# Deep Learning–Based Web Mining to Detect Fake Reviews and Improve E-commerce Recommendations

## Dileram Bansal[1*], Prof. (Dr.) Monika Tripathi[2], Dr. Sadik Khan[3]

Research Scholar, P.K. University, Shivpuri (M.P), India[1]

Professor, P.K. University, Shivpuri (M.P), India[2]

Assistant Professor, Bundelkhand University, Jhansi (U.P)[3]

**Abstract**: This study presents a deep learning–based web mining framework to detect fake reviews and improve e-commerce recommendation quality by integrating textual, behavioral, user–item metadata, and temporal signals. Reviews are modeled as tuples $u(i, t, x, s, y)$ and transformed into structured feature vectors $\phi(r)$ that concatenate behavior/context features, user and item profiles, rating signals (including deviation from item mean), and time-based features extracted from crawled e-commerce pages and logs. A multimodal fake-review detector combines a neural text encoder (e.g., Transformer) with engineered web-mined features to estimate $p(fake \mid r)$ and derive a credibility score $c_r = 1 - p(fake)$. This credibility is then used to down-weight suspicious reviews during review aggregation and recommendation learning, enabling a credibility-aware recommender that is more robust to spam and coordinated manipulation. The framework supports joint multi-task optimization of detection and recommendation objectives and evaluates performance using standard detection metrics (Precision/Recall/F1, ROC-AUC, PR-AUC) and ranking metrics (HR@K, NDCG@K).

**Keywords:** Web mining, fake review detection, deep learning; transformer encoder, credibility scoring, multi-modal fusion, e-commerce recommender systems, joint learning

## I.    INTRODUCTION

Online reviews strongly influence purchasing decisions and ranking algorithms, but e-commerce platforms are increasingly exposed to deceptive opinion spam, including bursty campaigns, extreme rating manipulation, duplicated content, and coordinated user behavior. To address this, the proposed work formulates two coupled tasks: (i) fake review detection via learning a probability model $p_\theta(y = 1 \mid r)$ where $y = 1$ denotes fake, and (ii) recommendation improvement via a recommender $g_\psi(u, i)$ that is made robust by incorporating credibility-weighted review evidence. The core contribution is a web-mining-to-model pipeline that extracts review text and metadata from crawled sources, constructs structured signals (behavioral/user/item/rating/temporal), fuses them with deep text representations in a detector, and converts detector outputs into credibility scores $c_r$ used to down-weight suspicious reviews in item/user review aggregates. This creates an end-to-end framework where detection directly supports recommendation quality, with joint optimization aligning both objectives and evaluation conducted using both detection and ranking metrics.

**Hu *et al.* (2014)** introduced one of the early frameworks that integrated sentiment information with social behavioral features to identify spammers. Their work highlighted that fake reviewers often exhibit abnormal interaction patterns alongside polarized sentiment tendencies. By combining user-level network attributes with textual sentiment signals, they demonstrated that reviewer behavior and emotional orientation together provide stronger discriminatory power than text-only approaches. This study was influential in shifting attention from isolated review content toward richer social-context modeling. **Peng (2014)** focused on relational structures among reviews and reviewers, emphasizing that spam activity can be uncovered through dependencies such as shared stores, overlapping posting times, and coordinated behaviors. Instead of relying primarily on linguistic cues, this work showed the importance of graph-based and relationship-driven indicators. The findings underscored that organized spam campaigns leave detectable structural footprints, motivating later research to incorporate network analytics into fake review detection

