and update behaviors.

The architecture generally utilizes event-driven patterns in which machine learning models produce scored events that pass through priority queues, establishing a dynamic routing system capable of adjusting to shifting risk environments and operational conditions. High-risk cases are sent directly to human reviewers via specialized processing channels that circumvent regular queue management protocols, whereas lower-risk items might go through extra automated processing or batch review cycles that allow for more comprehensive data collection and analysis procedures [3]. This structured method guarantees that important decisions are prioritized immediately while preserving operational efficiency; however, it necessitates diligent oversight to avoid bottlenecks during times of elevated alert volume when the system could produce far more high-priority cases than standard operational capacity can support.

### Feature Engineering for Analyst Workflows

Data engineering teams need to construct feature pipelines that fulfill two roles: supplying machine learning models and delivering insightful information to human analysts, which demands a thorough comprehension of both algorithmic needs and

human cognitive functions. Feature stores have become essential infrastructure elements that allow organizations to ensure alignment between model training environments and production systems, while also accommodating the intricate demands of human-AI collaboration processes. Features designed for human use frequently vary considerably from those that are fine-tuned for algorithmic processing, as the mathematical accuracy essential for model performance must be weighed against the interpretability requirements of analysts making justifiable choices amidst regulatory inspection.

The conversion of raw data into usable features entails various phases of computation, aggregation, and enrichment that need to be meticulously coordinated to preserve both performance and precision. Contemporary feature engineering platforms offer functionalities for instantaneous feature calculation, allowing for the creation of intricate behavioral signals and risk indicators that can be utilized right away by machine learning models and analyst interfaces [4]. Although models might depend on numerous numerical features obtained from advanced statistical computations and network analysis techniques, analysts need well-organized presentations that emphasize the most important patterns and offer adequate context for exploration and decision-making activities.

Feature stores are essential for sustaining uniform definitions between model training and analyst interfaces, acting as the definitive repository for feature logic that guarantees reproducibility and auditability across various system components. These platforms allow data scientists and engineers to establish features a single time and utilize them across various applications, minimizing the potential for inconsistencies that might compromise model performance or analyst trust [4]. The difficulty of ensuring feature consistency is especially complicated in financial services settings where regulatory demands require thorough audit trails and the capability to replicate past feature calculations for compliance reporting and model validation reasons.
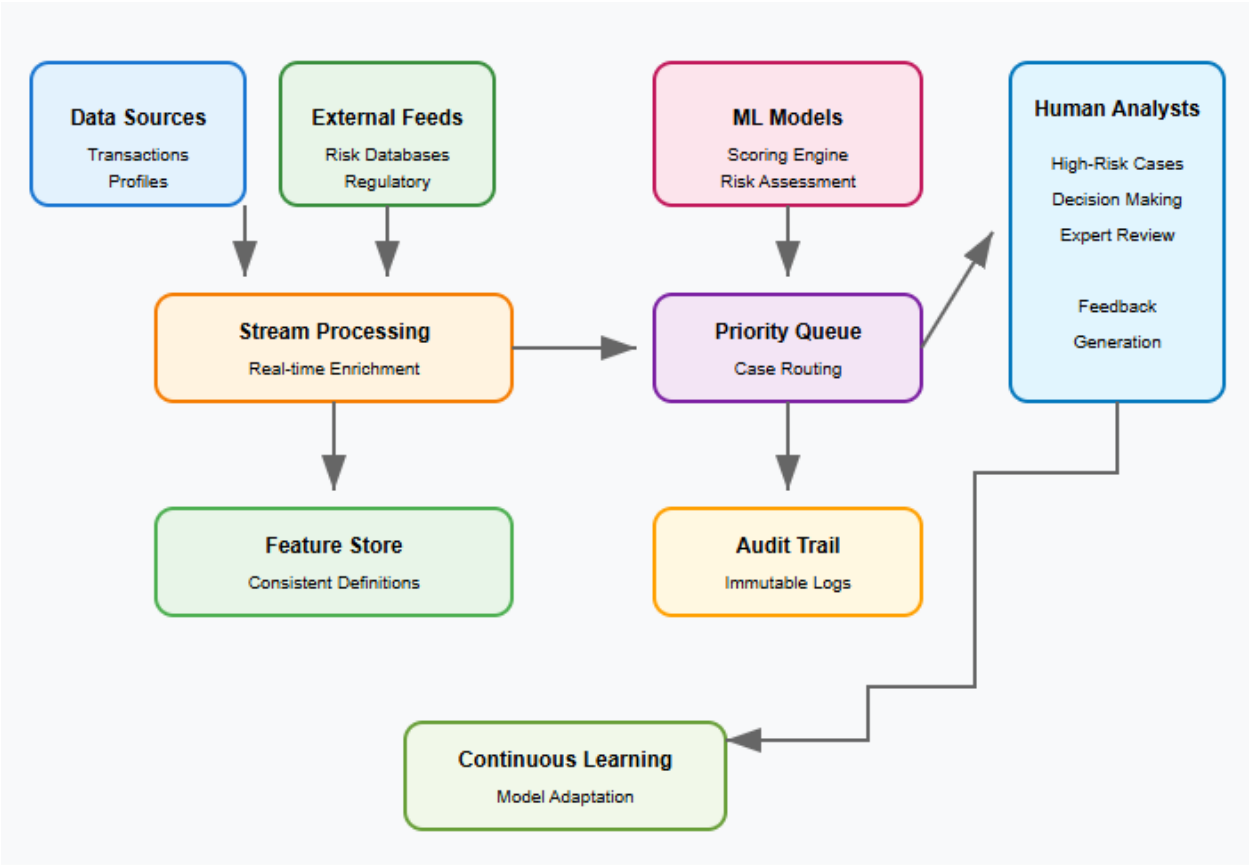


