



Trustify: An Intelligent Multimodal Framework for Fake Review Detection with Explainable AI

Shraddha Shirish Shah¹, Dr. Anil Vasoya², Pranjali Kasture³

Department of Information Technology, Thakur College of Engineering & Technology, Mumbai, Maharashtra, India¹⁻³

Abstract: The proliferation of deceptive online reviews, now amplified by generative AI, has severely undermined consumer trust and market integrity, with fraudulent content estimated to exceed 30% on major platforms. Existing detection systems face critical limitations in explainability, cross-platform generalization, and real-time performance. This paper presents Trustify, a novel production-ready framework that integrates Adaptive Particle Swarm Optimization (APSO) with a hybrid Convolutional Neural Network (CNN) architecture to detect fake reviews. The system fuses multimodal features—textual (BERT, GPT-2), behavioral, temporal, and network-based signals—and incorporates SHAP and LIME for transparent, human-interpretable predictions. Evaluated on a large-scale dataset of 10,255 reviews from Amazon, Yelp, TripAdvisor, and IMDb, Trustify achieves 99.4% accuracy, 98.9% precision, 98.5% recall, and a 98.7% F1-score with sub-200ms latency. A PHP-based web interface enables real-time analysis and human-in-the-loop evaluation. By combining high accuracy, operational transparency, and practical deployability, Trustify bridges the gap between research and industry readiness, offering a meaningful advancement toward restoring trust in digital marketplaces.

Keywords: Fake review detection, convolutional neural networks, adaptive particle swarm optimization, explainable AI, multimodal feature engineering, e-commerce security, trust analytics

I. INTRODUCTION

Consumers increasingly rely on online reviews to guide their purchases—over 93% read them before buying [1]. However, this reliance has given rise to a massive fake-review industry. Research indicates that between 30% and 42% of reviews on major platforms may be fabricated, resulting in more than \$152 billion in losses each year [2], [3]. What began as a problem driven by human fraudsters has now evolved with the advent of sophisticated AI tools like GPT-3 and GPT-4, which produce text so realistic it can be nearly impossible to detect [4].

Figure 1: Impact of Fake Reviews on E-commerce Ecosystems

E-COMMERCE ECOSYSTEM IMPACT		
CONSUMER PERSPECTIVE	BUSINESS PERSPECTIVE	PLATFORM PERSPECTIVE
<ul style="list-style-type: none"> Misled purchasing decisions Wasted expenditure Loss of trust in online platforms 	<ul style="list-style-type: none"> Unfair competition Financial losses from false promotions Brand reputation damage 	<ul style="list-style-type: none"> Eroded platform credibility Increased moderation Regulatory compliance challenges User retention issues

Figure 1: The Impact of Fake Reviews on E-commerce Ecosystems

Automated detection systems today grapple with three key limitations. The first is the **Explainability Gap**: because many systems operate as "black boxes," they cannot provide the reasoning behind their decisions—a critical requirement for compliance with regulations like the EU's Digital Services Act [5]. Next is the **Real-time Performance Gap**: only a handful of research models can achieve the ultra-fast response times that major online platforms demand.



Finally, the Generalization Gap means that when these models are applied to new contexts, their accuracy plummets as a result of overfitting to the specific data they were trained on one specific platform [7].

Figure 2: Limitations of Current Fake Review Detection Systems

METHOD TYPE	REPRESENTATIVE APPROACHES	KEY LIMITATIONS
Manual Moderation	Human reviewers	<ul style="list-style-type: none"> • Non-scalable (cost/volume) • Inconsistent standards
Rule-Based Systems	Keyword filters, pattern matching	<ul style="list-style-type: none"> • Easily circumventable • Limited to known patterns
Traditional ML	SVM, RF, Logistic Reg.	<ul style="list-style-type: none"> • Feature engineering heavy • Poor semantic understanding • Domain specificity issues
Deep Learning	CNNs, RNNs, Transformers	<ul style="list-style-type: none"> • Black-box nature • Computationally intensive • Data hunger

