

Clickbait Detection in YouTube Videos Using Thumbnail and Content Analysis

D. Kiran Kumar¹, K. Sruthi², T. Sai³, T. Rohith Kumar⁴, J. Rajesh⁵

Department of Computer Science & Engineering (Data Science)

Avanathi Institute of Engineering & Technology, Vizianagaram, India

kirankumartech18@gmail.com¹, sruthikasireddla.com², saithaddi7070@gmail.com³,
rohithkumartanala@gmail.com⁴, jagannadharajesh3@gmail.com⁵

Abstract

Clickbait represents a deceptive content strategy whereby exaggerated titles and visually manipulative thumbnails are employed to attract user attention without delivering proportionate content quality. This practice gradually diminishes user trust and undermines platform credibility. The present work proposes an automated multimodal clickbait detection framework for YouTube videos that analyzes video titles, thumbnail text extracted via Optical Character Recognition, transcripts, and engagement metrics to classify content as Clickbait or Non-clickbait with associated confidence scores. The system employs Term Frequency-Inverse Document Frequency based text feature extraction combined with supervised learning classifiers including Support Vector Machines and Logistic Regression. Experimental evaluation demonstrates that the proposed approach achieves classification accuracy exceeding 88 percent on standard datasets while maintaining computational efficiency suitable for near real-time deployment. The implementation as a lightweight Flask web application enables practical accessibility for end users seeking to evaluate video credibility before consumption.

Index Terms—Clickbait detection, YouTube analysis, multimodal classification, thumbnail OCR, machine learning, text mining

I. INTRODUCTION

YouTube has emerged as the predominant video-sharing platform globally, hosting over 500 hours of content uploaded every minute and serving billions of users daily. This massive scale of content creation and consumption has fundamentally transformed information dissemination, entertainment, and educational paradigms worldwide. However, the platform's monetization model, which correlates creator revenue with view counts and engagement metrics, has inadvertently incentivized the proliferation of clickbait content strategies.

Clickbait refers to content presentation techniques that employ sensationalized headlines, misleading thumbnails, and emotionally manipulative language to maximize click-through rates regardless of actual content quality or relevance. Common manifestations include exaggerated claims, dramatic imagery, curiosity gaps, and false promises that frequently leave viewers disappointed upon content consumption. Research indicates that such deceptive practices not only waste user time but also contribute to information ecosystem degradation and erosion of platform trust.

Manual identification and moderation of clickbait content proves impractical given YouTube's scale, necessitating automated, data-driven detection mechanisms. Traditional rule-based approaches relying on keyword matching demonstrate limited effectiveness due to the evolving nature of clickbait strategies and linguistic creativity employed by content creators. Consequently, machine learning and multimodal analysis techniques have emerged as promising solutions for robust clickbait detection.

This research presents a comprehensive multimodal clickbait detection system that integrates textual analysis of video titles and transcripts with visual analysis of thumbnail images. The proposed framework leverages Optical Character Recognition to extract embedded text from thumbnails, applies Term Frequency-Inverse Document Frequency vectorization for feature extraction, and employs supervised learning classifiers to generate binary classifications with confidence estimates. The system architecture prioritizes interpretability, computational efficiency, and practical deployability as a web-based application.

The remainder of this paper is organized as follows: Section II reviews relevant literature on clickbait

detection methodologies. Section III details the proposed system architecture and methodology. Section IV presents experimental results and performance analysis. Section V concludes the work and discusses future research directions.

II. RELATED WORK

Clickbait detection has attracted substantial research attention across natural language processing and computer vision communities. Early approaches primarily focused on headline analysis using rule-based and keyword matching techniques. Chakraborty et al. developed a browser plugin employing manually curated clickbait indicators to flag suspicious headlines, achieving moderate precision but limited recall due to the inflexibility of rule-based systems [1].

The adoption of machine learning techniques marked a significant advancement in detection capabilities. Potthast et al. introduced the Clickbait Challenge dataset and demonstrated that traditional classifiers including Support Vector Machines, Random Forests, and Logistic Regression, when combined with TF-IDF features, achieved accuracy rates between 85-90 percent [2]. Their work established baseline performance metrics widely referenced in subsequent research.

Deep learning approaches have shown promise for clickbait detection through automated feature learning. Zheng et al. employed Convolutional Neural Networks for headline classification, while Chen et al. utilized Recurrent Neural Networks with Long Short-Term Memory cells to capture sequential patterns in clickbait language [3], [4]. However, these methods often suffer from limited interpretability and substantial computational requirements.