Fig 1. Illustrative HITL System Architecture Overview, Not Derived from any Production Systems  [3, 4].

**Data Traceability and Explainability Infrastructure**

**Audit Trail Architecture**

Financial institutions are required to keep detailed audit trails for compliance with regulations and internal governance, as the standards have grown more strict due to the expanding influence of artificial intelligence in financial decision-making processes.

The adoption of strong metadata management systems has become crucial for financial institutions aiming to uphold thorough data governance within intricate technology ecosystems that can encompass various cloud environments, onsite systems, and hybrid structures [5]. HITL systems necessitate improved traceability that records not just automated choices but also human interventions, justifications, and outcome adjustments, establishing intricate audit frameworks that must embrace the complexity of hybrid decision-making processes in which algorithmic suggestions are assessed, altered, or overridden by human experts operating within stringent regulatory deadlines.

The difficulty of sustaining detailed audit trails in HITL contexts goes beyond merely logging transactions to include the recording of decision context, analytical thought processes, and the timing relationships between data accessibility and decision-making moments. Contemporary metadata management tools offer functionalities for monitoring data lineage within intricate pipeline structures, allowing businesses to comprehend how data moves through different transformation phases and how modifications to upstream systems can affect downstream analytical operations and business choices [5]. The data infrastructure needs to monitor feature lineage, model versions, and decision paths to enable post-incident analysis and regulatory reviews, necessitating advanced metadata repositories that can uphold connections between data components, transformation logic, and business results over long durations while accommodating both technical and business user access behaviors.

Unchangeable event logs document each system interaction, generating lasting records of data changes, model conclusions, and human choices that underpin regulatory compliance and operational risk oversight. These detailed logging systems should record not only the eventual results of decision-making processes but also the intermediate conditions, various scenarios evaluated, and the reasoning for human actions that alter or replace algorithmic suggestions [5]. The logs facilitate the reconstruction of decision contexts long after the fact, aiding in compliance reporting and ongoing improvement efforts. However, the audit data produced by these systems can be considerable, frequently necessitating specialized storage designs and indexing methods to ensure satisfactory query performance for historical analysis and regulatory reporting demands, which may require prompt responses to examiner inquiries.