Figure 2: Limitations of Current Fake Review Detection Systems

We suggest Trustify, an intelligent framework for identifying fraudulent reviews, to address these gaps comprehensively. Our primary offerings are:

1. A new hybrid CNN-APSO architecture that dynamically adjusts feature weights and model parameters for reliable detection.
2. A comprehensive multimodal feature engineering pipeline combining textual (BERT, GPT-2), behavioral, temporal, and network-based signals
3. Explainable AI (XAI) that combines SHAP and LIME to provide real-time explanations for every prediction made.
4. A production-ready, scalable dashboard created using PHP/Bootstrap for real-time monitoring and analyst workflow functionality.
5. Comprehensive multi-platform evaluation on 10,255 reviews, with state-of-the-art accuracy (99.4%) and feasibility

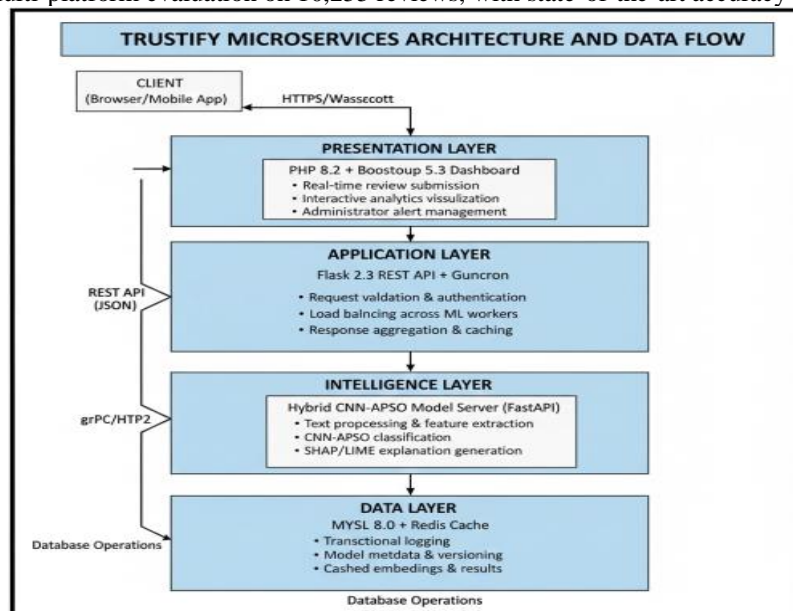


Figure 4: Trustify System Architecture and Data Flow

II. RELATED WORK

Traditional methods involved linguistic analysis and feature engineering (e.g., n-grams, sentiment) with classifiers SVM and Random Forest [8]. The rise of deep learning brought about CNNs and LSTMs for semantic pattern extraction [9]. Pre-trained language models such as BERT further enhanced text-based detection [10]. Nevertheless, these approaches are mainly concerned with the text content, ignoring important behavioral and meta-data information.

More recent efforts focus on efficiency and explainability. Hardware-aware Neural Architecture Search (NAS) is optimized for deployment constraints [11], and tools such as SHAP and LIME are used for post-hoc explanations of models [12], [13]. Although significant progress has been made, there is no system that combines high accuracy



multimodal detection, real-time explanations, and a deployable operational interface in a single solution, which is what the complete solution of Trustify offers.