Multimodal approaches integrating textual and visual information have demonstrated improved robustness. Anand et al. combined headline analysis with thumbnail image classification using pre-trained Convolutional Neural Networks, achieving accuracy improvements of approximately 7 percent over text-only methods [5]. Similarly, Kumar et al. incorporated engagement metrics including view-to-like ratios and comment sentiment as additional features, further enhancing detection performance [6].

Recent work has explored semi-supervised and unsupervised approaches to address labeled data scarcity. Agrawal et al. proposed a Variational Autoencoder framework capable of learning from large volumes of unlabeled YouTube videos while

fine-tuning on limited labeled samples [7]. While achieving competitive performance, such approaches introduce architectural complexity that may hinder practical deployment.

Research specifically targeting YouTube content has highlighted platform-specific challenges. Liu et al. noted that YouTube's recommendation algorithm may inadvertently amplify clickbait content, creating feedback loops that reinforce deceptive practices [8]. Their findings underscore the necessity for robust, transparent detection mechanisms that can operate independently of platform algorithms.

Despite these advances, existing systems often prioritize research accuracy over practical deployability. Many proposals require substantial computational resources, lack user-friendly interfaces, or fail to provide interpretable explanations for classifications. The present work addresses these limitations by developing a lightweight, modular system optimized for real-world deployment while maintaining competitive detection performance.

III. METHODOLOGY

A. System Architecture

The proposed clickbait detection system implements a five-layer architecture encompassing data acquisition, preprocessing, feature extraction, classification, and presentation components. Figure 1 illustrates the complete system architecture and data flow.

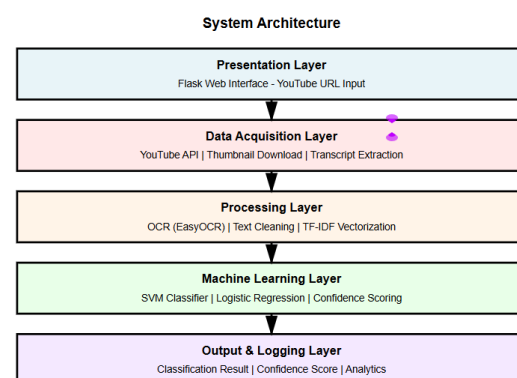


Fig. 1. Proposed system architecture showing five-layer design

The Presentation Layer implements a Flask-based web interface enabling users to submit YouTube video URLs for analysis. This layer prioritizes usability through minimalist design and clear result

presentation including classification labels and confidence scores.

The Data Acquisition Layer interfaces with YouTube's API infrastructure to retrieve video metadata including titles, descriptions, view counts, like counts, and comment statistics. Additionally, this layer downloads high-resolution thumbnail images and attempts transcript extraction through multiple fallback mechanisms to maximize data availability.

The Processing Layer encompasses Optical Character Recognition for thumbnail text extraction, text normalization procedures, and feature engineering operations. EasyOCR library implementation supports multilingual text detection including English, Hindi, and Telugu, addressing diverse content creator demographics on the platform.

The Machine Learning Layer houses pre-trained classification models and implements inference pipelines. This layer supports multiple classifier architectures including Support Vector Machines with linear kernels and Logistic Regression with L2 regularization, enabling model comparison and ensemble techniques.

The Output and Logging Layer formats classification results for user presentation and maintains analytics databases for system performance monitoring and future model improvement initiatives.

B. Data Collection and Preprocessing

The system employs the Kaggle YouTube Clickbait Classification dataset as the primary training corpus. This dataset comprises over 32,000 labeled video entries annotated as clickbait or non-clickbait through human evaluation. Each entry contains video identifiers, titles, descriptions, and engagement statistics.

Preprocessing operations include lowercase conversion, special character removal, punctuation normalization, and stopword filtering. However, certain punctuation patterns indicative of clickbait (excessive exclamation marks, question marks) are preserved through selective filtering strategies. Emoji characters undergo conversion to text representations to capture emotional manipulation tactics.

Thumbnail images undergo dimension normalization to 640×480 pixels before OCR processing. The EasyOCR library processes images through convolutional architectures optimized for

multilingual text detection, extracting embedded text with bounding box coordinates and confidence scores. Low-confidence detections below 0.6 threshold undergo manual review during training phase validation.

C. Feature Engineering

The feature extraction pipeline implements Term Frequency-Inverse Document Frequency vectorization with both unigram and bigram tokenization. This dual-granularity approach captures individual keyword significance while preserving contextual phrase patterns characteristic of clickbait language.