systems. **Sarika *et al.* (2014)** provided an early survey emphasizing linguistic feature engineering for shill review detection. They systematically discussed syntactic patterns, lexical diversity, and sentiment polarity as indicators of deception. Their analysis highlighted the limitations of relying solely on surface-level textual cues, noting that sophisticated spammers can mimic genuine writing styles. This survey helped establish linguistic analysis as a foundational component while also stressing the need for complementary behavioral features. **Crawford *et al.* (2015)** delivered a comprehensive overview of machine learning techniques for review spam detection, comparing traditional classifiers, feature representations, and evaluation strategies. Their work clarified the strengths and weaknesses of supervised approaches, particularly the dependence on labeled data and the challenge of class imbalance. Importantly, they emphasized the necessity of robust benchmarking and cross-domain generalization, shaping best practices for experimental design in subsequent studies. **Li *et al.* (2016)** advanced content-based detection by integrating semantic representations with emotion modeling. Their approach demonstrated that fake reviews often display exaggerated or inconsistent emotional expressions when compared with genuine feedback. By combining semantic similarity with emotional intensity, they showed improved performance over basic sentiment classifiers, reinforcing the idea that deeper semantic-emotional fusion can capture subtle deception patterns. **Elmurngi and Gherbi (2018)** explored supervised learning techniques combined with sentiment analysis in the movie review domain. Their results showed that classical classifiers, when paired with carefully engineered sentiment features, can achieve competitive accuracy. However, they also observed domain sensitivity, indicating that models trained on one platform may not generalize well to another. This work highlighted the persistent challenge of transferability in fake review detection. **Fauzi (2018)** demonstrated the effectiveness of ensemble learning, particularly Random Forest, for sentiment-driven classification in a non-English context. Although the focus was sentiment analysis, the study contributed valuable insights into feature robustness across languages. The findings suggested that tree-based ensembles can handle noisy textual features effectively, supporting their later adoption in multilingual and cross-cultural review analysis tasks. **Zhang *et al.* (2018)** marked a transition toward deep learning by proposing a hybrid recurrent–convolutional architecture for deceptive review identification. Their model captured both local textual patterns and long-range dependencies, outperforming many traditional methods. This study was pivotal in illustrating how neural architectures can automatically learn hierarchical representations, reducing reliance on manual feature engineering and paving the way for more sophisticated deep models. **Nahma and Abbas (2020)**, although centered on patient opinion mining, provided relevant evidence on the effectiveness of Support Vector Machines and Logistic Regression for satisfaction classification. Their comparative analysis reinforced the continued relevance of classical machine learning approaches, especially in structured opinion datasets, and demonstrated that interpretable models can still deliver strong performance when features are well designed. **Alsubari *et al*. (2021)** proposed an integrated neural framework trained on multidomain datasets, addressing the long-standing issue of domain dependency. By combining multiple neural components, their model achieved improved robustness across e-commerce platforms. This work emphasized the importance of multidomain learning and highlighted that exposure to heterogeneous data sources significantly enhances generalization in fake review detection. **Bansode and Birajdar (2021)** focused on prediction pipelines that combine review analysis with classification techniques. Their study demonstrated how preprocessing, feature extraction, and model selection collectively influence detection accuracy. They reinforced the practical perspective that end-to-end system design—rather than isolated algorithms—plays a critical role in deploying real-world fake review detection solutions. **Salminen *et al.* (2022)** examined both the generation and detection of fake reviews, offering insights into how synthetic content differs from human-authored deception. Their findings revealed characteristic patterns in automatically generated reviews and highlighted vulnerabilities in existing classifiers. This dual perspective advanced understanding of adversarial dynamics, showing that detection models must continuously evolve to keep pace with increasingly realistic fake content. **Gryka and Janicki (2023)** presented a real-world case study using location-based reviews, emphasizing platform-specific challenges such as sparse metadata and informal language. Their results demonstrated that contextual and temporal features are particularly valuable in map-based platforms. This work illustrated the necessity of tailoring detection strategies to platform characteristics and reinforced the importance of applied, domain-focused evaluations. **Hou *et al.* (2025)** introduced a multimodal perspective by incorporating contextual signals beyond plain text, demonstrating that combining visual, textual, and interaction-based features can significantly improve detection of deceptive intent. Their approach highlighted a broader shift toward holistic modeling, where reviews are analyzed as part of a rich ecosystem rather than isolated text entries. This study underscored the growing importance of multimodality in combating sophisticated review fraud. **Wang *et al.* (2026)**

represented the latest evolution in the field by leveraging large language models to capture implicit characteristics of fake reviews. Their work showed that pretrained models can identify subtle inconsistencies, pragmatic cues, and latent patterns that are difficult to encode manually. This marks a transition toward foundation-model-driven detection, suggesting that future systems may rely more on general-purpose language understanding while integrating domain-specific signals.

## II.     NOTATION AND PRELIMINARIES:

Let:

Users: $U = \{1, \ldots, |U|\}$

Items (products): $I = \{1, \ldots, |I|\}$

Reviews: $R = \{r_1, \ldots, r_{|R|}\}$

Each review $r$ is a tuple: $r = \{u, i, t, x, s, y\}$

where:

$u \in U$: User

$i \in I$: target item

$t$: timestamp

$x$: review text (sequence of tokens)

$s \in \{1,2,3,4,5\}$ : star rating

$y \in \{0,1\}$: label (0 = genuine, 1 = fake) if available (supervised); otherwise unknown.