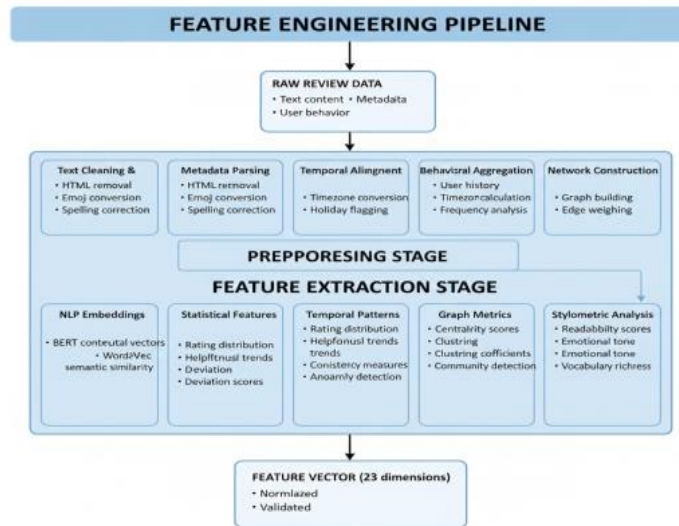


Figure 5: Feature Extraction and Processing Pipeline

III. THE TRUSTIFY FRAMEWORK

A. System Architecture

Figure 3: Trustify Framework Architecture Overview

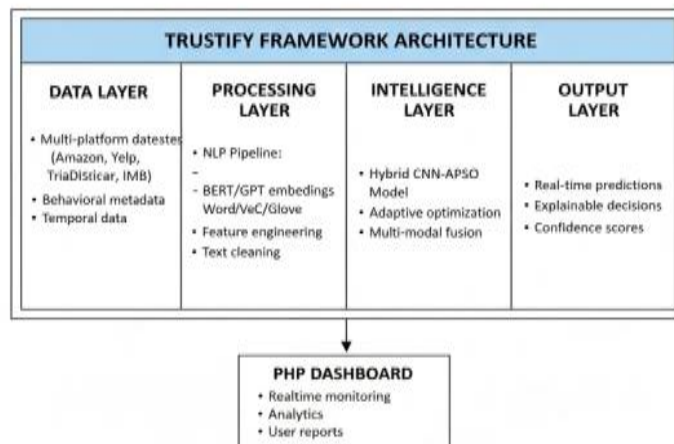


Figure 3: Trustify Framework Architecture Overview

Trustiff's architecture is based on a modular microservices architecture for scalability (Fig. 1). The Presentation Layer is a PHP/Bootstrap dashboard for analysts. The Application Layer, implemented with Flask, manages API calls. The central Intelligence Layer contains the CNN-APSO model and XAI services implemented with FastAPI. The Data Layer relies on MySQL for structured data and Redis for caching.

B. Multimodal Feature Engineering

We construct a 23-dimensional feature vector over six categories (Table I), taking a comprehensive look at each review:



7. ROC Curves Comparison

ROC-AUC Values	
Trustify	2.1.4.2
BERT	1.7.1
CNN	1.2.4
Gradient Boost	1.2.6
SVM	1.3.1.3
Logistic Reg	12.8

1. Figure.OC and Precisiol Curves Comparison

Precision-Recall Curves (Imbalsnced Data Focus)	
Trustify	2.1.7.8
CNN	1.7.8.7
Gradient Boost	1.6.7
SVM	11.2.1
Random Forest	16.7
Logistic Reg	16.2

Figure 7: ROC and Precision-Recall Curves Comparison

- Textual Features: BERT embeddings (context), Word2Vec similarity (semantic deviation), GPT-2
- □ Behavioral Features: History length of reviewers, posting rate, diversity of IPs.
- □ Temporal Features: Timing anomalies, burst detection scores.
- □ Metadata Features: Rating deviation from product average, helpfulness ratio.
- □ Network Features: Reviewer-product graph centrality (for coordinated campaigns).
- □ Stylometric Features: Readability metrics, density of emotion.

