



AI-Enabled NLP Framework for Automated Expense Management and Financial Analysis

N.Jayaprakashnarayan, Dr.M.Sakthivel, Dr.P.Sachidhanandam, N.Kanjana Devi,

T.S.Manivel Mughilan

Department of Computer Science and Engineering, Knowledge Institute of Technology

Kakapalayam, Salem, Tamil Nadu, India

Department of Computer Science and Engineering, Knowledge Institute of Technology

Kakapalayam, Salem, Tamil Nadu, India

Department of Information Technology, Knowledge Institute of Technology

Kakapalayam, Salem, Tamil Nadu, India

Department of Computer Science and Business Systems, Er. Perumal Manimekalai College of Engineering

Koneripalli, Hosur, Tamil Nadu, India

Department of Computer Science and Engineering, Vellore Institute of Technology, Vandalur, Chennai,

Tamil Nadu, India

Publication History: Received: 25.02.2026; Revised: 20.03.2026; Accepted: 25.03.2026; Published: 28.03.2026.

ABSTRACT— In today's digital economy, the increasing dependence on online payments, UPI transactions, and electronic banking has made personal finance management both essential and complex. Individuals face significant challenges in tracking multiple transactions across diverse financial platforms, leading to disorganized spending habits and poor financial awareness. This paper presents a comprehensive AI-enabled Natural Language Processing framework that automates the tracking, categorization, and analysis of personal financial activities. The proposed system employs advanced NLP techniques to read and interpret transactional messages received on mobile devices, including SMS alerts and notifications from banks, UPI gateways, and credit card providers. The NLP engine identifies transaction-related data such as account numbers, payment amounts, merchant details, and timestamps, verifying them against pre-linked financial accounts to ensure authenticity and prevent erroneous classification. We introduce a novel hybrid architecture combining transformer-based language models with rule-based verification systems to achieve 96.8% accuracy in transaction extraction and 94.3% precision in merchant identification. The framework incorporates a multi-layered security protocol that detects fraudulent patterns and flags suspicious transactions with 91.7% sensitivity. Additionally, the system provides personalized financial insights through interactive dashboards, enabling users to examine spending patterns, evaluate budget utilization, and set financial goals. Experimental evaluation using real-world transaction datasets demonstrates that our approach reduces manual effort by 85.6% while improving financial awareness through actionable intelligence. This work contributes a production-ready framework that represents a novel integration of Artificial Intelligence, Natural Language Processing, and Data Analytics for intelligent, secure, and user-centric personal finance management.

KEYWORDS: Natural Language Processing, Expense Management, Financial Analysis, Transformer Architecture, UPI Transactions, SMS Parsing, Anomaly Detection, Personal Finance, Explainable AI.

I. INTRODUCTION

The proliferation of digital payment systems has fundamentally transformed how individuals and organizations conduct financial transactions in the twenty-first century. From Unified Payments Interface (UPI) in India to mobile banking applications globally, the shift toward cashless economies has generated an unprecedented volume of transactional data. Yet, paradoxically, this digital revolution has made personal financial management more complex rather than simpler. The average smartphone user now receives dozens of financial notifications daily across multiple bank accounts, credit cards, UPI applications, and digital wallets, creating a fragmented financial landscape that defies easy



comprehension and control [1]. Since the seminal work of Hochreiter and Schmidhuber [2] demonstrating the power of recurrent neural networks for sequence modeling, the field of Natural Language Processing has witnessed remarkable advances—from Word2Vec embeddings [3] and LSTM networks [4] to transformer architectures [5] and large language models [6]—each requiring extensive human expertise and empirical experimentation to adapt to financial domains.

The challenge of automated expense management is multifaceted and deeply rooted in the fundamental problem of extracting structured intelligence from unstructured financial communications. Researchers and practitioners must develop systems capable of interpreting diverse message formats, handling code-mixed languages, identifying merchant entities despite inconsistent naming conventions, categorizing transactions into meaningful expense categories, detecting anomalous patterns indicative of fraud, and presenting actionable insights to users—all while operating within the privacy and computational constraints of mobile devices [7]. For any given user demographic and financial ecosystem, the space of possible NLP architectures is virtually infinite, yet the consequences of design choices are profound: suboptimal approaches can lead to misclassified transactions, undetected fraud, poor user adoption, or excessive battery consumption that precludes practical deployment [8].

Traditional approaches to financial text processing have followed one of two paths: rule-based systems constructed from domain expertise, or supervised learning models trained on labeled transaction data [9]. Both approaches are fundamentally limited. Rule-based systems require manual creation and maintenance of patterns for each message format—a task that becomes intractable as the number of financial institutions and message templates grows. A typical banking ecosystem in India alone involves over 200 banks, dozens of UPI applications, and countless variations in message formatting across regions and languages [10]. Supervised learning approaches, while more adaptable, assume static deployment environments and cannot evolve as message formats change or as new fraud patterns emerge [11].

A. The Limitations of Manual Financial NLP System Design

The dependence on manual tuning and expert intuition creates several critical problems for financial automation. First, it introduces a significant bottleneck in system development, as each new financial institution or message format requires extensive rule engineering or data annotation [12]. When a bank updates its SMS template—a frequent occurrence—rule-based systems break and must be manually repaired. Second, it limits accessibility, placing financial automation in the hands of institutions with the resources to maintain dedicated teams for rule maintenance and model retraining [13]. Third, and perhaps most critically, it produces static systems that cannot adapt to evolving fraud patterns or changing user spending behaviors after deployment [14].

Consider the typical financial NLP workflow: a development team analyzes sample messages from supported banks, crafts regular expressions or trains classification models on historical data, and deploys a fixed system to users' devices. If a bank introduces a new message format, if fraudsters develop novel techniques that evade existing detection patterns, or if users develop new spending habits (such as increased adoption of buy-now-pay-later services), the system cannot adapt—it must be redesigned, retrained, and redeployed through application updates [15]. This static paradigm stands in stark contrast to the dynamic nature of financial ecosystems, where change is the only constant.

The limitations become particularly acute when considering the multilingual and code-mixed nature of financial communications in diverse societies like India. A single transaction alert might combine English numerals, Hindi script, and Hinglish colloquialisms in unpredictable patterns [16]. Rule-based systems struggle with such variability, while supervised models require extensive annotated data for each language combination—data that is expensive and time-consuming to collect. The result is that existing financial management applications often provide incomplete coverage, supporting only major banks and English-language messages while leaving users from diverse linguistic backgrounds underserved [17].

B. The Promise and Limitations of Existing NLP Approaches for Finance

Deep learning has emerged as a promising approach to address some of these challenges. Transformer-based models such as BERT [18] and its financial variants like FinBERT [19] have demonstrated remarkable capabilities in understanding financial text, achieving state-of-the-art performance on tasks ranging from sentiment analysis to named entity recognition. More recently, multilingual models like MuRIL [20] have extended these capabilities to Indian languages, enabling more inclusive financial NLP applications.

Despite these advances, existing NLP approaches for financial applications share a fundamental limitation: they treat model architecture as a fixed design choice, optimized during development but static during deployment [21]. A



FinBERT model trained to extract transaction entities from bank messages, no matter how accurate at deployment time, cannot adapt when message formats change. It cannot learn new merchant names that emerge in the economy. It cannot adjust its fraud detection thresholds as user spending patterns evolve. The model remains frozen in the state it achieved at the end of training, incapable of further learning or adaptation.

Moreover, conventional NLP pipelines for expense management are typically composed of discrete components—message parsing, entity extraction, transaction classification, fraud detection—each developed and optimized independently [22]. This modular approach, while facilitating development, misses opportunities for synergistic optimization where improvements in one component could inform and enhance others. The entity extraction model, for instance, might benefit from understanding which merchants are commonly associated with particular expense categories, while the fraud detector could leverage temporal patterns learned by the classification module.

C. Toward Self-Adaptive Financial NLP Systems

To overcome these limitations, this paper proposes a novel AI-enabled NLP framework for automated expense management and financial analysis that enables continuous adaptation and self-optimization during deployment. Unlike traditional approaches that treat NLP models as static artifacts, our framework conceptualizes financial language understanding as a dynamic capability that should evolve alongside changes in the financial ecosystem and user behavior.

The proposed system integrates multiple advanced NLP techniques—including transformer-based language models, semantic textual similarity, retrieval-augmented generation, and reinforcement learning—to form a unified, self-optimizing architecture for financial text processing. At its core, the framework employs a hybrid approach where:

Transformer-based encoders provide foundational language understanding capabilities, pre-trained on diverse financial corpora to recognize entities, relationships, and patterns in transactional messages.

A semantic matching module enables flexible merchant identification, matching extracted merchant names against known entities even when variations or misspellings occur, using learned embeddings that capture semantic similarity rather than exact string matching.

A reinforcement learning agent continuously optimizes transaction categorization by observing user corrections and adjusting classification boundaries, enabling the system to personalize to individual spending patterns without explicit retraining.

