

Violence Detection in Video Using Advanced CV Algorithms

CH. KAVITHA¹, P. AKSHITH², S. ANSHIKA³, P. POOJITHA⁴, S. HARSHA VARDHAN⁵

¹ASSISTANT PROFESSOR, DEPARTMENT OF CSM, RAGHU ENGINEERING COLLEGE.

²DEPARTMENT OF CSM, RAGHU INSTITUTE OF TECHNOLOGY.

³DEPARTMENT OF CSM, RAGHU INSTITUTE OF TECHNOLOGY.

⁴DEPARTMENT OF CSM, RAGHU INSTITUTE OF TECHNOLOGY.

⁵DEPARTMENT OF CSM, RAGHU INSTITUTE OF TECHNOLOGY.

Abstract

This Investigation introduces an innovative dual-model framework for detecting violent content in video by synergistically combining YOLO (You Only Look Once) and RCNN (Region-based Convolutional Neural Networks) architectures. Our approach addresses the critical challenge of automated violence detection for applications in public safety monitoring, content moderation platforms, and surveillance systems. The proposed methodology leverages YOLO's computational efficiency with RCNN's precision through a novel weighted fusion algorithm that automatically balances their complementary strengths. Extensive evaluations across diverse datasets demonstrate that our hybrid system achieves 92.4% accuracy, surpassing single-model implementations while maintaining processing speeds suitable for practical deployment. The framework was rigorously tested on Hockey Fights, Movie Scenes, and Surveillance datasets, consistently demonstrating robust performance across varying environmental conditions, camera angles, and violence types.

Key Words: Computer Vision, YOLO, CNN, RCNN.

1. Introduction

The exponential growth of video content across digital platforms has created an urgent need for automated systems that can recognize potentially harmful material. Violence detection represents a particularly challenging yet essential application domain with significant implications for content moderation, public safety, and security monitoring. Manual review of video content has grown progressively impractical in the volume of data generated daily across social media, surveillance networks, and entertainment platforms.

While recent advances in deep learning have revolutionized computer vision capabilities, violence detection presents unique challenges requiring both

temporal understanding and contextual awareness. This Investigation explores how the integration of two distinct deep learning paradigms—YOLO and RCNN—can overcome limitations inherent in single-model approaches to violence detection.

Our work makes the following contributions:

- Development of a dual-architecture framework that strategically combines YOLO and RCNN models
- Introduction of a novel temporal-aware fusion mechanism that adapts to scene complexity
- Implementation of specialized preprocessing techniques optimized for violence detection
- Comprehensive performance analysis across multiple real-world scenarios and benchmark datasets
- Exploration of deployment considerations for practical implementation in resource-constrained environments

2. Related Work

2.1 Evolution of Violence Detection Approaches

Early Investigation in violence detection primarily relied on hand-crafted features and traditional machine learning algorithms. These approaches typically extracted low-level visual cues such as motion intensity, colour characteristics, and trajectory designs, combined with statistical classifiers. While foundational, these methods struggled with generalization across diverse scenarios and complex real-world conditions.

2.2 Deep Learning Advancements in Video Analysis

The emergence of deep neural networks has transformed video understanding capabilities. Convolutional architectures have demonstrated an exceptional ability to

extract meaningful spatial features, while sequential models like LSTM networks have enhanced temporal relationship modelling. Recent work has explored various architectural innovations for video analysis, though few have specifically addressed the unique requirements of violence detection.

2.3 Object Detection Frameworks

YOLO represents a milestone in real-time object detection with its single-pass approach that simultaneously predicts bounding boxes and class probabilities across the entire image. Meanwhile, the RCNN family employs a two-stage approach that first generates region proposals followed by classification, generally achieving higher precision at the cost of increased computational requirements. While both frameworks have been researched a lot for general object detection, their application to violence detection remains relatively unexplored.

2.4 Multi-Model Integration Strategies

Research on combining multiple deep learning models has demonstrated potential advantages in various computer vision tasks. However, Research on the subject is scarce. Focus on comprehensive integration approaches for violence detection. Our work addresses this gap by developing a specialized fusion strategy tailored to the distinct features of violent content in videos.

3. Methodology

3.1 Problem Definition

We formulate violence detection as a sequence classification task operating on temporal Video clips. For a given sequence of frames $F = \{f_1, f_2, \dots, f_n\}$, Our goal is to determine the probability $p(V|F)$ that the sequence contains violent content. We define violence as physical actions that exhibit aggressive behavior with the potential to cause harm, including fighting, assault, property destruction, and explosive events.