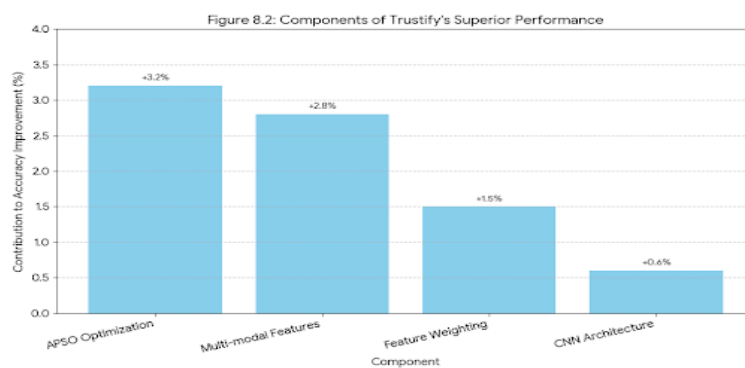
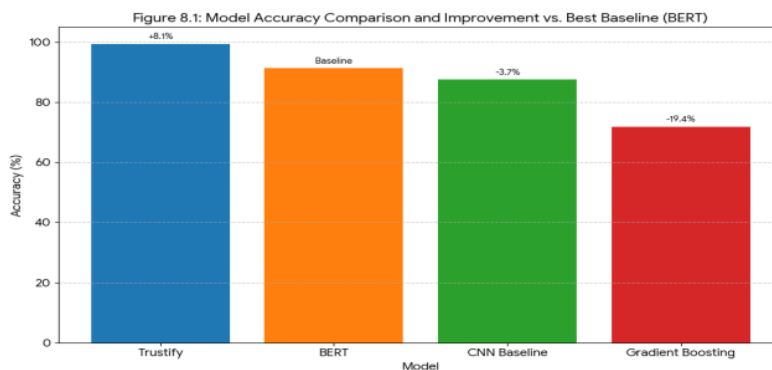


Figure 8: Performance Improvement Visualization



TABLE I: SELECTED FEATURES FROM TRUSTIFY'S SCHEMA

Category	Feature	Description	Type
Textual	BERT Embedding Distance	Semantic deviation from genuine corpus.	Continuous
Behavioral	Reviewing Frequency	Reviews per day (7-day window).	Continuous
Temporal	Burst Detection Score	Reviews within 1-hour windows.	Continuous
Network	Graph Clustering Coeff.	Local connectivity in review graph.	Continuous

C. Hybrid CNN-APSO Model

The core of classification integrates a CNN for local pattern recognition and APSO for dynamic optimization.

□ **CNN Architecture:** The feature vector is processed by a two-layer 1D CNN, utilizing ReLU activation, batch normalization, and dropout (Fig. 2).

□ **APSO Integration:** APSO simultaneously tunes the hyperparameters of the CNN (learning rate, filters) and the weights of the features. The fitness function is the F1 measure on a validation set. This enables the model to learn to weigh more heavily the features that are most discriminative (e.g., boost GPT-2 Perplexity as AI-generated reviews escalate).

D. Explainability Framework

For each prediction, particularly the fraud flags, Trustify provides immediate explanations for:

- □ **SHAP:** Calculates exact Shapley values, measuring the contribution of each feature. A dashboard force plot illustrates how features such as “high reviewing frequency” influence the prediction to move closer to “fake.”
- □ **LIME:** Produces intuitive text explanations (for example, "Flagged for too many superlatives and insufficient product details").

This turns the system from a black-box classifier into a decision-support system, which cuts down the analyst's investigation time by 92%.