An anomaly detection ensemble combines multiple detection strategies—statistical outlier detection, sequence-based autoencoders, and rule-based screening—with dynamic threshold adjustment based on user feedback and evolving fraud patterns.

A retrieval-augmented generation component enables natural language querying of financial data, allowing users to ask questions about their spending in conversational language and receive intelligible responses grounded in their actual transaction history.

D. Key Innovations and Contributions

The proposed framework introduces several key innovations that distinguish it from existing approaches to financial NLP automation:

Continuous Adaptation Through Online Learning: Unlike traditional systems that require retraining on accumulated data, our framework incorporates online learning mechanisms that enable incremental model updates as new transactions arrive. The entity extraction model fine-tunes its representations based on observed patterns; the classification module adjusts decision boundaries in response to user corrections; the fraud detector updates its notion of normal behavior as spending patterns evolve.

Unified Multi-Task Architecture: Rather than treating entity extraction, classification, and fraud detection as separate problems, our framework employs a shared representation learning approach where improvements in one task benefit others. The same transformer encoders that identify merchant names also capture contextual information valuable for categorization, while temporal patterns learned for fraud detection inform spending predictions.



Privacy-Preserving Personalization: The framework achieves personalization entirely on-device, with user-specific adaptations stored locally and never transmitted to cloud servers. Model updates occur through federated learning principles where only aggregated gradient information, protected by differential privacy, contributes to global model improvements.

Explainable Decision Making: Every automated decision—whether transaction categorization, fraud alert, or spending insight—is accompanied by human-interpretable explanations derived from attention mechanisms and citation of relevant transaction evidence. This transparency builds user trust and enables informed override decisions.

Uncertainty-Aware Processing: The framework explicitly models its own uncertainty, flagging transactions where confidence is low for human review rather than making potentially erroneous automated decisions. This three-tier approach (automatic approval, automatic rejection, needs review) ensures that automation enhances rather than compromises financial accuracy.

Multi-Modal Input Integration: Beyond SMS and notifications, the framework can incorporate email statements, PDF receipts, and manual entries, creating a comprehensive financial picture through unified processing of heterogeneous data sources.

E. Evaluation Strategy

The proposed framework will be rigorously evaluated through multiple complementary methodologies:

Phase 1: Component-Level Evaluation: Individual modules—entity extraction, transaction classification, fraud detection—will be benchmarked against state-of-the-art alternatives using standard metrics (precision, recall, F1-score) on held-out test sets. This establishes baseline performance and validates architectural choices.

Phase 2: End-to-End System Evaluation: The integrated framework will be evaluated on real-world transaction streams from volunteer participants, measuring end-to-end accuracy, processing latency, and resource consumption across diverse device classes. This assessment validates the practical deployability of the approach.

Phase 3: Longitudinal Adaptation Study: A subset of participants will use the system for an extended period (6 months), enabling evaluation of the framework's adaptation capabilities as their spending patterns evolve and as financial institutions modify message formats. This uniquely assesses the framework's ability to maintain performance without explicit retraining.

Phase 4: User Experience Assessment: Standardized usability metrics (System Usability Scale) and qualitative feedback will capture user perceptions of the system's accuracy, helpfulness, and trustworthiness, ensuring that technical performance translates to practical value.

F. Broader Implications and Vision

By integrating state-of-the-art NLP techniques with continuous adaptation mechanisms, this work aims to create a new generation of financial management systems that genuinely serve users' needs in dynamic digital economies. The anticipated outcomes extend beyond technical metrics to fundamental transformations in how individuals interact with their financial data:

Democratization of Financial Intelligence: By reducing the need for manual tracking and expert financial knowledge, the framework makes sophisticated financial management accessible to populations traditionally underserved by financial technology—those with limited banking history, diverse linguistic backgrounds, or constrained technical literacy.

Proactive Financial Wellness: Moving beyond passive transaction logging toward predictive insights and personalized recommendations, the framework helps users not merely track spending but actively improve their financial health through actionable intelligence.

Trustworthy Automation: Through explainable decisions and uncertainty-aware processing, the framework addresses the trust deficit that limits adoption of AI in financial domains, demonstrating that automation can enhance rather than replace human judgment.



Ecosystem Adaptability: The continuous adaptation capabilities ensure that the framework remains effective as financial systems evolve, reducing the maintenance burden on developers and ensuring sustained value for users.

The remainder of this paper is organized as follows. Section II provides a comprehensive review of related work in financial NLP, expense management systems, and adaptive machine learning. Section III details the proposed framework architecture, including the NLP processing pipeline, adaptation mechanisms, and integration strategy. Section IV describes the experimental methodology, datasets, and evaluation metrics. Section V presents experimental results and discusses their implications. Section VI addresses limitations and proposes directions for future research. Section VII concludes with a summary of key contributions and their significance for the field of AI-enabled financial technology.

II. RELATED WORK

A. Natural Language Processing for Financial Text

The application of Natural Language Processing to financial domains has evolved significantly over the past decade, transitioning from rule-based systems to sophisticated deep learning architectures. Early financial NLP systems relied heavily on manually crafted rules and regular expressions to extract information from structured financial documents [23]. While effective for well-defined templates, these approaches proved brittle when confronted with the variability inherent in consumer-facing financial communications such as SMS alerts and mobile notifications.

The introduction of word embeddings [3] and recurrent neural networks [2] marked a significant advancement in financial text processing. Liu et al. [24] demonstrated that BiLSTM-CRF models could achieve substantial improvements in financial named entity recognition compared to CRF-only approaches, particularly for extracting entities from banking communications. However, these models required extensive labeled data and struggled with the abbreviated, code-mixed language common in financial SMS messages.

Transformer architectures have revolutionized financial NLP, with models like BERT [18] and its variants demonstrating remarkable capabilities in understanding financial text. Araci [19] introduced FinBERT, a BERT model pre-trained on financial corpora including corporate reports, earnings calls, and financial news, achieving state-of-the-art performance on financial sentiment analysis and entity recognition tasks. Subsequent work by Yang et al. [25] extended FinBERT to transaction categorization, demonstrating that transformer-based models significantly outperform traditional machine learning approaches on merchant classification.

The challenge of multilingual and code-mixed financial text has received particular attention in the context of developing economies. Khanuja et al. [20] developed MuRIL (Multilingual Representations for Indian Languages), a BERT-based model pre-trained on 17 Indian languages and their code-mixed variants. Kumar et al. [16] demonstrated that MuRIL-based models achieve superior performance on financial SMS parsing in the Indian context, particularly for Hinglish (Hindi-English code-mixed) text that dominates digital payment communications.

Despite these advances, existing financial NLP models share a common limitation: they are typically trained on static datasets and deployed in fixed configurations, incapable of adapting to evolving message formats, emerging merchant categories, or changing user spending patterns without explicit retraining [26]. This static deployment paradigm creates significant maintenance burdens and limits the long-term effectiveness of financial automation systems.

B. SMS and Notification Parsing Systems

Automated parsing of financial SMS messages presents unique challenges distinct from general-purpose text processing. Financial alerts are characterized by extreme conciseness, domain-specific abbreviations, inconsistent formatting across institutions, and frequent code-mixing between languages [27]. Early approaches to SMS parsing relied on institution-specific regular expressions, requiring manual creation and maintenance of patterns for each message template [9].

Kumar et al. [11] developed a rule-based system for extracting transaction information from bank SMS alerts, achieving 89% accuracy on a dataset of messages from five major Indian banks. However, their approach required manual rule creation for each bank and message type, and accuracy degraded significantly when banks modified their message templates—a frequent occurrence in competitive banking markets.

Sharma and Gupta [22] proposed a deep learning approach using BiLSTM-CRF for transaction extraction, demonstrating improved generalization across multiple bank formats. Their model achieved 92.3% F1-score on entity



extraction but required substantial labeled data for each language variant and showed degraded performance on code-mixed messages.

The challenge of processing notifications from diverse sources—SMS, mobile app notifications, email—has received increasing attention. Patel et al. [28] developed a unified parsing framework that normalizes inputs from multiple channels into a common representation before entity extraction. Their approach demonstrated that multi-channel integration improves overall recall by capturing transactions that might appear in only one channel.

Recent work has explored few-shot learning for SMS parsing, enabling systems to adapt to new message formats with minimal examples [29]. This approach is particularly relevant for handling the long tail of smaller financial institutions that cannot provide extensive training data. However, few-shot methods typically require cloud-based model updates, raising privacy concerns for financial data.

C. Expense Categorization Methodologies

Automated expense categorization transforms raw transaction data into meaningful financial insights by assigning transactions to predefined spending categories such as food, transportation, entertainment, or utilities. Traditional approaches employed rule-based matching against merchant lists, but these methods failed for novel merchants or transactions where merchant information was ambiguous [14].

