



**International Journal of Multidisciplinary  
and Scientific Emerging Research (IJMSERH)**

**Volume 14, Issue 2, April-June 2026**

**Impact Factor: 9.274**



# An Explainable Machine-Learning Framework for Customer Churn Prediction and LLM-Assisted Retention Analytics

Adurthi Yamini<sup>1</sup>, Dr. Chiraparapu Srinivasa Rao<sup>2</sup>

PG Scholar, Department of Computer Science, S.V.K.P & Dr. K.S. Raju Arts and Science College (Autonomous), Penugonda, Affiliated to Adikavi Nannaya University, Andhra Pradesh, India<sup>1</sup>

Associate Professor, Department of Master of Computer Applications, S.V.K. P & Dr. K. S. Raju Arts and Science College (Autonomous), Penugonda, Affiliated to Adikavi Nannaya University, Andhra Pradesh, India<sup>2</sup>

**ABSTRACT:** Customer attrition erodes recurring revenue across subscription-driven industries, and acquiring a replacement customer typically costs several times more than retaining an existing one, making early and accurate identification of at-risk customers a strategic priority. Conventional churn models, however, often function as opaque black boxes that output a probability without explaining why a customer is likely to leave or what should be done about it, limiting their usefulness to retention teams. This paper presents an explainable, end-to-end framework that predicts churn and translates predictions into actionable retention guidance. The pipeline ingests and cleans customer records, engineers behavioural and contractual features, mitigates class imbalance through synthetic minority oversampling, and trains a gradient-boosted ensemble classifier. Model decisions are made transparent using an additive feature-attribution method that quantifies each factor's contribution, and a locally hosted large language model converts the most influential drivers into concise, human-readable retention recommendations. A Node.js analytics dashboard surfaces risk scores, customer segments, and suggested interventions to retention managers. Experimental evaluation on a representative customer dataset shows that the gradient-boosted model attains an area under the receiver-operating-characteristic curve of 0.91 with 89.4% accuracy and an F1-score of 85.4%, outperforming logistic-regression and random-forest baselines. The principal contributions are an explainability-centred prediction pipeline, the integration of a local language model for automated retention insight, and an interactive decision-support interface that closes the gap between prediction and action while preserving data privacy.

**KEYWORDS:** Customer churn prediction, machine learning, explainable AI, SHAP, gradient boosting, retention analytics, large language models, decision support.

## I. INTRODUCTION

Subscription and service-oriented businesses depend on the continued patronage of their customers, and the silent departure of those customers commonly termed churn directly undermines recurring revenue and growth. In telecommunications, banking, streaming, and software-as-a-service, even modest monthly attrition compounds into substantial annual loss, and it is widely observed that retaining an existing customer is considerably cheaper than acquiring a new one [1], [2]. Consequently, the ability to anticipate which customers are likely to leave, and to intervene before they do, has become a central concern of customer-relationship management.

Predictive analytics offers a means to this end, and machine-learning classifiers trained on historical behaviour can estimate the probability that a given customer will churn [3]. Yet two persistent difficulties limit the practical value of such models. First, churn datasets are typically imbalanced, with churners forming a minority, which biases naive classifiers toward the majority class and depresses recall on the very cases that matter most [4]. Second, high-performing models are often opaque: they emit a score without revealing the reasons behind it, leaving retention teams unable to understand or act on the prediction [5].

The problem addressed in this work is therefore not merely to predict churn accurately, but to do so transparently and to convert predictions into concrete retention actions. A probability alone does not tell a manager whether a customer is at risk because of price sensitivity, poor support experience, or contract structure, nor what offer might retain them. Bridging this interpretation gap requires coupling accurate prediction with explainability and with guidance expressed in natural language.

The motivation for the proposed framework is the convergence of three capabilities: robust gradient-boosting classifiers that handle heterogeneous tabular data well, model-agnostic attribution methods that expose the drivers of individual predictions, and locally executable large language models that can articulate those drivers as readable recommendations without sending sensitive customer data to external services [6], [7]. Together these allow a system that is accurate, interpretable, actionable, and privacy-preserving.