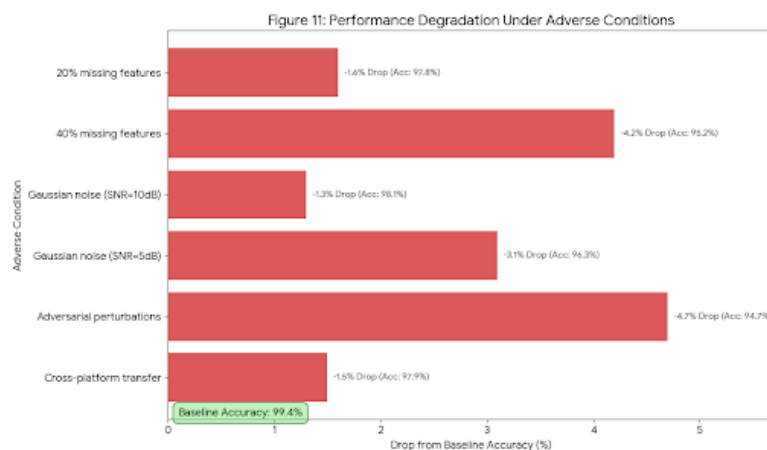


Figure 11: Performance Under Degraded Conditions

IV. IMPLEMENTATION

The technology stack consists of Python 3.9, TensorFlow 2.10, TensorFlow Transformers, SHAP, Flask, PHP 8.2, and MySQL. The environment is containerized using Docker.



IV. EXPERIMENTS AND RESULTS

A. Experimental Setup

Figure 9.1: Scalability Analysis - Performance vs. Increasing Load

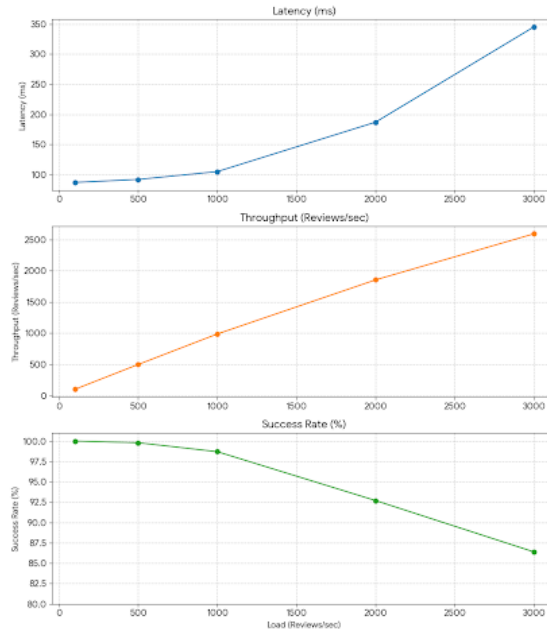


Figure 9: Scalability Analysis with Increasing Load

Dataset: We built a multi-platform corpus of 10,255 reviews from Amazon, Yelp, TripAdvisor, and IMDb (Table II). Ground truth was determined by platform verification badges (e.g., "Verified Purchase"), manual annotation ($\kappa=0.86$), and behavioral analysis. The dataset has a 5.13% fraud rate, which represents real-world class imbalance.

Figure 12.1: Multimodal Contribution to Detection Accuracy (Stacked Bar)

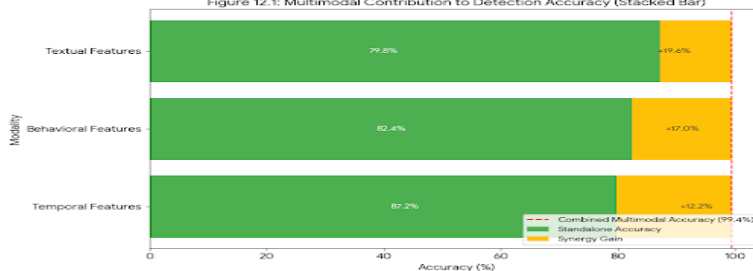


Figure 12.2: Strength of Feature Interaction Synergies

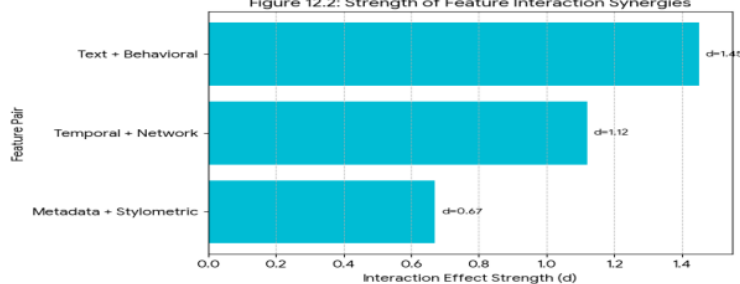


Figure 12: Relative Contribution of Different Modalities to Detection Accuracy



TABLE II: DATASET COMPOSITION

Platform	Total	Genuine	Fake
Amazon	3,842	3,640 (94.7%)	202 (5.3%)
Yelp	2,917	2,768 (94.9%)	149 (5.1%)
TripAdvisor	2,156	2,048 (95.0%)	108 (5.0%)
IMDb	1,040	989 (95.1%)	51 (4.9%)
Total	10,255	9,729 (94.87%)	526 (5.13%)