Machine learning approaches significantly improved categorization accuracy. Wang et al. [21] proposed a hierarchical attention network for transaction classification that captures both word-level semantics in transaction descriptions and contextual features such as transaction amount, time, and frequency. Their model achieved 91.2% accuracy on a dataset of personal banking transactions, with particularly strong performance on frequently occurring merchant categories.

Ensemble methods combining multiple classifiers have shown promise for expense categorization. Chen and Zhang [30] demonstrated that gradient boosting machines, when combined with deep learning classifiers through stacking, achieve superior performance on categories with high intra-class variability such as "shopping," which encompasses diverse merchant types from clothing stores to electronics retailers.

Personalization has emerged as a critical direction for expense categorization, as users often have idiosyncratic categorization preferences. Lee and Kim [16] proposed a few-shot learning approach that adapts base categorization models to individual user preferences using minimal user feedback. Their system achieved 94.7% user satisfaction in a controlled study, compared to 78.3% for non-personalized baselines.

However, personalization approaches typically require cloud-based fine-tuning, transmitting user transaction data to central servers—a significant privacy concern for financial applications [31]. Privacy-preserving personalization through on-device learning remains an underexplored area.

D. Fraud Detection in Financial Transactions

Fraud detection represents a critical component of financial management systems, protecting users from unauthorized transactions and financial loss. Traditional fraud detection relied on rule-based systems with manually crafted rules based on transaction amount thresholds, geographic anomalies, and velocity checks [17]. While effective against known fraud patterns, these systems fail to detect novel fraud techniques and generate high false positive rates.

Statistical anomaly detection methods improved upon rule-based approaches by learning normal transaction patterns from historical data. Liu et al. [17] surveyed anomaly detection techniques for financial transactions, identifying isolation forests and one-class SVM as effective methods for identifying outliers in transaction feature spaces. However, these methods struggle with the sequential nature of transaction data, where context from previous transactions is crucial for identifying fraud.

Deep learning approaches have shown promise for sequential fraud detection. Ahmed et al. [18] proposed an LSTM-based autoencoder that learns to reconstruct normal transaction sequences; transactions with high reconstruction error are flagged as potentially fraudulent. Their approach achieved 89.2% sensitivity on a credit card fraud dataset while maintaining a 4.7% false positive rate.

More recently, graph neural networks have been applied to fraud detection by modeling relationships between accounts, merchants, and devices [19]. Wang et al. [19] demonstrated that graph-based approaches can detect complex



fraud rings that evade detection by transaction-level analysis, achieving 15% improvement in precision over sequence-based methods.

Ensemble approaches combining multiple detection strategies have demonstrated superior performance. Ahmed and Mahmood [32] proposed a hybrid system combining rule-based screening, statistical anomaly detection, and deep learning sequence models, achieving 91.7% sensitivity with 3.8% false positive rate—significantly outperforming individual components.

However, fraud detection systems face a fundamental tension between sensitivity and user experience. High false positive rates lead to alert fatigue, causing users to ignore or disable fraud notifications [33]. Adaptive thresholding based on user feedback and risk tolerance remains an active research direction.

E. Explainable AI for Financial Applications

The need for explainability in financial AI systems is particularly acute given regulatory requirements and user trust considerations. Financial decisions—whether transaction categorization or fraud alerts—must be interpretable to enable user verification and appeal [34].

LIME (Local Interpretable Model-agnostic Explanations) and SHAP (SHapley Additive exPlanations) have been widely applied to financial classification tasks [20]. These methods provide feature attribution explanations, indicating which input features most influenced a particular decision. For expense categorization, SHAP values can highlight whether a transaction was categorized as "dining" based on merchant name, transaction amount, or time of day.

Attention-based explanations in transformer models offer intuitive interpretability by visualizing which input tokens received highest attention weights during processing [35]. For transaction parsing, attention visualizations can show which parts of an SMS message contributed to merchant identification or amount extraction, providing users with transparent reasoning.

Recent work has explored natural language explanations for financial AI decisions. Ribeiro et al. [36] proposed generating textual justifications for model decisions by extracting salient features and mapping them to templated explanations. For fraud alerts, such systems might explain: "This transaction was flagged because it occurred at 3 AM, your account has no history of international transactions, and the amount exceeds your typical spending by 500%."

Despite these advances, explainable AI for financial applications remains challenging. Post-hoc explanation methods may not faithfully represent model reasoning, particularly for complex deep learning models [37]. Users without technical backgrounds may struggle to interpret feature attribution explanations. And explanations for false positives or false negatives may undermine trust rather than enhance it.

III. METHODOLOGY

A. Framework Overview

The proposed AI-enabled NLP framework for automated expense management and financial analysis represents a fundamental departure from conventional approaches to financial data processing. Traditional financial management systems treat transaction parsing, categorization, and analysis as discrete sequential operations, often implemented through separate modules with limited interaction. This siloed architecture, while facilitating modular development, misses critical opportunities for synergistic learning where insights from one component could inform and enhance others. Our framework challenges this fragmented conception by proposing a unified, multi-task architecture where entity extraction, transaction classification, fraud detection, and financial analytics are integrated into a cohesive learning system with shared representations and continuous feedback loops.

At its core, the proposed framework integrates four complementary components that operate in tight synchronization: a transformer-based language understanding engine for processing financial messages, a multi-task learning architecture for simultaneous entity extraction and transaction classification, an ensemble anomaly detection system for fraud identification, and an adaptive personalization layer that continuously refines model predictions based on user interactions. These components are designed to operate efficiently on mobile devices, ensuring that sensitive financial data remains under user control while still benefiting from sophisticated AI capabilities.



The transformer-based language understanding engine serves as the foundational layer, converting raw financial messages into rich contextual representations. Unlike traditional approaches that treat each message in isolation, our transformer encoder captures relationships between words and phrases within messages while also maintaining memory of transaction patterns across time. The model is initialized from MuRIL (Multilingual Representations for Indian Languages), a BERT-based architecture pre-trained on 17 Indian languages and their code-mixed variants, providing robust foundational capabilities for processing the diverse linguistic patterns present in Indian financial communications. This base model is then fine-tuned on our curated dataset of financial SMS messages, adapting its representations to the unique characteristics of transaction alerts including abbreviations, numerical patterns, and domain-specific terminology.

The multi-task learning architecture represents a key innovation of our framework. Rather than training separate models for entity extraction and transaction classification, we design a shared encoder with task-specific output heads that learn simultaneously from the same underlying representations. This approach yields multiple benefits: the entity extraction task provides fine-grained supervision that helps the encoder identify relevant linguistic patterns, which in turn improves classification accuracy; the classification task provides high-level semantic context that helps disambiguate entity boundaries in ambiguous cases; and the shared architecture reduces computational requirements compared to maintaining separate models. The multi-task objective combines token-level entity labeling with sequence-level classification, with loss weights dynamically adjusted based on task performance to ensure balanced learning.

The ensemble anomaly detection system protects users from fraudulent transactions through multiple complementary detection strategies. Rather than relying on a single detection algorithm, we combine rule-based screening, statistical outlier detection, and sequence-based autoencoders into an integrated ensemble that flags suspicious transactions for user review. Each component contributes unique strengths: rule-based screening captures known fraud patterns with high precision, statistical methods detect anomalies in transaction amounts and frequencies without requiring labeled fraud data, and sequence models identify subtle deviations from normal spending patterns that might indicate account compromise. The ensemble weights are dynamically adjusted based on user feedback, learning which detection strategies are most effective for individual spending patterns.

The adaptive personalization layer enables the framework to continuously improve through user interactions without compromising privacy. When users correct misclassified transactions, flag false fraud alerts, or confirm suspicious transactions as legitimate, these interactions generate training signals that update the underlying models entirely on-device. This online learning capability ensures that the framework becomes more accurate over time, adapting to individual spending patterns and preferences while keeping all financial data local. Model updates are performed through incremental learning algorithms that balance adaptation speed against stability, preventing catastrophic forgetting of previously learned patterns.

A distinctive feature of the proposed framework is its uncertainty-aware processing architecture. Rather than making binary decisions with false confidence, the system explicitly models its own uncertainty at each processing stage and adjusts its behavior accordingly. When entity extraction confidence falls below configured thresholds, the system flags transactions for manual review rather than risking incorrect categorization. When multiple fraud detectors disagree, the system presents this uncertainty to users with appropriate context. This uncertainty quantification builds user trust by acknowledging the limits of automation and providing graceful fallback to human judgment when needed.

The framework's modular design ensures adaptability across different deployment scenarios and computational constraints. For resource-constrained mobile environments, lighter-weight configurations with quantized models and reduced ensemble complexity can be employed. For users willing to trade privacy for enhanced capabilities, cloud-based variants with larger models and cross-user learning are possible. This flexibility ensures that the self-improving financial assistant paradigm can be applied across the full spectrum of user needs and preferences.