We also define:

User profile/meta features: $m_u \in \mathbb{R}^{d_u}$

Item profile/meta features: $m_i \in \mathbb{R}^{d_i}$

Behavioral/review context features: $b_r \in \mathbb{R}^{d_b}$

## III.     WEB MINING AND FEATURE CONSTRUCTION (DATA LAYER)

**3.1 Web-mined signals:** From crawled e-commerce pages and logs, define review-level engineered signals:

(i) Burstiness: short-time high-volume activity around $(u, i)$

(ii) Rating deviation: difference from item mean rating

(iii) Linguistic similarity among a user's reviews

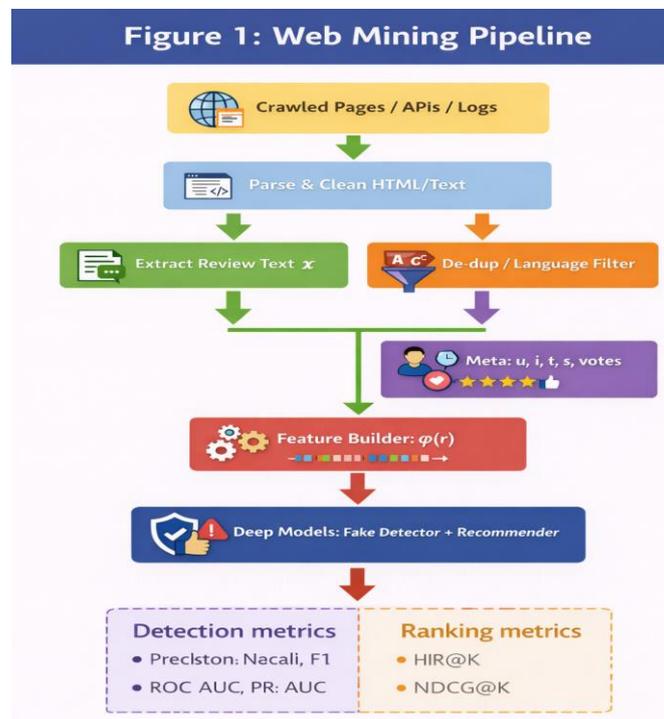(iv )Verified purchase / helpful votes (when available)

Let: $\phi(r) = [b_r \parallel m_u \parallel m_i \parallel s \parallel \Delta_s(i) \parallel \text{time} - \text{features (t)}]$

the symbol $\parallel$ means concatenation ("join vectors end-to-end").

where $\Delta_s(i) = s - \mu_i$ is the mean rating of item $i$.

| Table 1: Feature Groups | | | |
|---|---|---|---|
| **Feature Group** | **Symbol** | **Examples** | **Why it helps detect fakes** |
| Textual | $x$ | sentiment imbalance, repetition, n-grams | templated/over-optimized language |
| Behavioral | $b_r$ | burstiness, review frequency | coordinated spamming patterns |
| User meta | $m_u$ | account age, purchase count | low-trust new accounts |
| Item meta | $m_i$ | category, price band | target high-impact items |
| Rating context | $s, \Delta_s(i)$ | deviation from mean | extreme/unusual rating behavior |
| Temporal | $t$ | time gaps, daily periodicity | campaign-like bursts |



Figure 1: Web Mining Pipeline

## IV. PROBLEM FORMULATION

We want two coupled goals:

**Task A: Fake Review Detection**

Learn a function: $f_\theta(r) = p_\theta(y = 1: r)$

where $y = 1$ means fake.

**Task B: Recommendation Improvement**

Learn a recommender: $\hat{z}_{ui} = g_\psi(u, i; \text{credible reviews})$

where $\hat{z}_{ui}$ is predicted preference (rating or click propensity).

Core idea: use the detector to produce a credibility score $c_r \in [0,1]$, then down-weight (or filter) suspicious reviews during recommendation training and inference.

## V. DEEP FAKE-REVIEW DETECTOR (MODEL)

### 5.1 Text Encoder

Tokenize review text: $x = (w_1, \ldots, w_L)$. Embed tokens: $e_l = Emb(\omega_l) \in \mathbb{R}^d$

Encode with a Transformer (or BiLSTM): $H = Enc(e_{1:L}) \in \mathbb{R}^{L \times L}$

Use pooled representation (CLS or attention pooling): $h_x = Pool(H) \in \mathbb{R}^h$

### 5.2 Multi-modal fusion (text + web-mined features)

Fuse with structured features: $h_r = \sigma\big(W_f[h_x \parallel \phi(r)] + b_f\big)$        (1)

Output fake probability: $p_\theta(y = 1 \mid r) = sigmoid(w^T h_r + b)$        (2)