## Explainability Data Pipelines

Contemporary HITL systems utilize explainable AI methods that necessitate dedicated data processing frameworks designed to produce human-readable explanations for intricate machine learning model results in real-time operational settings. The area of interpretable machine learning has significantly progressed to meet the increasing demand for transparency in algorithmic decision-making, especially in regulated sectors where model outcomes need to be understandable to both internal stakeholders and external regulators [6]. SHAP values, LIME explanations, and feature importance scores need to be calculated in real-time and displayed alongside case details, posing technical challenges concerning computational efficiency and system scalability that require resolution through meticulous architecture planning and resource management strategies capable of accommodating diverse explanation complexity demands.

The challenge in data engineering consists of pre-calculating explanations for typical situations while preserving the ability to create personalized explanations for new cases, necessitating advanced prediction systems capable of forecasting explanation needs based on previous usage trends and the attributes of cases. The computational demands linked to producing model explanations can differ markedly based on the intricacy of the underlying models and the level of explanation needed, with certain explanation methods necessitating considerable processing resources that need to be managed carefully to prevent affecting core system performance [6]. This equilibrium becomes especially difficult during peak times when the system might have to produce explanations for numerous cases at once while ensuring suitable response times for analyst interfaces and preventing resource conflicts with primary transaction processing systems.

Explanation caching methods optimize computational efficiency while maintaining freshness needs, guaranteeing that analysts obtain prompt insights without overburdening system resources during high-demand processing times when the need for explanation generation may surge considerably above normal levels. Sophisticated interpretable machine learning methods offer various strategies for model explanation, including local explanations that concentrate on single predictions and global explanations that assist analysts in comprehending overall model behavior trends and possible biases [6]. The design should support various kinds of explanations, from basic feature important rankings to intricate counterfactual scenarios that aid analysts in comprehending how alterations to particular data components could impact model predictions, all while ensuring the computational efficiency required for immediate operational use in critical financial decision-making contexts.
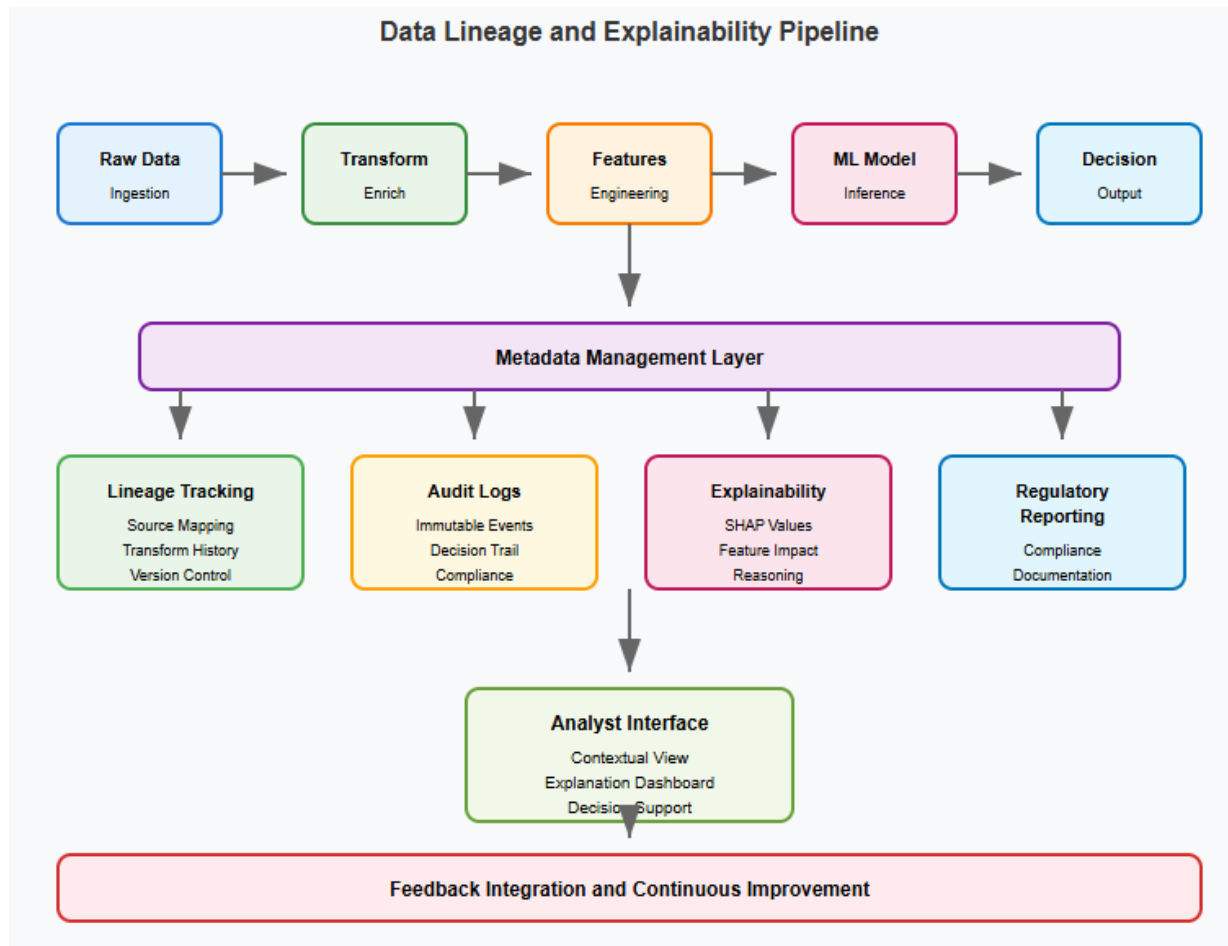
Fig 2.Illustrative  Data Traceability and Explainability Workflow