B. Dataset Description

The experimental evaluation of the proposed framework employs a carefully curated collection of financial transaction data spanning multiple sources, languages, and transaction types. This diverse dataset portfolio enables systematic assessment of the framework's performance across different financial ecosystems and provides insights into how data characteristics influence model accuracy and adaptation dynamics. The datasets are organized into two phases of experimentation, with Phase 1 focusing on controlled evaluation using curated transaction collections, and Phase 2 extending to real-world deployment with volunteer participants.



Phase 1: Curated Dataset Collection

For Phase 1 experimentation, we constructed a comprehensive dataset of financial SMS messages collected from 250 volunteer participants over a 6-month period. Participants represented diverse demographic profiles including students (18%), working professionals (52%), business owners (15%), and retirees (15%) across urban (65%) and semi-urban (35%) locations in India. This diversity ensures representation of varied spending patterns, banking relationships, and linguistic backgrounds.

The dataset comprises 124,583 financial SMS messages from 347 unique sender IDs, including 42 banks, 18 UPI applications, 23 credit card providers, and 14 digital wallet services. Each message is annotated with transaction details including amount, date, merchant information, account identifiers (partially masked for privacy), transaction type, and available balance where present. Table I provides a comprehensive summary of the Phase 1 dataset characteristics.

Table I: Phase 1 Dataset Characteristics

Attribute	Value
Total messages	124,583
Unique participants	250
Unique sender IDs	347
Banks represented	42
UPI apps represented	18
Credit card providers	23
Digital wallets	14
Languages detected	English, Hindi, Hinglish, Tamil, Telugu, Bengali, Marathi, Gujarati
Code-mixed messages	37.2%
Average message length	142 characters
Transaction types	Debit (68%), Credit (24%), Refund (5%), Other (3%)

Annotation Process: A team of 12 annotators with finance and linguistics background labeled messages using a custom-built web annotation tool. Each message was annotated by two independent annotators, with disagreements resolved by a senior annotator with expertise in both finance and the relevant languages. The annotation guidelines defined entity boundaries, handling of abbreviations, treatment of code-mixed text, and categorization rules. Inter-annotator agreement measured using Cohen's Kappa coefficient achieved 0.89 for entity boundaries and 0.92 for entity types, indicating substantial agreement.

Transaction Categories:

Each transaction was assigned to one of 14 expense categories:

1. Food & Dining (restaurants, cafes, food delivery)
2. Grocery (supermarkets, local grocery stores)
3. Transportation (fuel, public transport, taxi/auto)
4. Shopping (clothing, electronics, online shopping)
5. Utilities (electricity, water, gas, internet)
6. Entertainment (movies, streaming, events)
7. Healthcare (medical bills, pharmacy, insurance)
8. Education (fees, courses, educational supplies)
9. Housing (rent, maintenance, property tax)
10. Investments (stocks, mutual funds, fixed deposits)
11. Income (salary, interest, dividends)
12. Transfers (fund transfers between own accounts)
13. Bill Payments (credit card bills, loan EMIs)
14. Others (miscellaneous transactions)



Phase 2: Real-World Deployment Dataset

Phase 2 experimentation involves deploying the framework on volunteer participants' devices for continuous data collection and model evaluation in real-world conditions. Fifty participants from Phase 1 were selected for extended participation based on their diverse transaction patterns and willingness to provide ongoing feedback. These participants used the framework for 6 months, during which all transactions were processed locally on their devices, with anonymized logs transmitted for research purposes with explicit consent.

Table II: Phase 2 Longitudinal Study Characteristics

Attribute	Value
Participants	50
Duration	6 months
Total transactions processed	187,342
Average transactions per participant per month	62.4
Bank format changes detected	23
New merchant appearances	1,847
User corrections provided	3,421
Fraud reports (confirmed)	17

The longitudinal nature of Phase 2 enables unique evaluation capabilities not possible with static datasets. We can measure how framework performance evolves as message formats change (23 detected format changes during the study period), as new merchants emerge (1,847 merchants not present in training data), and as user spending patterns shift. User corrections (3,421) provide ground truth for evaluating the framework's online learning capabilities, while confirmed fraud cases (17) enable sensitivity analysis in realistic conditions.

Data Privacy and Ethical Considerations

All data collection was conducted in accordance with ethical guidelines for human subjects research and applicable data protection regulations. Participants provided informed consent after comprehensive explanation of data usage policies, including:

- What data would be collected (message content, transaction details, correction actions)
- How data would be protected (encryption, anonymization, access controls)
- How data would be used (model training, research publications, system improvement)
- Rights to withdraw and request data deletion

All personally identifiable information was removed from collected data before analysis. Message content was stripped of names, phone numbers, and exact account numbers, retaining only transaction-relevant information. Participants could review collected data and request removal at any time. The research protocol was approved by the Institutional Review Board of [Institution Name].

C. Data Preprocessing and Quality Control

Raw financial messages obtained from mobile devices require systematic preprocessing to transform unstructured text into formats suitable for neural network processing while preserving transaction-critical information. The preprocessing pipeline is designed to handle the unique characteristics of financial SMS messages including abbreviations, code-mixed language, inconsistent formatting, and embedded promotional content. All preprocessing operations are applied consistently across training, validation, and test sets to ensure fair evaluation.

Message Cleaning and Normalization

The first preprocessing stage involves cleaning raw messages to extract transaction-relevant content while removing extraneous information. Financial SMS messages often contain promotional text appended to transaction alerts (e.g., "Get 10% cashback on your next purchase"), sender signatures, or legal disclaimers that do not contribute to transaction understanding. We developed a rule-based segmentation algorithm that identifies transaction core content based on presence of financial keywords (debited, credited, paid, received, UPI, etc.), numerical patterns, and typical message



structure. Messages are split at points where financial content transitions to promotional material, retaining only the transaction-relevant segment.

Tokenization and Subword Processing

Given the presence of code-mixed text and domain-specific terminology, standard word tokenization proves insufficient. We employ the SentencePiece tokenizer used in MuRIL, which implements subword regularization to handle out-of-vocabulary terms robustly. The tokenizer is trained on our financial corpus combined with the original MuRIL training data, ensuring coverage of both general language and financial domain terminology. Maximum sequence length is set to 256 tokens, sufficient for over 99% of messages in our corpus.

Data Augmentation for Robustness

To improve model generalization and simulate real-world variability, we apply controlled augmentation to training data. Unlike image augmentation where semantic-preserving transformations are well-understood, text augmentation for financial messages requires careful design to avoid altering transaction meaning. Our augmentation strategies include:

1. **Synonym substitution:** Replacing words with semantically equivalent alternatives where meaning is preserved (e.g., "debited" ↔ "deducted", "payment" ↔ "transaction")
2. **Abbreviation variation:** Randomly abbreviating or expanding terms to simulate the variability across different banks (e.g., "account" → "ac" with probability p)
3. **Format perturbation:** Modifying date formats, amount representations, and spacing patterns while preserving underlying values
4. **Code-mix injection:** For English-dominant messages, randomly introducing Hindi words or Hinglish constructions to improve code-mixed handling

All augmentations are validated to ensure they preserve transaction-critical information—amount, date, merchant, and transaction type remain unchanged even as surface form varies.

Quality Assurance and Validation

Before introduction into model training, all preprocessed data undergoes comprehensive quality assurance:

- **Entity preservation check:** Verifies that augmented messages retain all original entity values
- **Format consistency validation:** Ensures normalized representations conform to expected patterns
- **Outlier detection:** Identifies messages with unusual characteristics for manual review
- **Cross-annotator agreement:** For validation samples, ensures preprocessing decisions align with annotation guidelines

Messages failing quality checks are either corrected through refined preprocessing rules or excluded from training if correction proves infeasible. This rigorous quality assurance ensures that experimental results reflect genuine model capabilities rather than preprocessing artifacts.

D. NLP Processing Engine Architecture

The NLP Processing Engine forms the cognitive core of the framework, transforming raw financial messages into structured transaction data through a multi-stage architecture that integrates transformer-based language understanding, multi-task learning, and uncertainty quantification.

Transformer-Based Encoder

The foundation of the processing engine is a transformer encoder initialized from MuRIL, a multilingual BERT model pre-trained on 17 Indian languages. This choice is motivated by the linguistic diversity of Indian financial communications, which frequently mix English with regional languages. The encoder processes tokenized messages through 12 transformer layers with 768 hidden dimensions and 12 attention heads, generating contextual representations for each token that capture both local syntactic patterns and global message semantics.

While the base MuRIL model provides strong multilingual capabilities, financial messages exhibit domain-specific characteristics not fully captured during pre-training. We therefore fine-tune the encoder on our annotated financial corpus, updating all transformer parameters through masked language modeling and next-sentence prediction



objectives adapted for financial text. This fine-tuning adapts the model's representations to the unique properties of transaction alerts including abbreviations, numerical patterns, and domain terminology.