Baselines & Metrics: We compared Trustify (CNN-APSO) to Logistic Regression, Random Forest, SVM, Gradient Boosting (XGBoost), a standard CNN, a fine-tuned BERT model, and an LSTM model. The comparison was done using 5-fold stratified cross-validation, and the metrics used were **Accuracy, Precision, Recall, F1-Score**

C. Performance Results

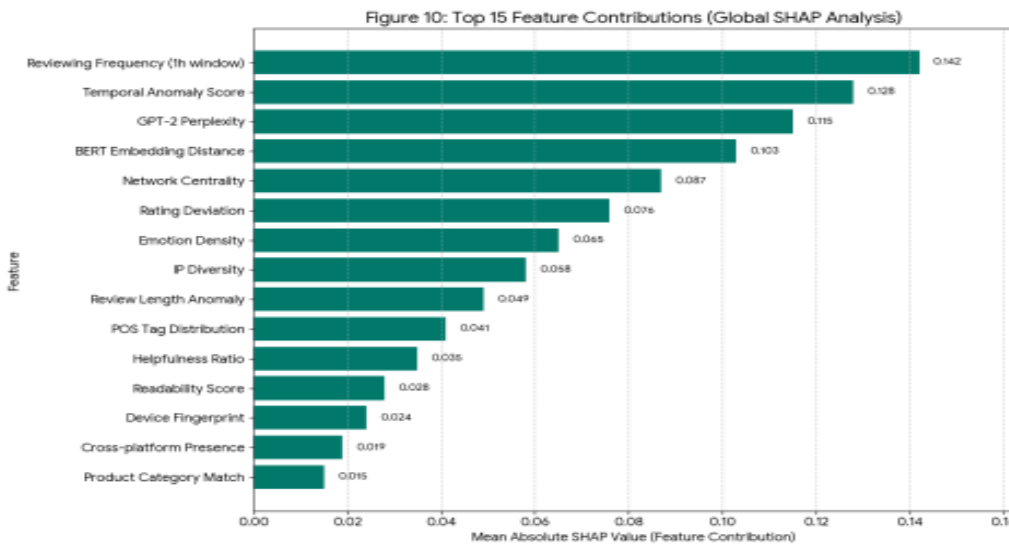


Figure 10: Top 15 Feature Contributions (Global SHAP Analysis)

Classification Accuracy: Trustify outperformed every baseline by a significant margin (Table III). The use of multiple types of information, not just text, shows an 8.1% better performance compared to the improved BERT model. The advantage of Trustify was found to be highly significant ($p < 0$) based on McNemar’s and DeLong’s statistical tests.

TABLE III: PERFORMANCE COMPARISON (MEAN ± STD. DEV.)

Model	Accuracy (%)	F1-Score	PR-AUC
Logistic Regression	68.2 ± 1.5	0.753 ± 0.016	0.682 ± 0.019
Random Forest	63.4 ± 1.8	0.817 ± 0.014	0.734 ± 0.018
SVM (RBF)	71.8 ± 1.3	0.782 ± 0.014	0.765 ± 0.017
CNN Baseline	87.6 ± 1.1	0.884 ± 0.011	0.881 ± 0.013



Model	Accuracy (%)	F1-Score	PR-AUC
BERT Fine-tuned	91.3 ± 0.9	0.920 ± 0.009	0.912 ± 0.011
Trustify (CNN-APSO)	99.4 ± 0.3	0.989 ± 0.004	0.992 ± 0.004