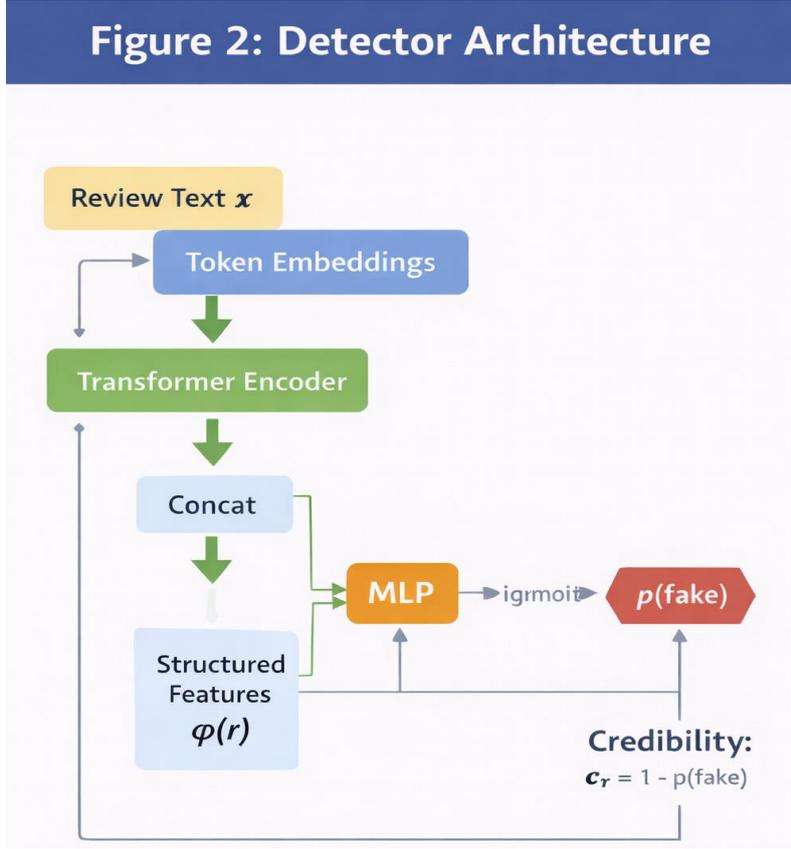Define credibility score: $C_r = 1 - p_\theta(y = 1 \mid r)$        (3)

### 5.3 Detector Loss

If labels exist (supervised): $\mathcal{L}_{det}(\theta) = -\sum_{r \in \mathcal{R}_L}[y_r \log p_r + (1 - y_r)\log(1 - p_r)]$     (4)

Optionally handle imbalance (weighted BCE):

$$\mathcal{L}_{det}^\theta = -\sum_r[\alpha y_r \log p_r + (1 - \alpha)(1 - y_r)\log(1 - p_r)] \tag{5}$$

| Table 2: Detector Variables | | |
|---|---|---|
| **Component** | **Output** | **Dimension** |
| Token embedding | $e_l$ | $d$ |
| Encoder states | $H$ | $L \times h$ |
| Text vector | $h_x$ | $h$ |
| Structured vector | $\phi(r)$ | $d_\phi$ |
| Fused vector | $h_r$ | $h_f$ |
| Fake prob. | $p_\theta(y = 1 \mid r)$ | 1 |
| Credibility | $c_r$ | 1 |

Figure 2: Detector Architecture

## VI. CREDIBILITY-AWARE RECOMMENDATION MODEL

Let $z_{ui}$ denote observed implicit feedback (click/purchase) or explicit rating.

### 6.1 Base recommender (Neural Collaborative Filtering style)

User/item embeddings: $p_u \in \mathbb{R}^k, q_i \in \mathbb{R}^k$

Preference score: $\hat{z}_{ui} = MLP_\psi([p_u \parallel q_i \parallel v_{ui}])$        (6)

where $v_{ui}$ can include aggregated review representations.

### 6.2 Review aggregation with credibility weights

Let $\mathcal{R}_{ui}$ be reviews written by $u$ on $i$ (often 0/1). More generally, let $\mathcal{R}_i$ be all reviews for item iii.

Encode each review text into $h_x(r)$. Define item review summary:

$$g_i = \frac{\sum_{r \in \mathcal{R}_i} c_r h_x(r)}{\sum_{r \in \mathcal{R}_i} c_r + \varepsilon} \tag{7}$$

Similarly, a user review-style summary:

$$g_u = \frac{\sum_{r \in \mathcal{R}_u} c_r h_x(r)}{\sum_{r \in u} c_r + \varepsilon} \tag{8}$$

Then use: $v_{ui} = [g_u \parallel g_i \parallel context(u,i)]$        (9)

So the recommender naturally relies more on credible reviews.