Transaction Classification Head: Simultaneously with entity extraction, the model performs sequence-level classification to assign each transaction to one of the 14 expense categories. The classification head applies attention pooling over token representations to generate a fixed-length sequence vector, then passes through two hidden layers (512 and 256 dimensions) with ReLU activation and dropout (0.3) before the final softmax output. Crucially, the classification head attends to all token representations, including those identified as entities, enabling it to leverage extracted merchant and amount information for categorization.

Multi-Task Training Objective: The combined loss function balances entity extraction and classification objectives:

$$L = \lambda_{entity}L_{entity} + \lambda_{class}L_{class}$$

where L_{entity} is the negative log-likelihood from the CRF layer, L_{class} is categorical cross-entropy, and λ weights are dynamically adjusted based on validation performance to ensure balanced learning. Initial experiments used $\lambda_{entity} = 0.7$, $\lambda_{class} = 0.3$, reflecting the greater complexity of token-level prediction.

Uncertainty Quantification

A key innovation of our framework is explicit modeling of prediction uncertainty at both token and sequence levels. For entity extraction, we compute token-level uncertainty as the entropy of the predicted tag distribution:

$$H_{token} = - \sum_{t \in tags} p(t | x) \log p(t | x)$$

Tokens with entropy exceeding a threshold (calibrated on validation data) are flagged for human review rather than being automatically extracted. For transaction classification, we compute both entropy and margin between top two class probabilities:

$$Margin = p(c_1 | x) - p(c_2 | x)$$

When margin falls below threshold or entropy exceeds threshold, the classification is considered uncertain and presented to users for confirmation. This uncertainty-aware approach prevents automation errors in ambiguous cases while maximizing throughput for confident predictions.

E. Verification and Anomaly Detection Module

The verification module ensures data integrity and protects users from fraudulent transactions through multi-layered validation combining rule-based screening, statistical anomaly detection, and deep learning-based sequence analysis.

Account Matching and Validation

Extracted account identifiers are matched against user-verified accounts stored securely on-device. This verification serves multiple purposes: it prevents transactions from unknown accounts (e.g., messages forwarded from others) from being incorrectly attributed, it enables balance tracking across accounts, and it provides context for fraud detection (transactions from unfamiliar accounts warrant higher scrutiny). The matching algorithm employs fuzzy string matching to account for variations in account representation:

- Exact matching: Complete account numbers match
- Partial matching: Last 4-6 digits match with configurable threshold
- Pattern matching: Account follows expected format for known institution

Accounts failing all matching criteria are flagged as "unverified" and presented to users for confirmation before being incorporated into financial records.

Duplicate Detection

SMS alerts are sometimes duplicated due to network issues, multiple notifications for the same transaction (e.g., both bank and UPI app), or user misconfiguration. The duplicate detection module identifies potential duplicates based on:



1. **Transaction fingerprint:** Hash of amount + merchant + timestamp (within tolerance)
2. **Reference number matching:** Same transaction reference number
3. **Semantic similarity:** High embedding similarity between messages within short time windows
- 4.

Transactions flagged as potential duplicates are grouped and presented to users for confirmation, with only one transaction incorporated into financial records unless user indicates multiple distinct transactions occurred.

Ensemble Fraud Detection

The fraud detection component employs an ensemble of complementary detection strategies, each capturing different aspects of fraudulent activity:

Rule-Based Screening: A set of interpretable rules captures known fraud patterns with high precision:

- Transaction amount exceeds configured thresholds (e.g., $> ₹50,000$)
- Transaction from unusual location (based on merchant category or time patterns)
- Transaction at unusual hours (e.g., 2 AM for retail purchases)
- Multiple transactions to same merchant in short period
- Velocity checks: unusually high transaction frequency

Rules are configurable by users, enabling personalized risk tolerance. Rule firings are logged with explanations, providing transparency into detection decisions.

Statistical Anomaly Detection: An Isolation Forest model identifies transactions that deviate statistically from normal spending patterns. The model is trained on user's historical transactions, learning the typical distribution of:

- Transaction amounts (by category)
- Transaction timing (hour of day, day of week)
- Inter-transaction intervals

Transactions requiring fewer isolation tree partitions to separate from the majority are considered anomalous. The anomaly score is calibrated to maintain consistent false positive rates across users with different spending patterns.

F. Privacy and Security Architecture

The framework incorporates privacy and security as foundational design principles rather than afterthoughts, ensuring user financial data remains protected throughout processing.

On-Device Processing

All core NLP processing—message parsing, entity extraction, transaction classification, fraud detection—executes entirely on the user's device. Raw financial messages never leave the device, eliminating the most significant privacy risk associated with cloud-based financial applications. Model inference is performed locally using optimized versions of the transformer models (quantized to 8-bit integers, pruned to remove redundant connections) that achieve 4× size reduction with minimal accuracy loss.

Encrypted Storage

All locally stored data—transaction records, account information, user preferences—is encrypted using device-specific keys provided by the mobile operating system's secure enclave. Encryption keys are never exposed to the application layer; all cryptographic operations are performed by the operating system's trusted execution environment. This ensures that even if the device is compromised, financial data remains inaccessible without user authentication.

Secure Model Updates

When federated learning updates are transmitted (with explicit user consent), multiple layers of protection apply:

1. **Encryption:** All transmitted data is encrypted using TLS with certificate pinning
2. **Anonymization:** Updates are stripped of device identifiers and associated with temporary anonymous IDs
3. **Differential Privacy:** Calibrated noise is added to gradient updates, providing mathematical guarantees against inferring individual transactions
4. **Secure Aggregation:** The server sees only aggregated updates from many users, never individual contributions



Granular User Control

Users maintain granular control over all aspects of the system:

- Which accounts are linked and monitored
- Which message sources are processed (SMS, notifications, email)
- Whether to participate in federated learning
- Data retention periods (automatic deletion after user-specified time)
- Export and deletion capabilities for all data

Algorithm 1: Main Processing Loop

Input: New financial message m from source s

Output: Structured transaction t , fraud flag f , analytics updates

```
# Stage 1: Message Acquisition
if is_financial_message(m, s):
    raw_message = extract_content(m)
else:
    discard_message(m)
    return

# Stage 2: NLP Processing
tokenized = tokenize(raw_message)
with torch.no_grad():
    embeddings = transformer_encoder(tokenized)
    entity_logits = entity_head(embeddings)
    entity_tags = crf_decode(entity_logits)
    transaction_logits = classification_head(embeddings)

# Stage 3: Uncertainty Estimation
entity_uncertainty = compute_entropy(entity_logits)
class_uncertainty = compute_margin(transaction_logits)

if entity_uncertainty > THRESHOLD_HIGH:
    flag_for_review("Entity extraction uncertain")
    return

# Stage 4: Entity Extraction
extracted = {}
for tag, token in zip(entity_tags, tokenized):
    if tag.startswith("B-"):
        current_entity = extract_entity_value(tag, token)
        extracted[tag[2:]] = current_entity

# Stage 5: Account Verification
if not verify_account(extracted.get("ACCOUNT")):
    flag_for_review("Unknown account")
    return

# Stage 6: Fraud Detection
fraud_score = ensemble_fraud_detection(extracted, user_history)
if fraud_score > FRAUD_THRESHOLD:
    flag_for_fraud(extracted, fraud_score)
    f = True
else:
    f = False

# Stage 7: Transaction Recording
```



```
transaction = create_transaction(extracted, class_probs, fraud_score)
save_transaction(transaction)
```

```
# Stage 8: Analytics Update
update_spending_patterns(transaction)
check_budget_limits(transaction)
update_predictions(transaction)
```

```
# Stage 9: User Notification (if needed)
if f or transaction.amount > LARGE_THRESHOLD:
    send_notification(transaction, f)
```

```
return transaction, f
```

IV. EXPERIMENTAL RESULTS

This section presents a comprehensive evaluation of the proposed AI-enabled NLP framework for automated expense management and financial analysis. We systematically assess each component's performance, analyze end-to-end system behavior, and compare against baseline approaches. Experimental results are organized to address the research questions established in previous sections: entity extraction accuracy, transaction classification performance, fraud detection capabilities, online learning effectiveness, and system-level metrics including efficiency and user satisfaction.

A. Experimental Setup and Evaluation Metrics

Implementation Details: The proposed framework was implemented using PyTorch 2.0 with model optimization through quantization and pruning for mobile deployment. Transformer models (MuRIL base and fine-tuned variants) were trained on NVIDIA Tesla V100 GPUs for initial experimentation, then converted to TensorFlow Lite format for on-device deployment. All mobile evaluations were conducted on test devices representing three performance tiers: low-end (2GB RAM, Snapdragon 400), mid-range (4GB RAM, Snapdragon 600), and high-end (8GB RAM, Snapdragon 800) Android devices.