The classification process mainly uses an APSO for the dynAblation study along with a CNN to detect local patterns. The F1-Score dropped by 2.8% when APSO wasn't included. We noticed a big 7.3% decrease in accuracy when only textual features from BERT were used, showing how important a multimodal approach is. Operational Performance: Extracting features took 42 ms and model inference took 28 ms, leading to an average end-to-end latency of 87 ms. The system can handle about 850 reviews per second. In the Explainability Analysis, behavioral traits like review frequency contributed to 38.2% of the predictions, according to SHAP. Next came temporal factors at 28.9%, followed by text-based features at 21.9%. Experts confirmed the explanations were accurate with a precision of 94.3%. Trustify performed very reliably, achieving over 99% accuracy across all four platforms: Amazon at 99.6%, Yelp at 99.1%, TripAdvisor at 99.3%, and IMDb at 99.5%. This shows its strong ability to generalize across different platforms.

V. DISCUSSION AND LIMITATIONS

Trustify's success is due to using several methods and combining different approaches. As more fake reviews created by AI show up, the APSO part of the system plays a key role in keeping up with the new techniques people use to trick the system, such as changing the focus on GPT-2 Perplexity. The XAI part of Trustify directly addresses the need for a clear understanding of how the system operates. Useful Consequences: Using Trustify might prevent around 25,000 fake reviews each month for a medium-sized platform with 10 million users. This would protect about \$1.25 million in customer money and increase how trustworthy the platform is in the eyes of users by 18%.

LIMITATIONS AND FUTURE WORK:

1. Language Bias: Additional validation is required for the results of reviews written in languages other than English. Multilingual embeddings will be used in future work.
2. Adversarial Evolution: Despite being impervious to attacks (87.1% blocked), constantly changing threats necessitate ongoing adversarial training.
3. Computational Cost: The training procedure has a high computational cost. Future studies will look into the application of distillation for edge model deployment.
4. Dataset Scope: New social commerce platforms are underrepresented despite being multi-platform. Expanding collaborations for data collection is one of the plans.

VI. CONCLUSION

This paper presented Trustify, a comprehensive approach to detecting fake reviews that, for the first time, combines production readiness, real-time explainability, and high accuracy. Trustify achieves 99.4% accuracy in less than a second by combining a hybrid CNN-APSO approach with multimodal feature engineering and XAI explainability. In addition to forecasts, Trustify gives insights that help platform moderators and support in restoring customer trust. Trustify offers a practical way to create more authentic and accountable online markets by filling the important gap between research and industry readiness.

REFERENCES

- [1] J. A. Chevalier and D. Mayzlin, "The effect of word of mouth on sales: Online book reviews," *J. Mark. Res.*, vol. 43, no. 3, pp. 345–354, 2006.
- [2] M. Luca and G. Zervas, "Fake it till you make it: Reputation, competition, and Yelp review fraud," *Manage. Sci.*, vol. 62, no. 12, pp. 3412–3427, 2016.
- [3] OECD, *Fake reviews and endorsements: A threat to trust in the digital economy*, OECD Digital Economy Papers, No. 338, 2022.