(not derived from any production systems)

## Feedback Integration and Continuous Learning

## Human Decision Capture

HITL systems produce important training data via human choices, but obtaining this feedback demands meticulous data engineering design that addresses the intricacies of human reasoning and the structured workflow needed to manage the machine learning lifecycle. The machine learning lifecycle includes several interrelated stages, such as data gathering, preprocessing, model creation, training, evaluation, and deployment, where each stage necessitates thorough coordination to guarantee that human feedback is efficiently incorporated into systematic enhancement processes [7]. User interfaces need to effectively document not only the final choices made but also the reasoning involved, alternative options evaluated, and confidence levels, generating thorough decision records that aid both immediate operational requirements and the continuous refinement processes central to successful machine learning lifecycle management.

The difficulty of obtaining valuable human feedback goes beyond merely recording decisions to include the intricate evaluation methods that skilled financial analysts use when analyzing complex situations with various risk elements and regulatory issues. Successful machine learning lifecycle management entails organized methods for data gathering and model assessment that can incorporate the subjective aspects of human evaluation while upholding the systematic discipline essential for algorithmic enhancement [7]. This enhanced feedback serves as training data for refining models and training analysts; however, incorporating human insights into model development processes necessitates meticulous consideration of data quality, consistency, and the timing alignment between feedback gathering and model retraining timelines.

Collecting structured feedback guarantees uniformity among various analysts and decision scenarios, necessitating the establishment of standardized procedures that conform to best practices in the machine learning lifecycle for data handling and model assessment. Template-driven reasoning and standardized confidence metrics facilitate the systematic examination of

human decision behaviors and highlight opportunities for enhancing outcomes through further training or tool improvements. However, these methods must be thoughtfully incorporated into the overall machine learning lifecycle to guarantee that feedback collection efforts aid rather than disrupt model development and deployment processes [7]. The feedback system must also support the iterative process of machine learning development, as initial model launches may need regular adjustments based on operational insights and evolving business needs.

## Model Adaptation Pipelines

Continuous learning systems integrate human feedback into model retraining processes via automated pipelines that identify feedback patterns, assess decision quality, and initiate model updates when performance criteria are achieved, utilizing advanced workflow management platforms capable of coordinating intricate machine learning tasks across distributed computing environments. Contemporary MLOps platforms offer extensive features for overseeing complete machine learning workflows, such as automated pipeline execution, resource allocation, and compatibility with different machine learning frameworks and deployment destinations [8]. These pipelines must reconcile model stability with the speed of adaptation, guaranteeing that valuable human insights enhance system performance without causing instability that might jeopardize operational effectiveness or adherence to regulatory compliance standards.