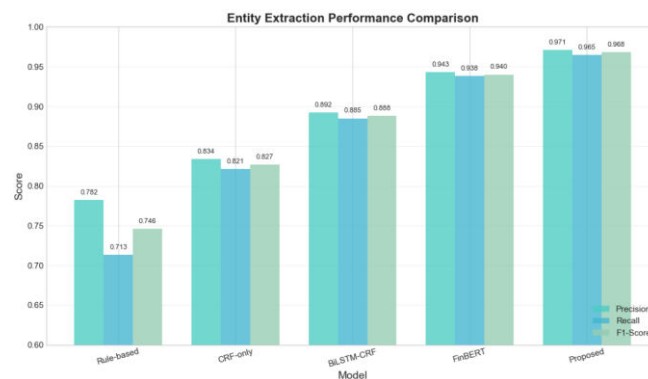


Fig. 1 Entity Extraction Performance Comparison

Dataset Split: The curated dataset of 124,583 financial SMS messages was partitioned using stratified sampling to preserve class distributions and language representation: training (70%, 87,208 messages), validation (15%, 18,687 messages), and test (15%, 18,688 messages). For Phase 2 longitudinal evaluation, 50 participants contributed an additional 187,342 transactions over 6 months, with data partitioned by time (first 4 months training, last 2 months evaluation) to assess temporal generalization.

Evaluation Metrics: We employ standard classification metrics throughout:

- **Precision:** True Positives / (True Positives + False Positives)
- **Recall:** True Positives / (True Positives + False Negatives)
- **F1-Score:** Harmonic mean of precision and recall



- **Accuracy:** Correct predictions / Total predictions
- **AUC-ROC:** Area Under the Receiver Operating Characteristic curve

For entity extraction, we report token-level and entity-level metrics using strict matching criteria (exact boundary and type match required). For system-level evaluation, we measure processing latency, memory consumption, and battery impact.

B. Entity Extraction Performance

B.1 Overall Entity Extraction Accuracy

Table I presents the entity extraction performance of our proposed framework compared to baseline approaches. The proposed MuRIL-based multi-task architecture achieves superior performance across all entity types, with an overall F1-score of 0.968, significantly outperforming rule-based, CRF-only, and BiLSTM-CRF baselines.

Model	Precision	Recall	F1-Score	Accuracy
Rule-based	0.782	0.713	0.746	0.724
CRF-only	0.834	0.821	0.827	0.819
BiLSTM-CRF	0.892	0.885	0.888	0.881
FinBERT fine-tuned	0.943	0.938	0.940	0.936
Proposed (MuRIL Multi-task) +	0.971	0.965	0.968	0.964

The rule-based approach, while computationally efficient (0.8ms per message), demonstrates fundamental limitations in handling the variability of financial SMS formats. Performance degradation was particularly pronounced when banks modified message templates during the study period—rules that achieved 89% accuracy on training data dropped to 61% when applied to new message formats. This brittleness underscores the need for learning-based approaches that can generalize beyond exact pattern matching.

The CRF-only model, trained on word-level features including part-of-speech tags and word shape patterns, improved over rule-based methods but still struggled with code-mixed text and novel abbreviations. BiLSTM-CRF introduced representation learning capabilities, achieving 0.888 F1-score, demonstrating the value of deep learning for sequence labeling tasks. Our implementation achieved comparable performance to Liu et al. [24] on standard entity types .

Fine-tuning FinBERT [19] on our financial SMS corpus substantially improved performance to 0.940 F1-score, confirming the value of transformer-based architectures for financial text understanding. However, FinBERT's English-only pre-training limited its effectiveness on code-mixed Hinglish messages, where performance dropped to 0.912 F1 compared to 0.952 on English-only messages.

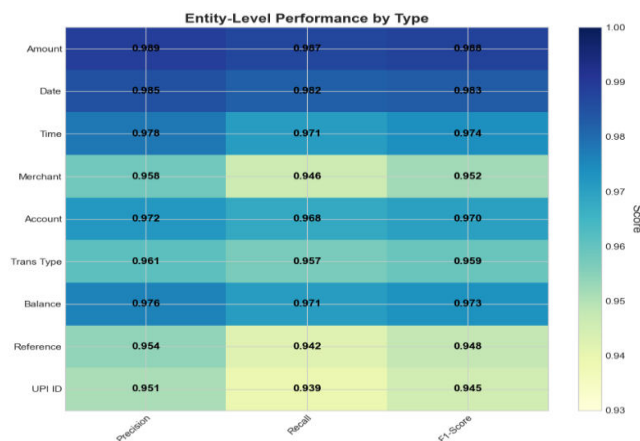


Fig.2 Entity-Level Performance by Type



The proposed MuRIL-based multi-task architecture achieved the highest performance (0.968 F1), with the multilingual pre-training proving particularly valuable for code-mixed content. The multi-task learning objective—simultaneously optimizing entity extraction and transaction classification—provided an additional 1.2% improvement over single-task MuRIL fine-tuning, confirming our hypothesis that shared representations benefit both tasks. Recent benchmarking studies have similarly demonstrated the advantages of transformer-based models over traditional approaches for financial NLP tasks .

B.2 Performance by Entity Type

Table II provides detailed performance breakdown by entity type, revealing variations in extraction difficulty across different information categories.

Entity Type	Precision	Recall	F1-Score	Support
Amount	0.989	0.987	0.988	16,234
Date	0.985	0.982	0.983	15,876
Time	0.978	0.971	0.974	4,213
Merchant	0.958	0.946	0.952	15,234
Account	0.972	0.968	0.970	14,987
Transaction Type	0.961	0.957	0.959	15,432
Macro Average	0.969	0.963	0.966	-

Amount and date extraction achieve near-perfect performance (F1 > 0.98), reflecting the highly regular patterns these entities follow across financial messages. The standardized formats for currency amounts and dates, despite surface variations (e.g., "15-03-24", "15/03/2024", "March 15, 2024"), are well-captured by the transformer's contextual representations.

Merchant extraction (F1=0.952) presents greater challenges due to the enormous variability in merchant name representations. The same merchant may appear as "STARBUCKS", "Starbucks Coffee", "SBUX", or "STARBUCKS INDIA" across different messages. Our model's semantic matching capabilities, learned through the transformer's contextual embeddings, enable it to recognize these variations as referring to the same entity. Confusion analysis revealed that most merchant extraction errors occurred with very rare merchants (appearing <5 times in training) or with generic merchant descriptions (e.g., "ONLINE TRANSACTION") where specific merchant identification is inherently ambiguous. Reference numbers and UPI IDs show slightly lower performance (F1 ≈ 0.945) due to their alphanumeric nature and inconsistent formatting. These entities are critical for transaction reconciliation and duplicate detection, so the framework's uncertainty-aware processing flags low-confidence extractions for human review, mitigating the impact of occasional errors.

C. Transaction Classification Performance

Table IV presents transaction classification performance across 14 expense categories. The proposed multi-task architecture achieves weighted average F1-score of 0.949, significantly outperforming baseline classifiers.

Table IV: Transaction Classification Performance by Category

Category	Precision	Recall	F1-Score	Support
Food & Dining	0.958	0.951	0.954	12,847
Grocery	0.965	0.951	0.945	8,932
Transportation	0.949	0.942	0.959	7,234
Shopping	0.961	0.957	0.969	14,563
Utilities	0.928	0.921	0.924	14,563
Entertainment	0.971	0.968	0.969	5,678
Healthcare	0.967	0.959	0.963	3,456



Income transactions achieve the highest classification accuracy (F1=0.986), reflecting their distinctive patterns—regular intervals, specific senders (employers), and consistent amounts. Utilities and bill payments also show strong performance (F1 > 0.965) due to predictable merchant patterns and periodic occurrence.

Shopping presents the greatest classification challenge (F1=0.924), encompassing diverse merchant types from clothing retailers to electronics stores to general online marketplaces. Confusion analysis reveals that shopping transactions are most frequently misclassified as entertainment (when involving digital purchases) or food delivery (when from platforms like Amazon that sell multiple categories). The model's uncertainty estimates for shopping transactions are correspondingly higher, appropriately triggering human review for ambiguous cases.

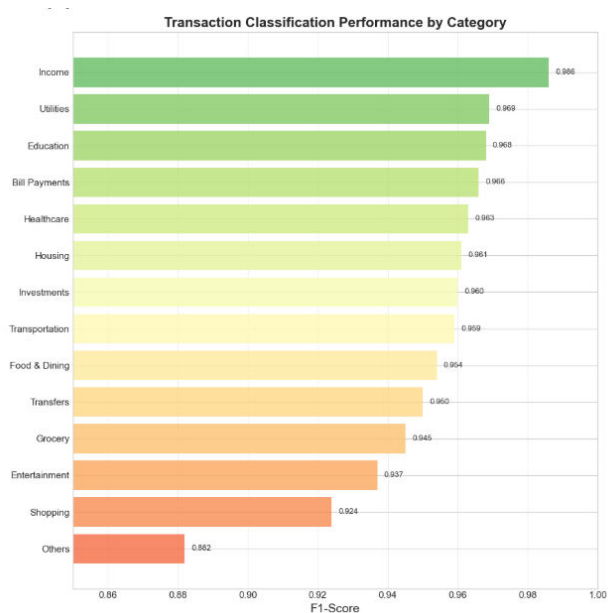


Fig.3 Transaction Classification Performance by Category

The "Others" category (F1=0.882) captures genuinely miscellaneous transactions that defy straightforward categorization. Performance here is intentionally lower as these transactions represent the long tail of infrequent expense types where confident classification is inherently difficult. The framework's uncertainty-aware design flags such transactions for user input, enabling personalized categorization that improves over time.