3.2 System Architecture

Our proposed framework consists of four primary components working in coordination:

3.2.1 Adaptive Video Processing Module

Unlike standard preprocessing pipelines, our system implements an adaptive approach that dynamically

adjusts processing parameters based on video characteristics:

- Scene-aware frame sampling that increases sampling rate during high-motion segments
- Multi-resolution frame processing (416×416 pixels for YOLO and 600×600 pixels for RCNN)
- Illumination normalization for challenging lighting conditions
- Overlapping temporal window extraction (50% overlap) with variable window sizes

3.2.2 Enhanced YOLO Detection Pipeline

We employ a customized implementation of YOLOv5 with several key modifications:

Integration of a specialized temporal feature extractor that captures motion patterns across frames

Reconfigured anchor box distribution optimized for human interaction detection

Implementation of dynamic confidence thresholding based on scene complexity

Focal loss implementation with automatic class weighting to address dataset imbalance

3.2.3 Context-Aware RCNN Pipeline

Our RCNN implementation utilizes Faster RCNN with ResNet-50 backbone, enhanced with:

- Self-attention mechanisms that highlight regions of potential interaction
- Custom feature pyramid network architecture for improved multi-scale detection
- Temporal context integration through the sequential region of interest alignment
- Specialized non-maximum suppression algorithm for tracking consistent detections

3.2.4 Adaptive Fusion Mechanism

Rather than using static weighting, our fusion approach implements an adaptive strategy:

$$S(t) = \alpha(t) \cdot S_YOLO(t) + (1 - \alpha(t)) \cdot S_RCNN(t)$$

where $\alpha(t)$ is determined dynamically based on:

Current frame characteristics (illumination, motion complexity, occlusion level)

Historical detection confidence trends, Computational resource availability, Scene context classification. Furthermore, a temporal consistency enforcement mechanism suppresses isolated false detections by requiring evidence across multiple frames.

3.3 Implementation Framework

Our implementation utilizes the PyTorch framework with the following configuration: Training conducted on a distributed system with 4 NVIDIA RTX 3090 GPUs

Batch sizes: 16 for YOLO, 8 for RCNN (optimized for GPU memory utilization)

Optimizer: Adam using a cosine annealing schedule for the learning rate

Extensive data augmentation suite: random flips, rotations, colour transformations, and motion simulation

Gradient accumulation for Large and impactful batch training

4. Experimental Evaluation

4.1 Dataset Construction and Preparation

We evaluated our approach on a comprehensive set of datasets:

Hockey Fights Dataset: 1,000 video clips balanced between violent and non-violent hockey sequences

Movie Violence Dataset: Manually annotated scenes from 200 diverse films spanning multiple genres

Surveillance Video Collection: 250 hours of CCTV footage with annotated violent incidents

Custom Multi-Environment Dataset: 5,000 clips gathered from diverse sources with varied conditions

For each dataset, we implemented:

- Standardized annotation protocol for consistent violence labeling
- Cross-dataset validation to ensure generalization capability
- Balanced training/validation/testing splits (60%/20%/20%)
- Scenario categorization for specialized performance analysis

4.2 Evaluation Protocol

We employed a comprehensive evaluation methodology using the following:

Standard metrics: Accuracy, Precision, Recall, F1-score

Temporal metrics: Mean Time to Detection (MTD), Temporal Intersection over Union (TIOU).

Computational efficiency: Frames per second (FPS), GPU memory utilization

User experience metrics: Perceived accuracy from human evaluators

4.3 Comparative Systems

We benchmarked against multiple state-of-the-art approaches:

- Single-architecture implementations (YOLO-only, RCNN-only)
- 3D convolutional networks (C3D)
- Two-stream CNN architectures
- Optical flow-based violence detection systems
- Transformer-based video analysis frameworks

5. Results and Analysis

5.1 Performance Benchmarks

Our hybrid approach demonstrated consistent performance advantages across all datasets:

Table 1: Comparative Performance Analysis

Method	Accuracy	Precision	Recall	F1-score	FPS
YOLO-only	88.7%	86.3%	89.5%	87.9%	45
RCNN-only	90.2%	91.7%	84.6%	88.0%	12
C3D	85.3%	82.1%	87.0%	84.5%	25
Two-stream CNN	89.1%	88.3%	87.9%	88.1%	18
Proposed Hybrid	92.4%	91.8%	91.5%	91.6%	20

Significantly, our system demonstrated exceptional strength performance in challenging scenarios

comprising low-light conditions, partial occlusions, and complex multi-person interactions.