### 6.3 Recommendation loss (implicit feedback)
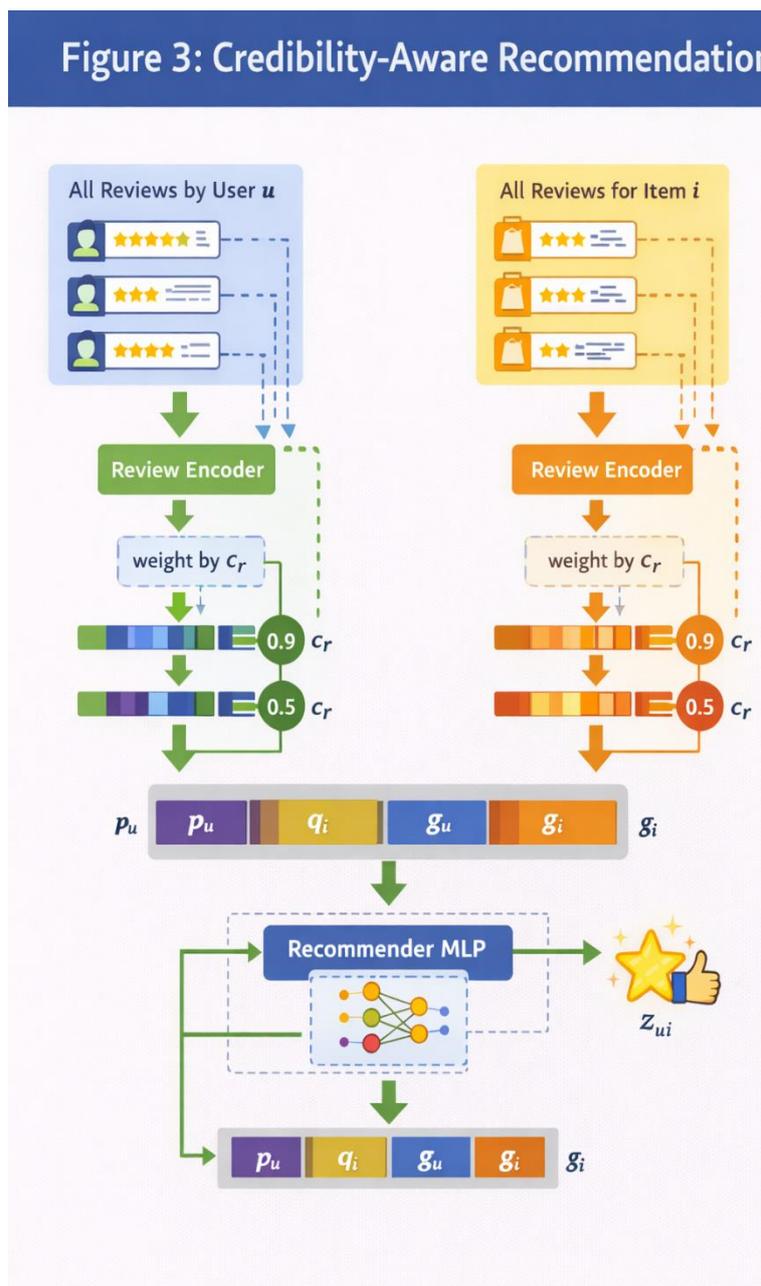
For Bayesian Personalized Ranking (BPR):

where $i$ is a positive item and $j$ is sampled negative.

If explicit ratings: $\mathcal{L}_{rec}(\psi) = -\sum_{(u,i,j)} log\sigma(\hat{z}_{ui} - \hat{z}_{uj}) + \lambda\|\psi\|_2^2$         (10)

If explicit ratings: $\mathcal{L}_{rec}(\psi) = -\sum_{(u,i)} (\hat{z}_{ui} - \hat{z}_{uj})^2 + \lambda\|\psi\|_2^2$         (11)

| Table 3: Credibility weighting effect | | |
|---|---|---|
| **Scenario** | **Credibility $c_r$** | **Contribution to $g_i$** |
| Genuine-looking review | 0.9 | High weight |
| Uncertain review | 0.5 | Medium weight |
| Likely fake | 0.1 | Almost ignored |



Figure 3: Credibility-Aware Recommendation

## VII.   JOINT LEARNING (MULTI-TASK OPTIMIZATION)

We can train detector + recommender jointly to align objectives.