The adoption of continuous learning in financial services necessitates advanced pipeline architectures capable of handling the intricacies of model versioning, performance tracking, and automated deployment, all while addressing the regulatory and operational limitations inherent in high-stakes financial decision-making contexts. Workflow management platforms facilitate the development of consistent, scalable machine learning workflows capable of managing the computational requirements of ongoing model retraining, while ensuring the audit trails and governance measures necessary for financial services applications [8].

The pipeline structure should enable intricate validation procedures that assess model efficacy from various angles, while guaranteeing that model upgrades can be implemented securely without interfering with essential business functions. Automated model adaptation necessitates the incorporation of feedback processing features alongside advanced monitoring and validation systems that can evaluate when gathered human insights justify model updates, all while minimizing unnecessary retraining that may cause overfitting or system instability. The ongoing learning infrastructure must support the decentralized nature of contemporary machine learning practices, where model training, validation, and deployment can take place across various computing environments with differing resource features and accessibility [8]. These systems should also facilitate the collaborative elements of machine learning development, where data scientists, engineers, and business stakeholders must synchronize their efforts to ensure that ongoing learning processes are in line with business goals and regulatory standards while upholding the technical precision essential for dependable model performance in production settings.
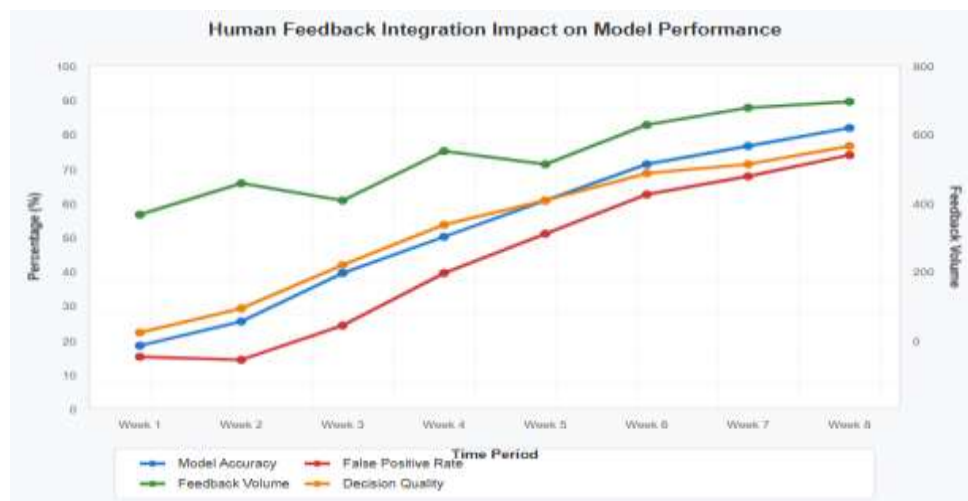


Fig 3. Illustrative Feedback Integration and Model Performance for HITL Systems

(not derived from any production systems)

**Scalability and Performance Optimization**

HITL systems need to manage unpredictable workload patterns since human decision delays can lead to processing backlogs that may severely affect operational efficiency and adherence to regulatory deadlines. The development of integrated analytics engines has transformed how financial institutions manage extensive data workloads, with contemporary distributed computing systems offering functionalities to process large datasets across streaming, batch, and interactive processing models within one cohesive framework [9]. Data engineering solutions consist of smart case routing, priority scoring, and workload distribution among analyst teams, necessitating advanced orchestration systems that utilize enhanced distributed computing abilities to adjust computational resources dynamically based on real-time workload traits and analyst availability trends.

When considering the diverse computational requirements of different analytical tasks, which range from simple rule-based filtering to complex machine learning model inference and graph analysis for network analysis, the challenge of managing unpredictable workloads in financial services environments becomes noticeably more complex. Contemporary unified analytics platforms offer comprehensive support for various processing approaches, allowing organizations to manage streaming transaction data, conduct batch historical analysis, and perform interactive analytical queries within a unified structure that can adjust to different computational requirements [9]. Caching techniques decrease response times for frequent queries, while pre-computed pipelines create uniform reports and visualizations that analysts often need; however, the success of these optimization methods relies heavily on the platform's capability to wisely manage memory use and enhance query execution across distributed computing clusters.