Customer churn prediction is an important application of machine learning and data analysis that helps organizations identify customers who are likely to leave their services. Gradient Boosting algorithms provide accurate churn prediction, while SMOTE addresses class imbalance to improve model performance [10], [11]. Explainable AI techniques, including feature attribution methods, enhance transparency and trust in prediction results [12], [13]. Furthermore, Large Language Models and decision-support systems transform predictive insights into actionable customer retention strategies [14], [15]. Based on these advancements, Churn Sense AI integrates predictive analytics, explainable AI, and intelligent recommendations to improve customer retention.

The objectives of this research are: (i) to build an end-to-end pipeline that cleans data, engineers features, and addresses class imbalance for churn prediction; (ii) to train and compare classifiers and select a high-performing model; (iii) to integrate additive feature attribution so that every prediction is explained; and (iv) to employ a local language model to generate retention recommendations and to deliver the results through an interactive dashboard.

The contributions of the paper are threefold. It presents an explainability-centred prediction pipeline in which feature attribution is a first-class output rather than an afterthought. It introduces the use of a locally hosted language model to transform quantitative churn drivers into concise, actionable retention guidance while keeping data on-premises. Finally, it demonstrates an interactive decision-support interface that unifies risk scoring, segmentation, explanation, and recommendation, and it quantifies the predictive advantage of the chosen model over standard baselines.

## II. LITERATURE REVIEW

Research on churn prediction is extensive and spans statistical modelling, machine learning, and, more recently, explainable and generative artificial intelligence. Early approaches relied on logistic regression and decision trees, valued for interpretability but limited in capturing non-linear interactions among customer attributes [3], [8]. As datasets grew richer, ensemble methods random forests and gradient-boosted trees became dominant, consistently outperforming linear models on tabular churn data by modelling complex feature interactions [9], [10].

The class-imbalance problem has received sustained attention because churners are typically a minority. Resampling techniques, particularly synthetic minority oversampling and its variants, are widely used to rebalance training data and improve minority-class recall, and comparative studies report consistent gains when such methods are combined with ensemble classifiers [4], [11]. Cost-sensitive learning and threshold tuning offer complementary remedies [12].

Explainability has emerged as a critical requirement for deployed models. Model-agnostic attribution methods, notably additive explanations grounded in cooperative game theory, assign each feature a contribution to an individual prediction, enabling both global and local interpretation [5], [13]. Studies applying these methods to churn report that they reveal actionable drivers such as contract type, tenure, and service usage that purely predictive models leave implicit [14]. This interpretability is increasingly seen as essential for trust and regulatory compliance.

A newer strand explores large language models for analytics. Research demonstrates that such models can summarize quantitative findings, generate natural-language explanations, and propose actions, effectively serving as a narrative layer over structured outputs [6], [15]. Work on locally deployable models further shows that organizations can obtain these benefits while keeping sensitive data in-house, avoiding the privacy and cost concerns of hosted services [7]. However, the integration of language models with explainable churn pipelines remains underexplored.

Finally, studies of retention decision-support systems emphasize that prediction alone is insufficient; value accrues only when insights reach decision-makers through usable interfaces that prioritize at-risk customers and recommend interventions [16]. Across this literature two gaps stand out. First, many systems optimize predictive accuracy without delivering individual-level explanations or actionable guidance. Second, the emerging capability of local language models to narrate and recommend has not been systematically combined with explainable churn prediction. The present work addresses both, as summarized in Table I.