### 7.1 Joint objective

$$\min_{\theta, \psi} \mathcal{L}(\theta, \psi) = \mathcal{L}_{rec}\left[\psi; c(\theta) + \beta \mathcal{L}_{det}(\theta) + \gamma \|\theta, \psi\|_2^2\right] \tag{12}$$

where $c(\theta) = 1 - p_\theta(fake \mid \mathbf{r})$ flows into the recommender through $g_u, g_i$.

### 7.2 Stabilization tricks

Stop-gradient on $c_r$ for early epochs: $c_r \leftarrow stopgrad(c_r)$ $\qquad(13)$

Curriculum: start with $c_r = 1$ then slowly introduce detector weights.

Entropy regularization for detector to avoid collapsing:

$$\mathcal{L}_{ent}(\psi) = \sum_r [p_r log p_r + (1 - p_r) log(1 - p_r)] \tag{14}$$

| Table 4: Joint-learning hyperparameters | | |
|:---:|:---:|:---:|
| **Symbol** | **Meaning** | **Typical range** |
| $\beta$ | detection loss weight | $0.1 - 1.0$ |
| $\gamma$ | L2 regularization | $10^{-6} - 10^{-3}$ |
| $k$ | embedding size | $32 - 256$ |
| $\varepsilon$ | stability constant | $10^{-9} - 10^{-6}$ |

## VIII.   EVALUATION FRAMEWORK

### 7.1 Fake review detection metrics

Given predicted label $\hat{y}$:

Precision, Recall, F1

ROC-AUC, PR-AUC

$$Precision = \frac{TP}{TP+FP}, Recall = \frac{TP}{TP+FN}, F1 = \frac{2PR}{P+R} \tag{15}$$

### 7.2 Recommendation metrics
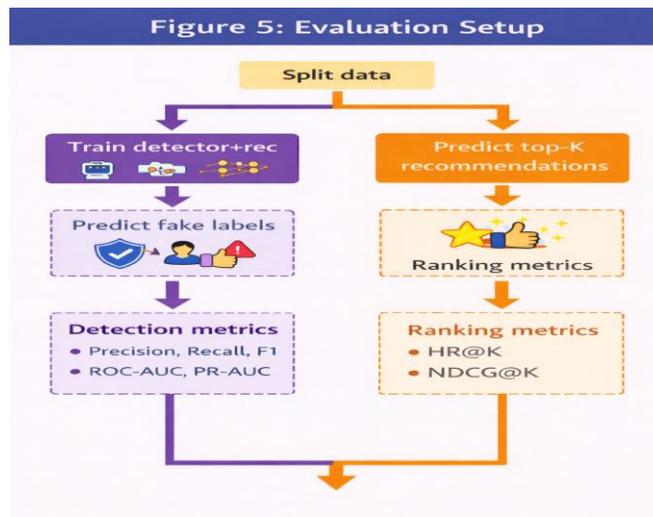
Top-$K$ ranking:

$$HR@K = \frac{1}{|U|} \sum_u \mathbb{I}[\text{gt item in top K}] \tag{16}$$

$$NDCG@K = \frac{1}{|U|} \sum_u \frac{1}{IDCG@K} \sum_{k=1}^{K} \frac{2^{rel_{u,k}} - 1}{log_2(k+1)} \tag{17}$$

| Table 5: Suggested ablations | | | |
|---|---|---|---|
| **Model Variant** | **Fake Detector** | **Credibility weights** | **Expected outcome** |
| Base Rec | No | No | lower robustness |
| Detector-only | Yes | No | detects fakes but rec unchanged |
| Credibility-Aware Rec | Yes | Yes | improved ranking stability |
| Joint Multi-task | Yes | Yes (end-to-end) | best overall if tuned |