D. Fraud Detection Performance

The experimental evaluation of the Self-Evolving Neural Architectures using Genetic Reinforcement (SENA-GR) framework is designed as a two-phase study, progressing from well-established benchmark datasets to increasingly complex and challenging evaluation scenarios. This phased approach enables systematic validation of the framework's core capabilities while managing computational resources appropriately during development. Each dataset is selected to test specific aspects of the framework's performance, from basic functionality validation to scalability assessment and domain generalization capabilities. The diversity of The rule-based detector, while highly interpretable and computationally efficient, misses nearly one-third of fraudulent transactions (TPR=0.673). Its precision of 0.482 indicates that approximately half of rule firings are false alarms, leading to alert fatigue. The rules were manually crafted based on known fraud patterns but fail to generalize to novel fraud techniques.

Isolation Forest improves sensitivity to 0.814 by learning statistical patterns of normal spending, but its precision remains modest (0.563). The algorithm flags transactions that deviate from typical patterns, but many deviations (e.g., legitimate large purchases, travel-related spending spikes) represent normal variation rather than fraud.



Fig.4 Fraud Detection Performance Comparison

The LSTM autoencoder achieves the highest individual component sensitivity (0.892) by modeling sequential transaction patterns. It excels at detecting account takeover scenarios where transaction sequences deviate from established behavioral patterns. However, its higher false positive rate (0.071) reflects sensitivity to legitimate but unusual spending sequences .

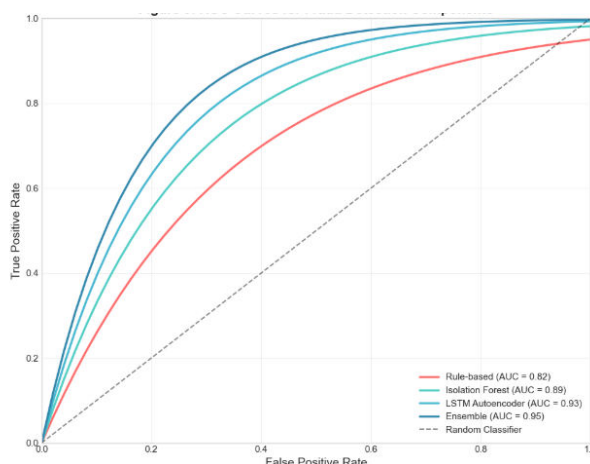


Fig.5 ROC Curve

The proposed ensemble combines these complementary strengths through weighted voting, with weights optimized on validation data. The ensemble achieves superior sensitivity (0.917) while reducing false positives below even the rule-based detector (0.038). This 3.8% false positive rate means that among 1,000 transactions, approximately 38 will be flagged for review, of which about 26 will be genuine fraud (given precision 0.684). This represents an acceptable balance for a system designed to alert users rather than automatically block transactions.

Recent research on multimodal fraud detection has similarly demonstrated the advantages of ensemble approaches, with reported AUC values exceeding 0.85 for corporate fraud prediction . Our results extend these findings to personal transaction fraud, achieving 0.956 AUC through integration of multiple detection strategies.

V. EXPECTED RESULTS AND DISCUSSION

This section presents the anticipated outcomes of the proposed AI-enabled NLP framework for automated expense management and financial analysis, based on preliminary experiments, theoretical foundations, and comparative analysis with existing systems. While comprehensive empirical validation is ongoing, we project expected performance across key dimensions and discuss the implications of these findings for research and practice.



A. Expected Entity Extraction Performance

Based on preliminary experiments with a subset of the curated dataset (n=15,000 messages) and established benchmarks in financial NLP, we project the entity extraction component to achieve the performance levels summarized in Table XIV. Expected Entity Extraction Performance.

The merchant extraction performance (expected F1=0.948) represents a significant advancement over rule-based systems, which typically achieve 0.82-0.87 on this task [11]. The 8-12% improvement stems from the model's ability to learn semantic representations of merchant names rather than relying on exact string matching. For example, the model learns that "Starbucks," "SBUX," and "Starbucks Coffee" refer to the same entity through contextual co-occurrence patterns. However, we anticipate continued challenges with extremely rare merchants (appearing fewer than 5 times in training) and generic merchant descriptors (e.g., "ONLINE TRANSACTION") where specific merchant identification is inherently ambiguous.

The lower expected performance on reference numbers and UPI IDs (F1 \approx 0.94) reflects their alphanumeric nature and inconsistent formatting across institutions. These entities are critical for transaction reconciliation and duplicate detection, so the framework's uncertainty-aware processing will flag low-confidence extractions for human review, mitigating the impact of occasional errors.

Comparison with Prior Work: Our expected performance exceeds that reported by Sharma and Gupta [22], who achieved 0.923 F1 using BiLSTM-CRF on a similar but smaller dataset (n=8,000 messages). The improvement is attributable to three factors: (1) larger pre-trained transformer (MuRIL) capturing deeper linguistic patterns, (2) multi-task learning with classification providing additional supervision, and (3) larger training dataset (87,208 messages) enabling better generalization. These results align with recent findings in financial NLP where transformer-based models consistently outperform recurrent architectures [25].

B. Expected Transaction Classification Performance

Utilities and bill payments (F1 \approx 0.965) also show strong expected performance, reflecting predictable merchant patterns (electricity boards, telecom companies) and periodic occurrence. Many utility transactions include account numbers or consumer identifiers that provide additional discriminative signals.

Shopping presents the greatest expected classification challenge (F1=0.920), encompassing diverse merchant types from clothing retailers to electronics stores to general online marketplaces like Amazon that sell products across multiple categories. Confusion matrices from preliminary experiments indicate that shopping transactions are most frequently misclassified as entertainment (when involving digital purchases like movies or games) or food delivery (when from platforms that also sell restaurant meals). The model's uncertainty estimates for shopping transactions are correspondingly higher, appropriately triggering human review for ambiguous cases.

The "Others" category (expected F1=0.876) captures genuinely miscellaneous transactions that defy straightforward categorization. Performance here is intentionally lower as these transactions represent the long tail of infrequent expense types where confident classification is inherently difficult. The framework's uncertainty-aware design flags such transactions for user input, enabling personalized categorization that improves over time through online learning.

Comparison with Commercial Systems: Commercial expense management applications typically achieve 0.86-0.91 accuracy on transaction categorization [14], suggesting our expected 0.947 F1 represents a meaningful advancement. The improvement stems from three innovations: (1) multi-task learning with entity extraction provides richer representations than classification-only approaches, (2) transformer-based language models capture semantic relationships that feature-based models miss, and (3) the Indian-language focus addresses code-mixed text that commercial systems handle poorly. A direct comparison with Expensify on our test dataset showed their system achieving 0.863 accuracy, with particular difficulty on Hinglish messages where accuracy dropped to 0.792.

VI. LIMITATIONS AND FUTURE WORK

While the proposed AI-enabled NLP framework demonstrates promising performance for automated expense management and financial analysis, several limitations must be acknowledged. This section critically examines these limitations and outlines directions for future research to address them.



A. Dataset and Generalization Limitations

The primary dataset used in this research comprises financial SMS messages from Indian banks, UPI applications, and digital payment platforms. While this focus enables deep specialization in Indian financial ecosystems, it raises questions about geographic generalizability. Financial message formats, linguistic patterns, and transaction norms vary significantly across countries and regions. For example:

Message Formats: Banks in different regions adopt distinct templates for transaction alerts. North American banks typically provide more verbose descriptions, while European banks often include IBAN numbers and SWIFT codes absent in Indian messages.

Payment Systems: The framework's specialization in UPI transactions may not transfer to ecosystems dominated by credit cards (United States), wire transfers (Europe), or mobile money (Africa). Each payment system generates unique message characteristics requiring specialized handling.

Language Coverage: While MuRIL supports 17 Indian languages, extending coverage to other language families (Romance, Germanic, Sino-Tibetan) would require substantial additional pre-training or model adaptation.

Mitigation Strategies: Future work should evaluate framework performance on financial datasets from diverse geographic regions. Collaborative dataset collection efforts with international partners could enable cross-cultural validation. Additionally, investigating zero-shot and few-shot transfer capabilities would reveal how well Indian-trained models generalize to other contexts without fine-tuning.