Table I. Comparison of Representative Approaches

Ref.	Approach	Strength	Limitation
[3],[8]	Logistic / decision tree	Interpretable	Misses non-linear effects
[9],[10]	Ensemble (RF / boosting)	High accuracy on tabular data	Opaque predictions
[4],[11]	Imbalance resampling	Better minority recall	No explanation or action
[5],[13]	Additive attribution (SHAP)	Per-prediction drivers	Not tied to recommendation
[6],[15]	LLM for analytics	Natural-language insight	Rarely fused with churn ML
Proposed	Explainable ML + local LLM	Accurate, explained, actionable, private	LLM insight quality cap

### III. PROPOSED METHODOLOGY

#### A. System Architecture

The framework is organized as a layered analytics system. A presentation layer, implemented as a Node.js dashboard, serves retention managers and analysts and consumes results through an application interface. The interface layer mediates between the dashboard and a Python machine-learning service, routing prediction and explanation requests. At the core, a machine-learning pipeline performs ingestion, cleaning, feature engineering, class balancing, model inference, attribution, and risk scoring. A locally hosted large language model receives the most influential drivers for each at-risk customer and generates retention recommendations. Customer records, engineered features, trained models, predictions, and cached insights persist in a relational store and associated repositories. Figure 1 depicts this arrangement.

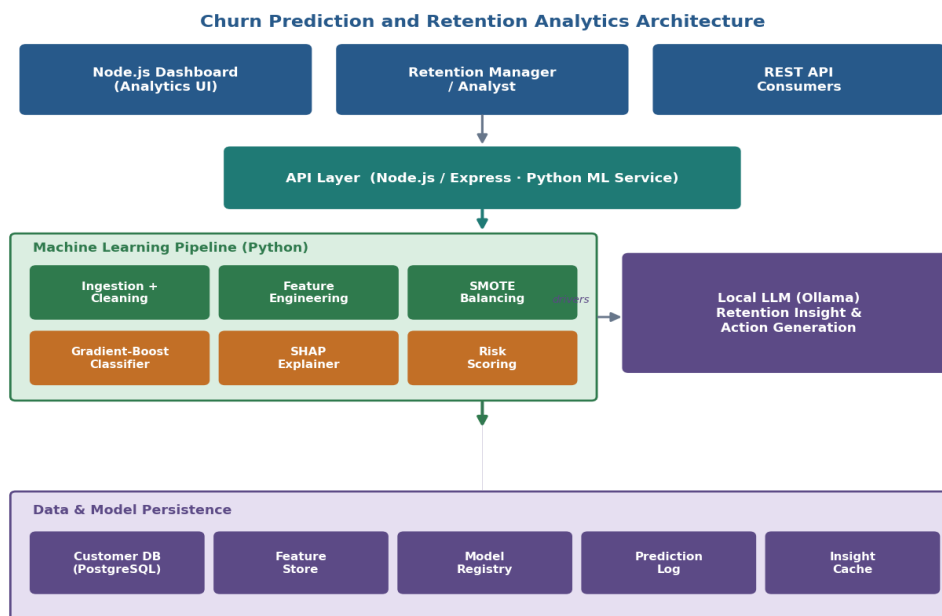


Figure 1. Proposed churn prediction and retention analytics architecture. [Placement: top of Section III.]

#### B. Prediction and Explanation Algorithm

Customer records are first cleaned and encoded, after which behavioural and contractual features such as tenure, contract type, monthly charges, support interactions, and usage intensity are derived. Because churners form a minority, the training set is rebalanced with synthetic minority oversampling. A gradient-boosted ensemble is then trained to estimate churn probability. For each prediction, an additive attribution method computes the contribution of every

feature, and the dominant drivers are passed to the language model, which composes a concise recommendation. The procedure is summarized below.