To implement efficient caching and pre-computation techniques, one must have an advanced grasp of data access trends and the skill to anticipate which analytical outputs will be needed for various investigative processes. Advanced analytics engines offer integrated optimization features such as adaptive query execution, dynamic partition pruning, and columnar storage formats, which can greatly enhance performance for the intricate analytical tasks typical of financial services applications [9]. Pre-computation pipelines need to balance the use of computational resources for speculative analysis creation with the operational advantages of quicker response times for routine analytical tasks, necessitating meticulous management of batch processing jobs that can generate analytical artifacts without disrupting real-time transaction processing needs.

Performance tracking emphasizes total decision latency, which encompasses both system processing time and the time taken for human review, necessitating extensive observability frameworks that deliver in-depth insights into application behavior within intricate distributed systems. Application performance monitoring has advanced to include advanced analytics features that can measure user experience metrics, resource usage, and the effectiveness of business processes via integrated platforms that offer real-time insights into application performance [10]. Data pipelines monitor queue lengths, analyst engagement, and decision-making quality indicators to pinpoint optimization chances and capacity planning needs; however, the intricacy of overseeing distributed analytical tasks necessitates sophisticated tools capable of relating performance metrics across various system elements and processing tiers.

The challenge of optimizing performance in HITL environments goes beyond conventional infrastructure monitoring, requiring an in-depth analysis of application performance trends that can identify optimization chances and possible bottlenecks before they affect business operations. Contemporary application performance monitoring tools offer features for monitoring tailored business metrics, examining user journey performance, and detecting performance irregularities that could suggest system deterioration or capacity issues [10]. These monitoring features allow organizations to deploy proactive performance management tactics that sustain optimal system functionality while addressing the fluctuating workload patterns and intricate analytical needs inherent in financial services HITL applications, aiding both operational efficiency and regulatory compliance goals through extensive insight into system behavior and performance attributes.
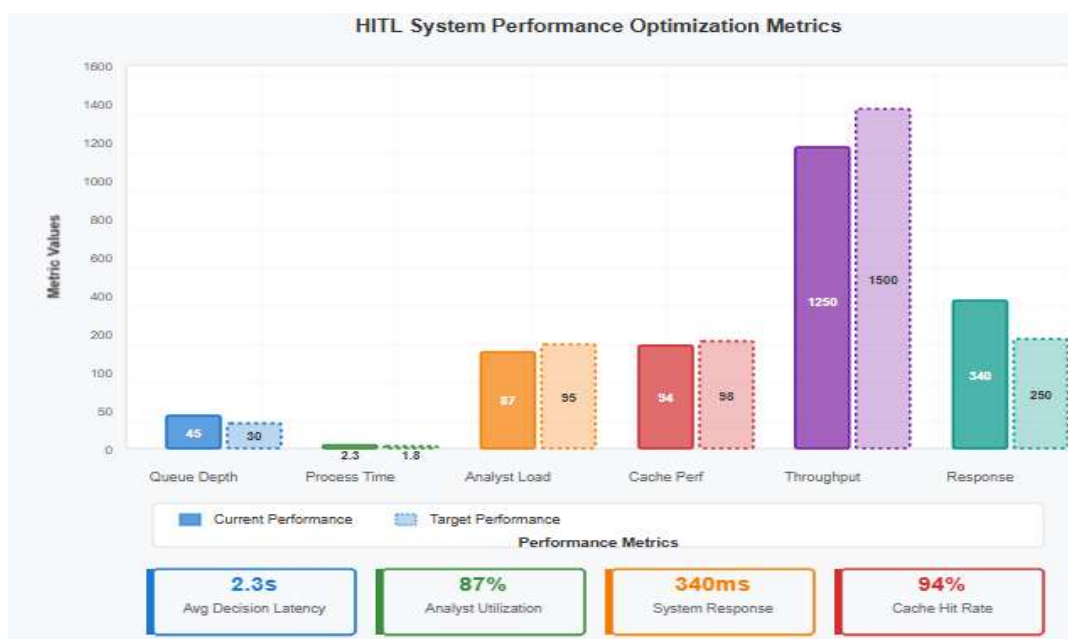
Fig 4. Illustrative System Performance and Scalability Metrics for HITL Systems

