

---

# Contextual Targeting for Video Advertisements: A Machine Learning Approach to Enhancing Relevance and Performance

Hrishikesh Desai

---

## Abstract

This research introduces an innovative method for contextual targeting in video advertising, leveraging advanced machine learning techniques. Our multi-modal model, which fuses visual, audio, and textual analysis, significantly surpasses traditional keyword-based approaches across various performance metrics. The 37% increase in click-through rates (CTR) and a 22% rise in view-through rates (VTR) highlight the practical advantages of our strategy for advertisers and content platforms. These results showcase the potential of machine learning to enhance the relevance and effectiveness of video ad placements while ensuring user privacy. Future studies could delve into real-time processing for live video streams and the inclusion of additional contextual signals, such as user engagement patterns. As the digital advertising landscape continues to change, our research lays a strong foundation for utilizing machine learning in contextual targeting for video ads.

Copyright © 2025 International Journals of Multidisciplinary Research Academy. All rights reserved.

---

## Keywords:

Contextual targeting  
Video advertising  
Machine learning  
Multi-modal analysis  
Performance optimization.

---

## Author correspondence:

Hrishikesh Desai,  
Email: hrishikeshd2@gmail.com

---

## 1. Introduction

The landscape of digital advertising is changing quickly, with video content consumption hitting new heights. As traditional targeting methods that rely on user data come under more scrutiny due to privacy issues, contextual targeting has made a strong comeback as an effective alternative. This paper introduces an innovative approach to contextual targeting for video ads, using advanced machine learning techniques to analyze and align video content with suitable advertisements. While earlier research has looked into contextual targeting in text-based settings, the intricate nature of video content brings its own set of challenges and opportunities. Our study aims to fill this gap by creating a multi-modal machine learning framework that processes visual, audio, and textual components of video content simultaneously to achieve highly accurate contextual matches. The method we propose seeks to address the shortcomings of keyword-based strategies, which often struggle to grasp the subtle context of video content. By employing deep learning models that can recognize complex patterns across various data types, we believe our approach will greatly enhance the relevance and effectiveness of video ad placements.

## 2. Research Method

### 2.1 Data Collection and Preprocessing

We compiled a diverse dataset of 1 million video-ad pairs from a major online video platform. The dataset included:

- Video content: Full-length videos (mean duration: 8.7 minutes, SD: 4.2 minutes)
- Associated advertisements: 15-30 second video ads
- Performance metrics: Click-through rates (CTR) and view-through rates (VTR) for each video-ad pair

Preprocessing steps included:

1. Frame extraction: Sampling video frames at 1 fps
2. Audio feature extraction: MFCC (Mel-frequency cepstral coefficients) with a 25ms window and 10ms step

3. Transcript generation: Using an automated speech recognition (ASR) system with a word error rate of 5.3%

## 2.2 Multi-Modal Machine Learning Architecture

Our contextual targeting system consists of three main components:

1. Visual Analysis Module:
  - Architecture: ResNet-152 CNN pretrained on ImageNet
  - Fine-tuning: Last two layers retrained on our video frame dataset
  - Output: 2048-dimensional feature vector for each frame
2. Audio Analysis Module:
  - Architecture: Bidirectional LSTM network
  - Input: MFCC features
  - Output: 512-dimensional audio context vector
3. Textual Analysis Module:
  - Architecture: BERT-base transformer model
  - Input: Video transcripts and ad descriptions
  - Output: 768-dimensional contextual embedding for text

The outputs from these three modules were concatenated and fed into a fully connected neural network for final contextual matching.

## 2.3 Training and Optimization

We used a two-step training approach:

1. Training individual modules: Each module was trained independently on its specific data type.
2. Fine-tuning the entire model: We fine-tuned the complete model with a contrastive loss function to enhance similarity between matching video-ad pairs while reducing similarity for non-matching pairs.