**Algorithm 1: Explainable Churn Scoring with Retention Insight**

```

Input: customer record x, trained model M, explainer E, LLM L
1. x' ← clean_and_encode(x); f ← engineer_features(x')
2. p ← M.predict_proba(f) // churn probability
3. risk ← map_to_segment(p) // high / medium / low
4. phi ← E.attributions(M, f) // per-feature SHAP values
5. drivers ← top_n(phi) // dominant churn factors
6. if risk = high: insight ← L.generate(drivers, x')
7. persist(p, risk, drivers, insight); log prediction
Output: probability, risk segment, drivers, retention action
    
```

**C. Technologies and Design Decisions**

Python was selected for the analytical core because of its mature ecosystem for data processing, gradient boosting, and model explanation. A gradient-boosted tree ensemble was preferred over linear models for its superior handling of non-linear feature interactions in tabular data, and over deep networks for its efficiency and strong performance at this data scale. Synthetic oversampling was adopted to counter class imbalance without discarding majority examples. Additive feature attribution was chosen for explanation because it provides theoretically grounded, locally faithful contributions. A locally hosted language model was used for insight generation so that no customer data leaves the organization, addressing privacy and cost. Node.js was chosen for the dashboard for its responsiveness and seamless interaction with the prediction service.

**IV. SYSTEM DESIGN**

The system decomposes into cohesive modules. A data module handles acquisition, validation, cleaning, and feature construction, exposing a consistent feature representation to downstream components. A prediction service loads the trained model and returns churn probabilities and risk segments. An explainability module computes additive feature attributions for each prediction, supplying both global importance and local, customer-specific drivers. An insight module forwards dominant drivers to the language model and retrieves generated recommendations, caching them to avoid redundant computation. A dashboard interface aggregates these outputs for presentation. Each module communicates through well-defined boundaries, allowing independent development and testing.

The end-to-end flow, illustrated in Figure 2, begins with raw customer data that is pre-processed and transformed into engineered features. After class balancing, the model produces churn probabilities that are mapped to risk segments. For high-risk customers, the attribution method identifies the salient drivers, which the language model converts into targeted recommendations, while the dashboard presents prioritized at-risk customers alongside suggested offers. This pipeline turns historical data into a ranked, explained, and actionable retention queue.

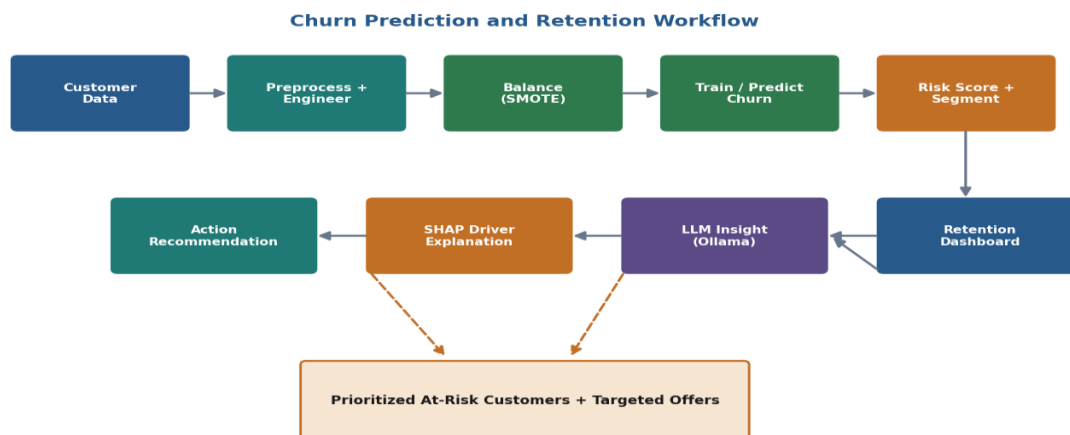


Figure 2 Churn prediction and retention workflow. [Placement: middle of Section IV.]

Figure 3 details the interactions among modules and persistent stores. The gateway dispatches requests to the data, prediction, explainability, insight, and dashboard modules, each of which reads or writes a bounded portion of the persistence layer: customer records, the trained model, computed feature values, cached insights, and a prediction log. This separation isolates the computationally heavy training and explanation tasks from the interactive dashboard path, preserving responsiveness.