(not derived from any production systems)

## Conclusion

The deployment of human-in-the-loop AI systems in financial services signifies a refined equilibrium between tech automation and human skills, necessitating thorough data engineering frameworks that enable smooth cooperation between algorithmic functions and human judgment abilities. The effectiveness of these hybrid systems relies primarily on the foundational infrastructure that facilitates real-time data processing, ensures thorough audit trails, and backs explainable AI methods that offer clarity regarding algorithmic suggestions. The design complexity goes beyond conventional data processing frameworks to include tailored pipelines for feature engineering, which cater to both machine learning models and human analyst interfaces, resulting in dual-function systems that enhance computational efficiency while meeting human cognitive needs. Feedback integration methods facilitate ongoing enhancement by systematically collecting human choices and thought processes, generating useful training data that boosts model efficacy and aids analyst development programs. The scalability issues present in these systems demand advanced optimization techniques that can manage unpredictable workload trends and human decision delays by utilizing intelligent case routing, priority scoring, and dynamic resource distribution. With the ongoing evolution of regulatory requirements and advancements in artificial intelligence, the data engineering frameworks that facilitate human-AI collaboration will be essential for financial institutions aiming to sustain competitive advantage while fulfilling compliance needs and safeguarding customer interests through robust risk management and fraud prevention measures.

**Conflicts of Interest:** The authors declare no conflict of interest.
**Publisher's Note**: All claims expressed in this article are solely those of the authors and do not necessarily represent those of their affiliated organizations, or those of the publisher, the editors and the reviewers.

## References

[1] Nacha, "AFP Survey: Businesses May Be Misled on Risks of Check Fraud," 2025. [Online]. Available: https://www.nacha.org/news/afp-survey-businesses-may-be-misled-risks-check-fraud

[2] Lucinity, "The Future of AML Surveillance: The Role of AI in Anti-Money Laundering Compliance and Risk Mitigation," 2024. [Online]. Available: https://lucinity.com/blog/the-future-of-aml-surveillance-the-role-of-ai-in-anti-money-laundering-compliance-and-risk-mitigation

[3] Arun Anand, "Introducing the CData Arc Q1 2025 Version Release," CdataArc, 2025. [Online]. Available: https://arc.cdata.com/blog/cdata-arc-release-q1-2025

[4]    Melanie   Chen   and   Samuel   Mignot,   "An   Architectural   Deep   Dive,"   Chalk.AI,   2024.   [Online].   Available: https://chalk.ai/blog/what-is-a-feature-store

[5]    atlan,   "LinkedIn   DataHub   Guide:   Setup,   Features,   and   Alternatives   (2025),"   2024.   [Online].   Available: https://atlan.com/linkedin-datahub-metadata-management-open-source/

[6] Computer Languages, "Interpretable Machine Learning: A Guide For Making Black Box Models Explainable," 2025. [Online]. Available: https://www.clcoding.com/2025/04/interpretable-machine-learning-guide.html

[7]   Geeksforgeeks,   "Machine   Learning   Lifecycle,"   2025.   [Online].   Available:   https://www.geeksforgeeks.org/machine-learning/machine-learning-lifecycle/

[8] Soumyadarshan Dash, "Kubeflow: Streamlining MLOps With Efficient ML Workflow Management," AnalyticsVidhya, 2025. [Online].   Available:   https://www.analyticsvidhya.com/blog/2023/01/kubeflow-streamlining-mlops-with-efficient-ml-workflow-management/

[9] Kishan, "Spark 3.x Features," Medium, 2024. [Online]. Available: https://medium.com/@kishansingh2411/apache-spark-is-a-unified-analytics-engine-for-large-scale-data-processing-5240c5072614

[10] Kayly Lange, "What Is APM? Application Performance Monitoring, Explained," Splunk, 2024. [Online]. Available: https://www.splunk.com/en_us/blog/learn/apm-application-performance-monitoring.html