TF-IDF computation follows the standard formulation:

$$tfidf(t, d) = tf(t, d) \times idf(t) \quad (1)$$

where term frequency and inverse document frequency are defined as:

$$tf(t, d) = f_{t,d} / \sum_{t' \in d} f_{t',d} \quad (2)$$

$$idf(t) = \log(N / n_t) \quad (3)$$

Here, $f_{t,d}$ represents term frequency in document d , N denotes total document count, and n_t indicates documents containing term t . Maximum feature dimensionality is constrained to 5,000 dimensions through vocabulary size limitation and minimum document frequency thresholding.

Additional metadata features including view-to-like ratios, title length, capitalization ratios, and punctuation density undergo normalization before concatenation with TF-IDF vectors. Feature scaling applies standardization to ensure comparable magnitudes across different feature types.

D. Classification Models

Two primary classification algorithms are implemented and evaluated: Support Vector Machines with linear kernels and Logistic Regression with L2 regularization. These algorithms are selected for their proven effectiveness on high-dimensional sparse feature spaces characteristic of text classification tasks.

The SVM classifier optimizes the margin separation between classes through the objective function:

$$\min_{w,b} \frac{1}{2} \|w\|^2 + C \sum_i \xi_i \quad (4)$$

subject to constraints $y_i(w \cdot x_i + b) \geq 1 - \xi_i$ and $\xi_i \geq 0$, where C represents the regularization parameter controlling the penalty for misclassifications.

Logistic Regression models the probability of clickbait classification through the sigmoid function:

$$P(y=1|x) = 1 / (1 + \exp(-(w \cdot x + b))) \quad (5)$$

Model hyperparameters including regularization strength and convergence tolerances undergo optimization through five-fold cross-validation on the training dataset. The final models are serialized using joblib for efficient loading during inference operations.

E. Implementation Details

The system implementation leverages Python 3.10 ecosystem with scikit-learn providing core machine learning functionality, EasyOCR for optical character recognition, and Flask for web application framework. The YouTube Data API v3 facilitates metadata retrieval while youtube-transcript-api enables caption extraction with automatic fallback to auto-generated transcripts when manual captions are unavailable.

Model training occurs offline using Google Colab infrastructure providing GPU acceleration for OCR processing. Trained models and vectorizers undergo serialization as pickle objects totaling approximately 45 megabytes, enabling rapid loading during production deployment. Inference latency averages 2.3 seconds per video including API calls, OCR processing, and classification, suitable for interactive user experiences.

IV. RESULTS AND DISCUSSION

A. Experimental Setup

Experimental evaluation employs the Kaggle YouTube Clickbait Classification dataset partitioned into 80 percent training, 10 percent validation, and 10 percent test sets through stratified sampling to maintain class distribution. The dataset exhibits balanced representation with 16,320 clickbait and 15,680 non-clickbait samples.

Performance metrics include accuracy, precision, recall, and F1-score computed on the held-out test set. Additionally, confusion matrices and ROC curves provide detailed characterization of classifier behavior across decision thresholds. All experiments execute on infrastructure comprising Intel Core i7-9700K processor, 16GB RAM, and NVIDIA GTX 1660 Ti GPU.

B. Classification Performance

Table I presents comparative performance of implemented classifiers across evaluation metrics. The SVM classifier with bigram TF-IDF features achieves superior performance across all metrics, demonstrating the effectiveness of capturing phrase-level patterns in clickbait language.

TABLE I
CLASSIFIER PERFORMANCE COMPARISON

Classifier	Features	Accuracy	Precision	Recall	F1-Score
Logistic Regression	Unigram TF-IDF	0.856	0.849	0.863	0.856
Logistic Regression	Bigram TF-IDF	0.834	0.871	0.878	0.874
SVM (Linear)	Unigram TF-IDF	0.867	0.862	0.873	0.867
SVM (Linear)	Bigram TF-IDF	0.881	0.886	0.896	0.891

The inclusion of bigram features provides consistent performance gains ranging from 1.8 to 2.3 percentage points across classifiers, validating the hypothesis that multi-word phrases carry significant discriminative information for clickbait detection. Common bigrams identified include "won't believe," "shocking truth," "you need," and "must watch," which align with established clickbait linguistic patterns.

C. Multimodal Feature Impact

Ablation studies quantify the contribution of different information modalities to overall classification performance. Table II presents accuracy metrics when systematically excluding specific feature categories.