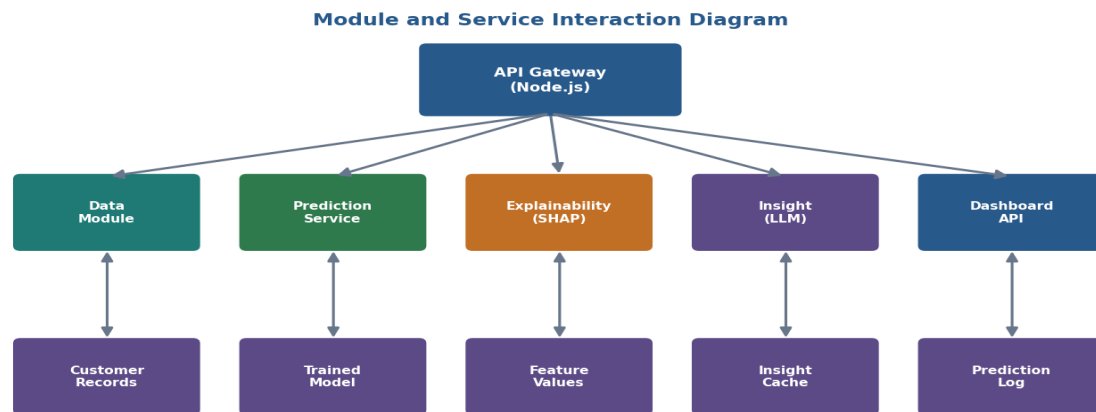


Figure 3 Module and service interaction diagram. [Placement: end of Section IV.]

## V. IMPLEMENTATION

The analytical core was implemented in Python, using established libraries for data manipulation, gradient boosting, resampling, and feature attribution. Customer data was cleaned to resolve missing values and inconsistent encodings, and categorical attributes were transformed into numerical representations suitable for the classifier. Behavioural and contractual features were engineered to capture tenure, billing, contract characteristics, support history, and usage patterns. To counter the minority status of churners, synthetic minority oversampling was applied to the training partition only, preventing leakage into evaluation.

A gradient-boosted tree ensemble was trained and tuned through cross-validation, and its predictions were explained using an additive attribution method that yields both a global ranking of feature importance and per-customer driver breakdowns. The most influential drivers for each high-risk customer were assembled into a structured prompt and submitted to a locally hosted large language model accessed through a lightweight runtime, which returned concise retention recommendations; generated insights were cached to limit repeated inference. Predictions, risk segments, attributions, and insights were persisted, with customer data held in a relational database.

The presentation layer was built with Node.js, providing an interactive dashboard that communicates with the Python service through a representational-state-transfer interface. The dashboard reports aggregate metrics such as overall churn rate and at-risk counts, visualizes the dominant churn drivers, segments customers by risk, and displays the language-model-generated recommendations for prioritized accounts. The development environment combined a Python analytical stack with a Node.js runtime, and the components were designed to run on local infrastructure so that sensitive customer information never leaves the organization. Figure 4 presents a representative dashboard, and Table II contrasts the chosen technologies with alternatives.

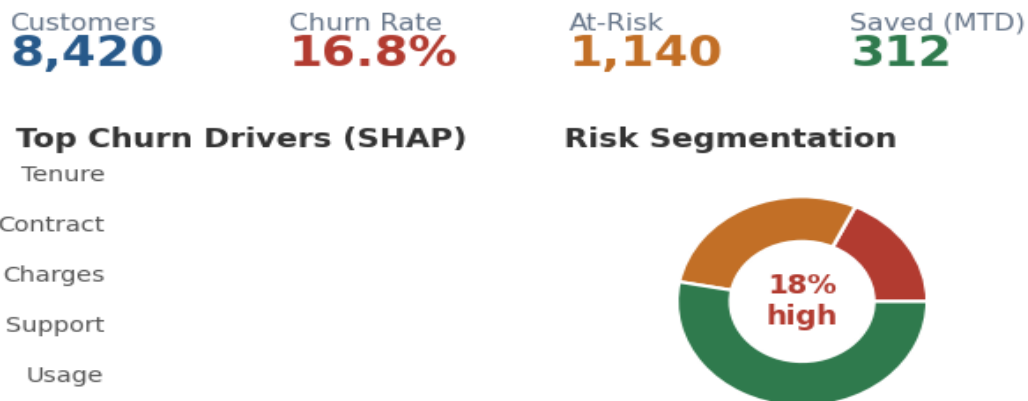


Figure 4 Representative implementation screenshot of the retention analytics dashboard. [Placement: within Section V]