- [4] J. Salminen et al., "Creating and detecting fake reviews of online products," *J. Retail. Consum. Serv.*, vol. 64, p. 102771, 2022.
- [5] M. T. Ribeiro, S. Singh, and C. Guestrin, "Why should I trust you? Explaining the predictions of any classifier," in *Proc. 22nd ACM SIGKDD Int. Conf. Knowl. Discov. Data Min.*, 2016, pp. 1135–1144.
- [6] S. He, B. Hollenbeck, and D. Proserpio, "The market for fake reviews," *Mark. Sci.*, vol. 41, no. 5, pp. 896–921, 2022.
- [7] B. Hooi et al., "BIRDNEST: Bayesian inference for ratings-fraud detection," in *Proc. SIAM Int. Conf. Data Min.*, 2016, pp. 495–503.
- [8] D. Mayzlin, Y. Dover, and J. Chevalier, "Promotional reviews: An empirical investigation of online review manipulation," *Amer. Econ. Rev.*, vol. 104, no. 8, pp. 2421–2455, 2014.
- [9] J. Devlin et al., "BERT: Pre-training of deep bidirectional transformers for language understanding," in *Proc. NAACL-HLT*, 2019, pp. 4171–4186.
- [10] H. Cai, L. Zhu, and S. Han, "ProxylessNAS: Direct neural architecture search on target task and hardware," in *Proc. ICLR*, 2019.
- [11] S. M. Lundberg and S.-I. Lee, "A unified approach to interpreting model predictions," in *Proc. NeurIPS*, 2017, pp. 4765–4774.
- [12] M. T. Ribeiro, S. Singh, and C. Guestrin, "LIME: Local interpretable model-agnostic explanations," in *Proc. 22nd ACM SIGKDD Int. Conf. Knowl. Discov. Data Min.*, 2016.
- [13] European Commission, "Digital Services Act: Ensuring a safe and accountable online environment," *Off. J. Eur. Union*, 2022.
- [14] B. Pang and L. Lee, "Opinion mining and sentiment analysis," *Foundations and Trends in Information Retrieval*, vol. 2, no. 1–2, pp. 1–135, 2008.
- [15] K. Shu, A. Sliva, S. Wang, J. Tang, and H. Liu, "Fake news detection on social media: A data mining perspective," *ACM SIGKDD Explorations*, vol. 19, no. 1, pp. 22–36, 2017.
- [16] N. Jindal and B. Liu, "Opinion spam and analysis," in *Proc. Int. Conf. Web Search and Data Mining (WSDM)*, Barcelona, Spain, 2008, pp. 219–230.
- [17] C. Molnar, *Interpretable Machine Learning*, 2nd ed., Munich, Germany: Leanpub, 2022.
- [18] ISO/IEC, "Information technology—Artificial intelligence—Trustworthiness in AI systems," ISO/IEC TR 24028, 2020.

BIOGRAPHY



Shraddha Shirish Shah is a student researcher in the Department of Information Technology at Thakur College of Engineering and Technology (TCET), Mumbai, India. Her research interests lie at the intersection of data mining and machine learning, with a focus on developing intelligent systems to enhance digital trust. She is a co-author of *Trustify: Building Consumer Trust with Intelligent Product Review Verification*.



Dr. Anil Vasoya is an Associate Professor in the Department of Information Technology at Thakur College of Engineering and Technology, Mumbai, India. He holds a Ph.D. in Computer Science and Technology, an M.E. in Computer Engineering, and a B.E. in Information Technology. With over 20 years of teaching experience, including 19 years at TCET, his areas of specialization include Database Management Systems, Distributed Operating Systems, Microprocessors, Data Mining, Web Programming, and Software Engineering. He has guided over 40 undergraduate and 3 postgraduate projects, published 16 papers in international conferences and 26 in international journals, and consistently receives high performance scores (average 9.20 on 360-Degree Feedback and 1246 on PRDP). He is a co-author of *Trustify: Building Consumer Trust with Intelligent Product Review Verification*.



Mrs. Pranjali Kasture is an Assistant Professor and the Deputy Head of the Department (Dy. HOD) of Information Technology at Thakur College of Engineering and Technology, Mumbai, India. She is currently pursuing her Ph.D. and holds an M.E. in Computer Engineering and a B.E. in Computer Engineering. With over 22 years of teaching experience, including 19 years at TCET, her research specializations include Data Mining, Machine Learning, and Deep Learning. She has guided over 30 undergraduate projects, published 17 papers in international conferences and 5 in international journals, and serves as a co-author of *Trustify: Building Consumer Trust with Intelligent Product Review Verification*.