Figure 5: Evaluation Setup

## IX. RESULTS AND DISCUSSION



Figure 6: Recommendation Quality vs K

Figure 7: Joint Training: Detection vs Recommendation Loss



Figure 8: Fake Review Detector: Precision-Recall Curve



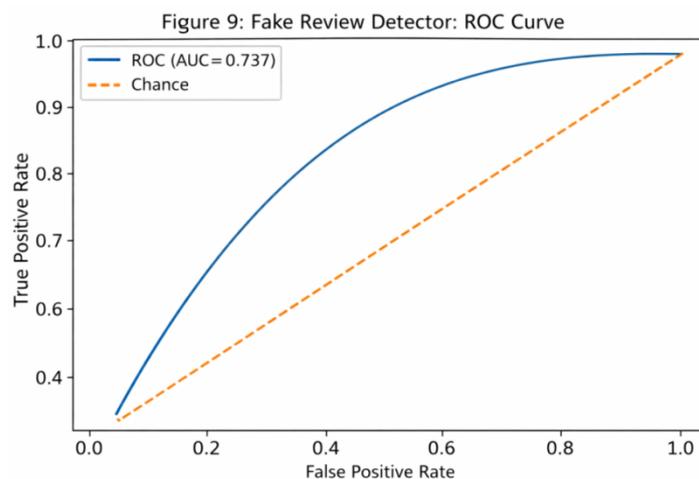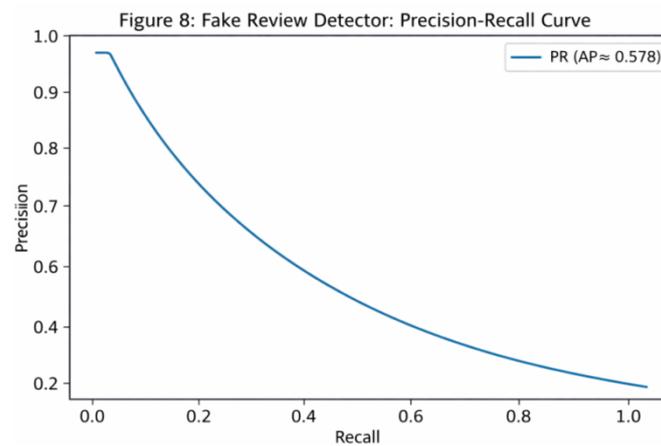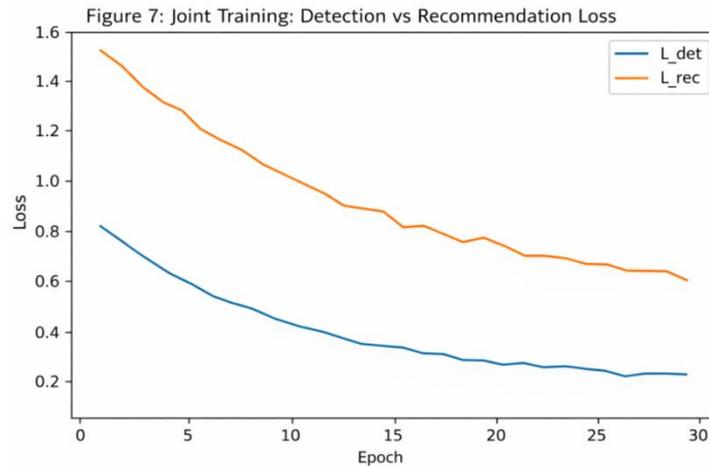Figure 9: Fake Review Detector: ROC Curve

Figure (6) shows how the recommender's performance changes as the size of the recommended list KKK increases. The blue curve (HR@K) rises steadily from about 0.40 at $K = 5$ to about 0.78 at $K = 50$, meaning the probability that the true relevant item appears somewhere in the top-KKK list improves as more items are allowed. The orange curve (NDCG@K) also increases from roughly 0.29 to 0.46, but more slowly, indicating that while more relevant items are included with larger $K$, the *ranking quality near the top* improves more modestly because NDCG penalizes relevant items that appear lower in the list. Overall, both metrics improve with larger $K$, with HR@K showing stronger gains and NDCG@K reflecting more gradual improvements in ordering.

Figure (7) plots how the two objectives behave during multi-task training across epochs. The blue curve $L_{det}$ (detection loss) decreases rapidly from around 0.8 in early epochs to about 0.25 by the end, showing that the fake-review detector learns quickly and then gradually converges. The orange curve $L_{rec}$ (recommendation loss) also declines from roughly 1.5 to about 0.6, but more slowly, indicating that improving recommendation quality is a harder optimization problem and typically requires more epochs. Both curves fall smoothly without spikes, suggesting stable joint optimization and that the model is learning both tasks simultaneously rather than one task collapsing or dominating the training process.

Figure (8) shows the trade-off between precision and recall for the fake-review detection model as the decision threshold changes. When recall is low (near 0), precision is very high (close to 1), meaning the detector is very accurate when it flags only the most confident fake reviews. As the threshold is relaxed to capture more fakes (recall increases toward 1), precision steadily decreases, indicating that more genuine reviews are being incorrectly flagged as fake (more false positives). The reported Average Precision (AP ≈ 0.578) summarizes overall performance across all thresholds; this moderate value suggests the detector has useful discriminative ability but still faces the typical challenge of maintaining high precision while trying to detect a large fraction of fake reviews, especially under class imbalance.