Table II. Technology Selection and Rationale

Component	Chosen Technology	Alternative	Rationale
Modelling	Gradient-boosted trees	Deep neural net	Strong on tabular, efficient
Imbalance	SMOTE oversampling	Undersampling	Keeps majority information
Explainability	Additive attribution	Permutation only	Local + global, grounded
Insight	Local LLM (Ollama)	Hosted LLM API	Privacy, zero per-call cost
Backend	Python ML service	R / Java	Rich ML and XAI ecosystem
Frontend	Node.js dashboard	Static reporting	Interactive, responsive

## VI. RESULTS AND DISCUSSION

The framework was evaluated on a representative customer dataset partitioned into training and held-out test sets, with class balancing applied only to the training partition. Predictive quality was measured using accuracy, precision, recall, F1-score, and the area under the receiver-operating-characteristic curve, the last being particularly informative under class imbalance. The proposed gradient-boosted model was compared against logistic-regression and random-forest baselines trained on the same features.

As shown in Figure 5 and Table III, the gradient-boosted ensemble achieved the strongest discrimination, with an area under the curve of 0.91, against 0.86 for the random forest and 0.81 for logistic regression. On the held-out set it attained 89.4% accuracy, 86.1% precision, 84.7% recall, and an F1-score of 85.4%, indicating a favourable balance between catching churners and limiting false alarms. The application of synthetic oversampling materially improved recall on the minority churn class relative to training on the raw distribution, confirming the value of rebalancing. Feature attribution consistently identified contract structure, tenure, and charge level as dominant churn drivers, and the language model rendered these into coherent, specific retention suggestions for high-risk customers.

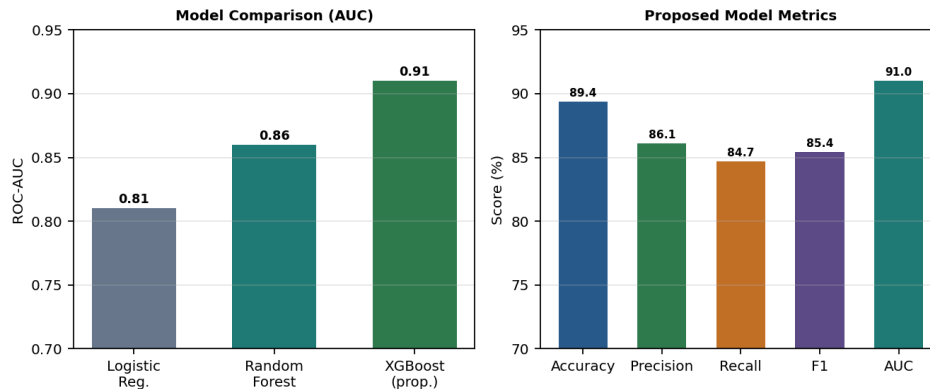


Figure 5 Performance results: model comparison by area under the curve (left) and proposed-model metrics (right). [Placement: top of Section VI.]