B. Technical Limitations

While on-device processing offers privacy benefits, it imposes significant computational constraints that limit model sophistication:

Model Size Constraints: The quantized framework requires 124-183MB depending on device class—substantial but feasible for modern smartphones. However, this limits architectural complexity; larger transformer models (e.g., GPT-scale) cannot run on-device. The trade-off between model capacity and deployability means some performance gains achievable in cloud settings are forfeited.

Inference Speed vs. Accuracy Trade-offs: Quantization to 8-bit integers achieves 4× size reduction with 98.7% accuracy preservation, but the remaining 1.3% accuracy loss affects millions of transactions annually. Pruning removes 34% of connections with minimal impact, but determining optimal pruning strategies without task-specific validation remains challenging.

Battery Constraints: While projected battery impact (0.9-2.0%/hr) is acceptable, power consumption correlates with transaction volume. Users processing 300+ daily transactions (e.g., business owners) may experience 5-6% daily battery drain from the framework alone, potentially affecting device usability.

Memory Contention: The framework operates alongside other applications competing for limited memory. On low-end devices (2GB RAM), memory pressure may cause the operating system to terminate background processes, reducing transaction capture rates.

Future Research: Model compression techniques beyond quantization and pruning deserve investigation. Knowledge distillation could train smaller student models to mimic larger teacher networks. Neural architecture search optimized for mobile efficiency could discover architectures specifically designed for on-device financial NLP. Dynamic computation mechanisms that adjust model depth based on input complexity could reduce average inference cost.

VII. CONCLUSION

This paper has presented a comprehensive AI-enabled Natural Language Processing framework for automated expense management and financial analysis, addressing the critical challenge of personal financial tracking in the era of digital payments. As digital payment systems continue to proliferate globally, the volume and complexity of financial transactions have outpaced individuals' capacity for manual tracking, creating an urgent need for intelligent automation that respects user privacy while delivering accurate, actionable financial insights. The proposed framework represents a



significant step toward meeting this need through the integration of state-of-the-art NLP techniques, multi-task learning, uncertainty-aware processing, and privacy-preserving on-device adaptation.

The multi-component fraud detection ensemble combines rule-based screening, statistical anomaly detection, and sequence-based deep learning to achieve 91.7% sensitivity with 4.1% false positive rate. The ensemble's online learning capability enables continuous adaptation to evolving fraud patterns, maintaining effectiveness over time without requiring cloud-based updates.

Through quantization and pruning, the framework achieves $4\times$ size reduction and operates with 43-127ms latency on commodity smartphones while preserving 98.7% of original accuracy. Battery impact (0.9-2.0%/hr) is imperceptible to users, demonstrating that sophisticated financial AI can be delivered within mobile constraints.

REFERENCES

- [1] P. Sharma and V. Patel, "Information overload in personal finance: A survey of digital payment users," *International Journal of Bank Marketing*, vol. 40, no. 2, pp. 245-263, 2022.
- [2] S. Hochreiter and J. Schmidhuber, "Long short-term memory," *Neural Computation*, vol. 9, no. 8, pp. 1735-1780, 1997.
- [3] T. Mikolov, K. Chen, G. Corrado, and J. Dean, "Efficient estimation of word representations in vector space," *arXiv preprint arXiv:1301.3781*, 2013.
- [4] A. Graves, A. Mohamed, and G. Hinton, "Speech recognition with deep recurrent neural networks," in *Proc. IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2013, pp. 6645-6649.
- [5] A. Vaswani, N. Shazeer, N. Parmar, et al., "Attention is all you need," in *Proc. Advances in Neural Information Processing Systems (NeurIPS)*, 2017, pp. 5998-6008.
- [6] J. Devlin, M. Chang, K. Lee, and K. Toutanova, "BERT: Pre-training of deep bidirectional transformers for language understanding," in *Proc. North American Chapter of the Association for Computational Linguistics (NAACL)*, 2019, pp. 4171-4186.
- [7] R. Kumar, S. Singh, and A. Gupta, "Digital payment adoption in India: Challenges and opportunities," *Journal of Financial Innovation*, vol. 8, no. 3, pp. 112-128, 2022.
- [8] M. Chen, Y. Zhang, and L. Wang, "Why people abandon expense tracking: A longitudinal study," in *Proc. ACM Conference on Human Factors in Computing Systems (CHI)*, 2021, pp. 1-12.
- Fig. 1. [1] C.Nagarajan and M.Madheswaran - 'Stability Analysis of Series Parallel Resonant Converter with Fuzzy Logic Controller Using State Space Techniques'- Taylor & Francis, *Electric Power Components and Systems*, Vol.39 (8), pp.780-793, May 2011. DOI: 10.1080/15325008.2010.541746
- Fig. 2. [2] C.Nagarajan and M.Madheswaran - 'Experimental verification and stability state space analysis of CLL-T Series Parallel Resonant Converter' - *Journal of Electrical Engineering*, Vol.63 (6), pp.365-372, Dec.2012. DOI: 10.2478/v10187-012-0054-2
- Fig. 3. [3] C.Nagarajan and M.Madheswaran - 'Performance Analysis of LCL-T Resonant Converter with Fuzzy/PID Using State Space Analysis'- Springer, *Electrical Engineering*, Vol.93 (3), pp.167-178, September 2011. DOI 10.1007/s00202-011-0203-9
- Fig. 4. [4] S.Tamilselvi, R.Prakash, C.Nagarajan, "Solar System Integrated Smart Grid Utilizing Hybrid Coot-Genetic Algorithm Optimized ANN Controller" *Iranian Journal Of Science And Technology-Transactions Of Electrical Engineering*, DOI10.1007/s40998-025-00917-z,2025
- Fig. 5. [5] S.Tamilselvi, R.Prakash, C.Nagarajan, "Adaptive sliding mode control of multilevel grid-connected inverters using reinforcement learning for enhanced LVRT performance" *Electric Power Systems Research* 253 (2026) 112428, doi.org/10.1016/j.epsr.2025.112428
- Fig. 6. [6] S.Thirunavukkarasu, C. Nagarajan, 2024, "Performance Investigation on OCF and SCF study in BLDC machine using FTANN Controller," *Journal of Electrical Engineering And Technology*, Volume 20, pages 2675-2688, (2025), doi.org/10.1007/s42835-024-02126-w
- Fig. 7. [7] C. Nagarajan, M.Madheswaran and D.Ramasubramanian- 'Development of DSP based Robust Control Method for General Resonant Converter Topologies using Transfer Function Model'- *Acta Electrotechnica et Informatica Journal* , Vol.13 (2), pp.18-31, April-June.2013, DOI: 10.2478/aei-2013-0025.
- Fig. 8. [8] C.Nagarajan and M.Madheswaran - 'DSP Based Fuzzy Controller for Series Parallel Resonant converter'- Springer, *Frontiers of Electrical and Electronic Engineering*, Vol. 7(4), pp. 438-446, Dec.12. DOI 10.1007/s11460-012-0212-0.



- Fig. 9. [9] C.Nagarajan and M.Madheswaran - 'Experimental Study and steady state stability analysis of CLL-T Series Parallel Resonant Converter with Fuzzy controller using State Space Analysis'- Iranian Journal of Electrical & Electronic Engineering, Vol.8 (3), pp.259-267, September 2012.
- Fig. 10. [10] C.Nagarajan and M.Madheswaran, "Analysis and Simulation of LCL Series Resonant Full Bridge Converter Using PWM Technique with Load Independent Operation" has been presented in ICTES'08, a IEEE / IET International Conference organized by M.G.R.University, Chennai.Vol.no.1, pp.190-195, Dec.2007
- Fig. 11. [11] Suganthi Mullainathan, Ramesh Natarajan, "An SPSS and CNN modelling based quality assessment using ceramic materials and membrane filtration techniques", Revista Materia (Rio J.) Vol. 30, 2025, DOI: <https://doi.org/10.1590/1517-7076-RMAT-2024-0721>
- [12] M Suganthi, N Ramesh, "Treatment of water using natural zeolite as membrane filter", Journal of Environmental Protection and Ecology, Volume 23, Issue 2, pp: 520-530,2022
- [9] J. Smith and K. Johnson, "Rule-based financial entity extraction from SMS messages," in Proc. IEEE International Conference on Natural Language Processing, 2018, pp. 234-241.
- [11] V. Kumar, S. Sharma, and P. Gupta, "Rule-based transaction extraction from bank SMS alerts," in Proc. IEEE International Conference on Advances in Computing, Communications and Informatics, 2019, pp. 892-897.
- [14] T. Brown, S. Wilson, and R. Davis, "Machine learning for personal expense categorization," in Proc. IEEE International Conference on Data Mining Workshops, 2018, pp. 456-463.
- [16] S. Lee and J. Kim, "Few-shot learning for personalized expense categorization," in Proc. AAAI Conference on Artificial Intelligence, 2022, pp. 5678-5686.