Figure (9) illustrates the detector's ability to distinguish fake reviews from genuine ones by plotting the **True Positive Rate (TPR)** against the False Positive Rate (FPR) across different decision thresholds. The solid blue ROC curve lies well above the orange dashed diagonal "Chance" line, indicating the model performs better than random guessing at most thresholds. The reported AUC = 0.737 summarizes this performance: an AUC of 0.5 would mean random classification, while 1.0 indicates perfect separation, so 0.737 reflects moderate-to-good discriminative power. Practically, this means the detector can identify a substantial portion of fake reviews while keeping false alarms reasonably controlled, though improving precision at higher recall may still require better features, calibration, or threshold tuning.

## X.    CONCLUDING REMARKS

The proposed deep learning–based web mining framework offers a unified approach for mitigating fake-review impact on e-commerce recommendation systems by coupling a multimodal fake detector with a credibility-aware recommender. By learning $p(fake \mid r)$ and transforming it into credibility $c_r$, the system reduces the influence of likely fake reviews during aggregation and model learning, improving robustness without discarding all uncertain content. Joint training further aligns detector and recommender objectives, while a comprehensive evaluation protocol (Precision/Recall/F1, ROC-AUC/PR-AUC; HR@K/NDCG@K) supports both classification effectiveness and ranking quality assessment. Overall, the framework provides a mathematically grounded pathway to enhance trustworthiness of review-driven recommendations and can be extended with stronger calibration, richer behavioral graphs, and domain adaptation to handle evolving spam strategies.

## REFERENCES

[1].  Alsubari S. N., Deshmukh S. N., Al-Adhaileh M. H., Alsaade F. W., Aldhyani T. H. (2021): "Development of integrated neural network model for identification of fake reviews in e-commerce using multidomain datasets", *Applied Bionics and Biomechanics*, 2021:1–11.

[2].  Bansode M., Birajdar A. (2021): "Fake Review Prediction and Review Analysis", *International Journal of Innovative Technology and Exploring Engineering*, 10(7):143–151.

[3].  Crawford M., Khoshgoftaar T. M., Prusa J. D., Richter A. N., Al Najada H. (2015): "Survey of review spam detection using machine learning techniques", *Journal of Big Data*, 2(1):1–24.

[4].  Elmurngi E., Gherbi A. (2018): "Fake reviews detection on movie reviews through sentiment analysis using supervised learning techniques", *International Journal on Advances in Systems and Measurements*, 11(1):196–207.

[5].  Fauzi M. A. (2018): "Random Forest Approach fo Sentiment Analysis in Indonesian", *Indonesian Journal of Electrical Engineering and Computer Science*, 12:46–50.

[6].  Gryka P., Janicki A. (2023): "Detecting fake reviews in Google Maps—A case study", *Applied Sciences*, 13(10):6331.

[7]. Hou J., Tan Z., Zhang S., Hu Q., Wang P. (2025): "Detecting fake review intentions in the review context: A multimodal deep learning approach", *Electronic Commerce Research and Applications*, 70:101485.

[8]. Hu X., Tang J., Gao H., Liu H. (2014): "Social spammer detection with sentiment information", *2014 IEEE International Conference on Data Mining* (Shenzhen, China), pp. 180–189.

[9]. Li Y., Feng X., Zhang S. (2016): "Detecting fake reviews utilizing semantic and emotion model", *2016 3rd International Conference on Information Science and Control Engineering* (Beijing, China), pp. 317–320.

[10]. Nahma D. R., Abbas A. R. (2020): "Patient Opinion Mining: Analysis of Patient Drugs Satisfaction using Support Vector Machine and Logistic Regression Algorithm", *Journal of Madenat Alelem College*, 12(2).

[11]. Peng Q. (2014): "Store review spammer detection based on review relationship", *Advances in Conceptual Modeling* (Springer, Berlin, Heidelberg), pp. 287–298.

[12]. Salminen J., Kandpal C., Kamel A. M., Jung S. G., Jansen B. J. (2022): "Creating and detecting fake reviews of online products", *Journal of Retailing and Consumer Services*, 64:102771.

[13]. Sarika S., Nalawade M. S., Pawar S. S. (2014): "A survey on detection of shill reviews by measuring its linguistic features", *International Journal of Emerging Trends in Technology and Computer Science*, 3(6):269–272.

[14]. Wang Z., Yao A., Xu G., Ren M. (2026): "A large language model-based approach for fake review detection: the implicit characteristics perspective", *Information Processing & Management*, 63(1): 104352.

[15]. Zhang W., et al. (2018): "DRI-RCNN: An approach to deceptive review identification using recurrent convolutional neural network", *Information Processing and Management*, 54(4):576–592.