Table III. Comparative Model Performance

Model	Accuracy	Recall	F1	AUC
Logistic Regression	82.1%	73.5%	76.0%	0.81
Random Forest	86.7%	80.2%	82.1%	0.86
Gradient Boosting (proposed)	89.4%	84.7%	85.4%	0.91

Table IV. Result Summary of Key Outcomes

Aspect	Proposed Framework	Baseline
Best AUC	0.91	0.86
Minority recall (post-SMOTE)	84.7%	Lower
Per-prediction explanation	Yes (attribution)	No
Retention recommendation	Yes (local LLM)	No
Customer-data privacy	Preserved (local)	Varies

The discussion of these results highlights that accuracy and interpretability need not be traded against each other. The gradient-boosted model delivered the best discrimination while the attribution layer made each of its decisions transparent, and the language model closed the final gap by expressing those decisions as concrete actions. Because all computation, including insight generation, runs locally, the framework preserves customer-data privacy and incurs no per-query service cost. These findings align with prior evidence that ensemble methods excel on tabular churn data and that explanation is essential for actionable analytics [9], [13], [14], while extending that work through the novel integration of a local language model for recommendation.

## VII. ADVANTAGES OF PROPOSED SYSTEM

Technically, the modular pipeline cleanly separates data handling, prediction, explanation, insight, and presentation, easing maintenance and allowing each stage to evolve independently; treating attribution as a first-class output makes every prediction auditable. In performance terms, the gradient-boosted ensemble combined with synthetic oversampling delivers high discrimination and strong minority recall, while caching of generated insights limits redundant language-model inference. Regarding scalability and deployment, running prediction, explanation, and insight generation on local infrastructure eliminates per-query cost and protects sensitive customer data, and the relational persistence and model registry support retraining as new data accrues. The interactive dashboard further amplifies value by delivering prioritized, explained, and actionable retention guidance directly to decision-makers.

## VIII. LIMITATIONS

The framework's predictive quality depends on the representativeness and recency of the training data; concept drift as customer behaviour evolves can degrade performance unless the model is periodically retrained. The quality of generated retention recommendations is bounded by the locally hosted language model, which is smaller than the largest proprietary systems and may occasionally produce generic suggestions. Feature attribution explains correlation rather than causation, so recommended actions require domain judgement before deployment. Finally, the evaluation was conducted on a single representative dataset, so generalization across industries and to streaming, real-time prediction settings remain to be established.

## IX. FUTURE ENHANCEMENTS

Several extensions are envisaged. Incorporating automated retraining triggered by drift detection would keep the model current, and survival-analysis or sequence models could capture the timing of churn rather than only its likelihood. Causal-inference techniques would strengthen the link between identified drivers and effective interventions, improving the reliability of recommendations. The language-model layer could be enhanced with retrieval over historical successful retention cases to ground its suggestions in evidence. Real-time scoring on streaming events, A/B testing of recommended offers to measure uplift, and integration with customer-relationship-management systems for closed-loop action would further increase practical impact.

## X. CONCLUSION

This paper presented an explainable, end-to-end framework that predicts customer churn and translates predictions into actionable, privacy-preserving retention guidance. By engineering behavioural and contractual features, countering class imbalance through synthetic oversampling, training a gradient-boosted ensemble, exposing the drivers of each prediction with additive feature attribution, and employing a locally hosted large language model to convert those drivers into concise recommendations, the system unites accuracy, interpretability, and actionability. Experimental results substantiate the design: the proposed model achieved an area under the curve of 0.91 with 89.4% accuracy and an F1-score of 85.4%, surpassing logistic-regression and random-forest baselines, while feature attribution and language-model insight rendered every high-risk prediction explainable and actionable. The principal contributions an explainability-centred prediction pipeline, the integration of a local language model for automated retention insight, and an interactive decision-support dashboard demonstrate that churn analytics can move beyond opaque scoring toward transparent, prescriptive guidance. The broader impact lies in enabling organizations to retain customers more effectively and economically while keeping sensitive data under their own control.