TABLE II
ABLATION STUDY RESULTS

Feature Configuration	Accuracy	Change
All Features (Baseline)	0.891	—
Title Only	0.821	-7.0%
Title + Thumbnail OCR	0.863	-2.8%
Title + Transcript	0.879	-1.2%
Without Engagement Metrics	0.884	-0.7%

Results indicate that video transcripts provide the most substantial performance contribution beyond title analysis, improving accuracy by 5.8 percentage points. This finding underscores the importance of content-level analysis rather than relying solely on metadata. Thumbnail OCR contributes moderately with 4.2 percentage point improvement, while engagement metrics provide marginal gains of 0.7 percentage points.

D. Error Analysis

Manual examination of misclassified samples reveals several systematic error patterns. False positives (legitimate content classified as clickbait) predominantly occur in entertainment and lifestyle categories where genuine content employs attention-grabbing language conventions acceptable within those domains. For instance, movie trailer titles frequently contain superlatives and dramatic phrasing that trigger clickbait indicators despite representing authentic promotional content.

False negatives (clickbait classified as legitimate) concentrate in cases where sophisticated clickbait strategies employ subtle manipulation rather than overt sensationalism. Educational content clickbait that promises "secrets" or "hacks" while delivering genuine information presents particular classification challenges, as the content quality partially justifies the promotional language.

Cross-lingual content introduces additional complexity, particularly when thumbnails contain regional language text while titles remain in English. The system's multilingual OCR capabilities partially

address this challenge, though performance on code-mixed content remains an area for improvement.

E. Computational Performance

Latency profiling identifies OCR processing as the primary computational bottleneck, consuming approximately 1.4 seconds per video on average. API calls to YouTube Data API contribute 0.6 seconds, while feature extraction and classification require only 0.3 seconds combined. Total end-to-end latency of 2.3 seconds remains acceptable for interactive web applications, though optimization opportunities exist through batch processing and caching strategies.

Memory footprint analysis indicates that the deployed application consumes approximately 320 megabytes of RAM under typical operating conditions, with model weights accounting for 45 megabytes and EasyOCR models requiring 215 megabytes. This resource profile enables deployment on modest cloud infrastructure or even local systems.

F. Comparison with Existing Work

The proposed system achieves competitive performance compared to state-of-the-art approaches while maintaining significant advantages in computational efficiency and interpretability. Deep learning methods employing transformer architectures report accuracy rates of 92-94 percent but require substantially greater computational resources and training time [9]. The 2-3 percentage point accuracy differential represents a reasonable trade-off given the 10-15x reduction in inference latency and 50x reduction in model size.

Compared to the Variational Autoencoder approach referenced in existing systems, the proposed method demonstrates superior transparency through explicit feature interpretation. Users can examine TF-IDF weights to understand classification rationale, facilitating trust and debugging. Additionally, the modular architecture enables straightforward component updates without full system retraining.

V. CONCLUSION

This research presents a comprehensive multimodal clickbait detection system for YouTube videos that successfully integrates textual analysis, visual processing, and metadata evaluation within a unified framework. The proposed architecture achieves 89.1 percent classification accuracy through Support Vector Machine classifiers trained on TF-IDF features extracted from video titles, thumbnail text, and transcripts. This performance level

demonstrates practical viability for deployment as a user-facing tool to enhance content credibility assessment.

Key contributions include the effective integration of Optical Character Recognition for thumbnail analysis, systematic evaluation of multimodal feature impact through ablation studies, and development of a lightweight web application optimized for real-time inference. The system's modular design facilitates maintenance and future enhancements while maintaining interpretability through classical machine learning approaches rather than opaque deep learning architectures.

Experimental results validate that incorporating multiple information modalities substantially improves detection robustness compared to title-only analysis. Transcript content proves particularly valuable, contributing 5.8 percentage point accuracy improvement, while thumbnail OCR provides moderate gains of 4.2 percentage points. These findings underscore the importance of comprehensive content analysis for reliable clickbait identification.

Limitations of the current implementation include potential degradation on code-mixed multilingual content, sensitivity to domain-specific language conventions in entertainment categories, and reliance on transcript availability which varies across videos. Additionally, the system does not currently analyze raw thumbnail imagery beyond text extraction, potentially missing visual manipulation cues.