For hyperparameter optimization, we utilized Bayesian optimization with a Gaussian process prior. The main hyperparameters we focused on included the learning rate, batch size, and the temperature parameter in the contrastive loss function.

## 2.4 Evaluation Metrics

We evaluated our model using the following metrics:

1. Area Under the Receiver Operating Characteristic curve (AUROC)
2. Mean Average Precision (MAP)
3. Normalized Discounted Cumulative Gain (NDCG)
4. Click-Through Rate (CTR) improvement
5. View-Through Rate (VTR) improvement

## 3. Results and Analysis

Our multi-modal machine learning approach demonstrated significant improvements in contextual targeting performance compared to traditional keyword-based methods.

### 3.1 Model Performance

Table 1 presents the performance metrics of our model compared to the baseline keyword-based approach:

Table 1. Performance Comparison

Metric	Our Model	Baseline	Improvement
AUROC	0.92	0.78	+17.9%
MAP	0.87	0.71	+22.5%
NDCG	0.89	0.73	+21.9%

The substantial improvements across all metrics indicate that our model more accurately captures the contextual relevance between video content and advertisements.

### 3.2 Impact on Advertising Performance

We observed significant improvements in key advertising performance metrics:

1. Click-Through Rate (CTR):
  - Our model achieved a mean CTR of 3.8% (SD: 0.7%)
  - Baseline method: mean CTR of 2.8% (SD: 0.6%)
  - Improvement: 37% increase in CTR ( $p < 0.001$ , paired t-test)
2. View-Through Rate (VTR):
  - Our model: mean VTR of 72.3% (SD: 5.2%)
  - Baseline method: mean VTR of 59.3% (SD: 6.1%)
  - Improvement: 22% increase in VTR ( $p < 0.001$ , paired t-test)

### 3.3 Feature Importance Analysis

To grasp the importance of each modality, we carried out an ablation study where we eliminated each component and analyzed its impact on performance. The visual features were found to have the most substantial effect on the model's performance, accounting for 45%, while textual features contributed 35% and audio features 20%. This emphasizes the value of a multi-modal approach in thoroughly understanding the context of video content.

### 3.4 Temporal Analysis

We analyzed the model's performance across different video durations and found a positive correlation between video length and targeting accuracy (Pearson's  $r = 0.68$ ,  $p < 0.001$ ). This suggests that longer videos provide more contextual information for our model to leverage.

## 4. Conclusion

This research introduces an innovative method for contextual targeting in video advertising, leveraging advanced machine learning techniques. Our multi-modal model, which fuses visual, audio, and textual analysis, significantly surpasses traditional keyword-based approaches across various performance metrics. The 37% increase in click-through rates (CTR) and a 22% rise in view-through rates (VTR) highlight the practical advantages of our strategy for advertisers and content platforms. These results showcase the potential of machine learning to enhance the relevance and effectiveness of video ad placements while ensuring user privacy. Future studies could delve into real-time processing for live video streams and the inclusion of additional contextual signals, such as user engagement patterns. As the digital advertising landscape continues to change, our research lays a strong foundation for utilizing machine learning in contextual targeting for video ads.

## References

- [1] Smith, J. et al., "Contextual Targeting in Text-Based Digital Advertising: A Comparative Analysis," *Journal of Digital Marketing*, vol. 45, pp. 78-95, 2022.
- [2] Johnson, A. and Lee, K., "Deep Learning Approaches for Video Content Analysis," *IEEE Transactions on Multimedia*, vol. 23, pp. 1567-1582, 2021.
- [3] Zhang, Y. et al., "Multi-Modal Deep Learning for Advertisement Recommendation," *Proceedings of the 15th International Conference on Web Search and Data Mining*, pp. 503-511, 2023.
- [4] Brown, T. and White, R., "Advancements in Natural Language Processing for Contextual Understanding," *Computational Linguistics Journal*, vol. 37, pp. 221-240, 2022.