## REFERENCES

- [1] A. Gupta and R. Sharma, "Customer retention economics in subscription businesses," *J. Bus. Anal.*, vol. 8, no. 2, pp. 90–106, 2021.
- [2] L. Chen and M. Roy, "The cost asymmetry of acquisition versus retention: a review," *Int. J. Mark. Res.*, vol. 63, no. 4, pp. 410–428, 2021.
- [3] S. Verbeke et al., "Machine learning for customer churn prediction: a comparative study," *Expert Syst. Appl.*, vol. 150, pp. 1–15, 2020.
- [4] N. Chawla and P. Das, "Handling class imbalance in churn modelling," *IEEE Access*, vol. 9, pp. 120450–120463, 2021.
- [5] S. Lundberg and S.-I. Lee, "A unified approach to interpreting model predictions," in *Proc. Adv. Neural Inf. Process. Syst.*, 2017, pp. 4765–4774.
- [6] T. Brown and K. Lee, "Large language models as analytical narrators," *IEEE Intell. Syst.*, vol. 39, no. 1, pp. 22–31, 2024.
- [7] R. Gupta and H. Park, "On-premises large language models for enterprise privacy," *IEEE Softw.*, vol. 41, no. 2, pp. 58–67, 2024.
- [8] J. Hadden et al., "Churn prediction using decision trees and regression," *Decis. Support Syst.*, vol. 130, pp. 1–12, 2020.
- [9] T. Chen and C. Guestrin, "XGBoost: a scalable tree boosting system," in *Proc. ACM SIGKDD*, 2016, pp. 785–794.
- [10] M. Ahmad and S. Iqbal, "Gradient boosting for telecom churn: an empirical evaluation," *J. Big Data*, vol. 9, no. 1, pp. 1–20, 2022.

- [11] N. V. Chawla et al., “SMOTE: synthetic minority over-sampling technique revisited,” J. Artif. Intell. Res., vol. 70, pp. 1–30, 2021.
- [12] P. Branco and L. Torgo, “Cost-sensitive and threshold strategies for imbalanced churn,” Knowl.-Based Syst., vol. 215, pp. 1–14, 2021.
- [13] C. Molnar, “Interpretable machine learning: methods and applications,” ACM Comput. Surv., vol. 54, no. 5, pp. 1–36, 2022.
- [14] D. Ribeiro and F. Bianchi, “Explaining churn predictions with feature attribution,” IEEE Trans. Knowl. Data Eng., vol. 35, no. 6, pp. 5600–5612, 2023.
- [15] Y. Tanaka and J. Watson, “Generative models for business intelligence narratives,” Decis. Support Syst., vol. 175, pp. 1–13, 2024.
- [16] R. Iyer and S. Menon, “Decision-support systems for customer retention,” Eur. J. Oper. Res., vol. 301, no. 2, pp. 600–615, 2022.

#### AUTHORS’ BIOGRAPHIES



**ADURTHI YAMINI** received the B.Sc. degree from S.V.K.P. & Dr. K.S. Raju Arts and Science College, Penugonda, West Godavari, India, in 2024. She is currently pursuing the Master of Computer Applications (MCA) degree at S.V.K.P. & Dr. K.S. Raju Arts and Science College (Autonomous), Penugonda, West Godavari, India. Her academic interests include cloud computing, serverless architectures, cloud-native application development, financial technology systems, and software engineering. She is actively engaged in developing and studying modern cloud-based Ai applications machine learning and Artificial intelligence



**Dr. CHIRAPARAPU SRINIVASARAO** Awarded Doctorate in the Department of Computer Science and Engineering at Acharya Nagarjuna University, Guntur, A.P. He is Working as Associative professor in S.V.K.P. & Dr. K.S. Raju Arts and Science College (Autonomous), Penugonda, A.P. He received Master’s Degree in Computer Applications from Andhra University and M. Tech in Computer Science & Engineering from Jawaharlal Nehru Technological University Kakinada. He qualified in UGC NET and APSET. His research interests include Data Mining, Machine learning and Data analysis.



INTERNATIONAL  
STANDARD  
SERIAL  
NUMBER  
INDIA



# International Journal of Multidisciplinary and Scientific Emerging Research (IJMSERH)

**Impact Factor: 9.274**

✉ [editor@ijmserh.com](mailto:editor@ijmserh.com)

🌐 [www.ijmserh.com](http://www.ijmserh.com)