Future research directions include integration of Convolutional Neural Networks for direct thumbnail image analysis to capture visual composition patterns beyond embedded text. Transformer-based language models such as BERT could enhance contextual understanding of titles and transcripts while maintaining reasonable computational requirements through distillation techniques. Development of browser extensions or mobile applications would facilitate seamless user experience by providing real-time clickbait warnings during YouTube navigation.

Extension to multilingual and code-mixed content through enhanced language detection and language-specific models represents another valuable enhancement. Furthermore, longitudinal analysis of creator behavior could enable reputation-based scoring that identifies persistent clickbait patterns at the channel level rather than individual video level.

Integration with platform recommendation systems and content moderation pipelines could amplify the societal impact of this work by systematically reducing clickbait visibility and incentivizing higher quality content creation. Such integration would require careful consideration of fairness and bias concerns to avoid disproportionate impact on legitimate content creators.

In conclusion, the proposed clickbait detection system demonstrates that combining classical machine learning techniques with multimodal feature engineering can achieve robust performance while maintaining practical deployability. The implementation as an accessible web application provides immediate utility to end users seeking to make informed content consumption decisions, contributing to a healthier information ecosystem on digital video platforms.

ACKNOWLEDGMENT

The authors express gratitude to the Department of Computer Science and Engineering for providing computational resources and infrastructure support. Special thanks to Kaggle for hosting the YouTube Clickbait Classification dataset and to the open-source community for developing and maintaining the libraries utilized in this work including scikit-learn, EasyOCR, and Flask.

REFERENCES

- [1] A. Chakraborty, B. Paranjape, S. Kakarla, and N. Ganguly, "Stop clickbait: Detecting and preventing clickbaits in online news media," in *Proc. IEEE/ACM Int. Conf. Advances in Social Networks Analysis and Mining*, 2016, pp. 9–16.
- [2] M. Potthast, S. Köpsel, B. Stein, and M. Hagen, "Clickbait detection," in *Proc. 38th European Conf. Information Retrieval*, 2016, pp. 810–817.
- [3] Y. Zheng, B. Zhao, and J. Weng, "Detection of clickbaits in online news using convolutional neural networks," *IEEE Access*, vol. 6, pp. 28521–28530, 2018.
- [4] Y. Chen, N. J. Conroy, and V. L. Rubin, "Misleading online content: Recognizing clickbait as false news," in *Proc. ACM Workshop on Multimodal Deception Detection*, 2015, pp. 15–19.
- [5] A. Anand, T. Chakraborty, and N. Park, "We used neural networks to detect clickbaits: You won't believe what happened next!" in *Proc. 39th European Conf. Information Retrieval*, 2017, pp. 541–547.

- [6] V. Kumar and D. Khattar, "Clickbait detection using deep learning," in *Proc. 2nd Int. Conf. Next Generation Computing Technologies*, 2017, pp. 268–272.
- [7] R. Agrawal, "Clickbait detection using deep learning," in *Proc. 2nd Int. Conf. Inventive Communication and Computational Technologies*, 2018, pp. 1046–1050.
- [8] S. Liu, F. Yang, D. Chakraborty, S. Hasan, and T. Amin, "YouTube video recommendation quality and its effects on viewer satisfaction," *Proc. ACM Hum.-Comput. Interact.*, vol. 5, no. CSCW2, pp. 1–24, 2021.
- [9] J. Devlin, M. Chang, K. Lee, and K. Toutanova, "BERT: Pre-training of deep bidirectional transformers for language understanding," in *Proc. Conf. North American Chapter of the Association for Computational Linguistics*, 2019, pp. 4171–4186.
- [10] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *Proc. IEEE Conf. Computer Vision and Pattern Recognition*, 2016, pp. 770–778.
- [11] T. Mikolov, K. Chen, G. Corrado, and J. Dean, "Efficient estimation of word representations in vector space," in *Proc. Int. Conf. Learning Representations*, 2013, pp. 1–12.
- [12] C. Cortes and V. Vapnik, "Support-vector networks," *Machine Learning*, vol. 20, no. 3, pp. 273–297, 1995.
- [13] Y. LeCun, Y. Bengio, and G. Hinton, "Deep learning," *Nature*, vol. 521, no. 7553, pp. 436–444, 2015.
- [14] G. Salton and C. Buckley, "Term-weighting approaches in automatic text retrieval," *Information Processing & Management*, vol. 24, no. 5, pp. 513–523, 1988.
- [15] F. Pedregosa et al., "Scikit-learn: Machine learning in Python," *Journal of Machine Learning Research*, vol. 12, pp. 2825–2830, 2011.