5.2 Component Contribution Analysis

Through systematic ablation studies, we isolated the impact of individual components:

Temporal feature integration: +2.1% F1-score improvement

Attention mechanisms: +1.8% F1-score improvement

Adaptive fusion strategy: +3.5% F1-score over best single model

Scene-aware preprocessing: +1.2% F1-score enhancement

These improvements were not merely additive, as we observed synergistic effects when components operated together.

5.3 Qualitative Performance Assessment

Detailed analysis revealed distinctive operational patterns:

- YOLO excelled at detecting obvious violent actions with distinctive motion signatures
- RCNN demonstrated superior performance with subtle violence cues and complex interactions
- The hybrid system successfully mitigated individual weaknesses, particularly in edge cases
- Frame-level analysis showed complementary detection capabilities across various scenarios

5.4 Challenge Scenarios

We identified specific conditions that presented difficulties:

- Rapid camera movement creates motion artifacts
- Extreme lighting variations between frames
- Distinguishing between staged (e.g., film/sports) and actual violence
- Culturally-specific expressions of aggression

5.5 Efficiency and Deployment Considerations

The hybrid system achieved 20 FPS on average while maintaining high detection accuracy. We explored several optimization approaches:

- Model quantization (8-bit and 16-bit precision)
- TensorRT implementation for inference acceleration
- Batch processing optimizations
- Adaptive resolution processing based on computational constraints

6. Conclusion and Future Directions

This Investigation has demonstrated that a strategically designed hybrid approach to violence detection can overcome limitations inherent in single-model implementations. Our system integrates YOLO and RCNN architectures using an adaptive fusion mechanism, delivering exceptional accuracy while preserving efficient processing speeds. The consistent performance across diverse datasets highlights the robustness and generalizability of our approach. Several promising directions for future Investigation include:

- Integration of audio analysis for multimodal violence detection
- Development of deployable edge variants optimized for resource-constrained environments
- Extension to fine-grained violence categorization and severity assessment
- Exploration of federated learning approaches to enhance privacy preservation
- Investigation of unsupervised anomaly detection for identifying novel violence patterns

Acknowledgments

This study was supported by the Raghu Institute of Technology. The authors thank the anonymous reviewers for their constructive feedback Which Substantially improved this work.

References

- [1] Ren, S., He, K., Girshick, R., & Sun, J. (2015). Faster r-cnn: Towards real-time object detection with region proposal networks. *Advances in Neural Information Processing Systems*, 28, 91-99.
- [2] Redmon, J., & Farhadi, A. (2018). YOLOv3: An incremental improvement. *arXiv preprint arXiv:1804.02767*.
- [3] Jocher, G., Stoken, A., & Borovec, J. (2020). YOLOv5. *GitHub repository*. <https://github.com/ultralytics/yolov5>.

- [4] Lin, T. Y., Goyal, P., Girshick, R., He, K., & Dollár, P. (2017). Focal loss for dense object detection. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 42(2), 318-327.
- [5] He, K., Zhang, X., Ren, S., & Sun, J. (2016). Deep residual learning for image recognition. *IEEE Conference on Computer Vision and Pattern Recognition*, 770-778.
- [6] Tran, D., Wang, H., Torresani, L., Ray, J., LeCun, Y., & Paluri, M. (2018). A closer look at spatiotemporal convolutions for action recognition. *IEEE Conference on Computer Vision and Pattern Recognition*, 6450-6459.
- [7] Wu, Y., & He, K. (2018). Group normalization. *European Conference on Computer Vision*, 3-19.
- [8] Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., Kaiser, Ł., & Polosukhin, I. (2017). Attention is all you need. *Advances in Neural Information Processing Systems*, 30, 5998-6008.
- [9] Loshchilov, I., & Hutter, F. (2017). Decoupled weight decay regularization. *International Conference on Learning Representations*.
- [10] Zhang, H., Cisse, M., Dauphin, Y. N., & Lopez-Paz, D. (2018). Mixup: Beyond empirical risk minimization. *International Conference on Learning Representations*.
- [11] Howard, A. G., Zhu, M., Chen, B., Kalenichenko, D., Wang, W., Weyand, T., Andreetto, M., & Adam, H. (2017). MobileNets: Efficient convolutional neural networks for mobile vision applications. *arXiv preprint arXiv:1704.04861*.
- [12] Li, Y., Li, W., Danelljan, M., & Zhang, K. (2022). Deep spatiotemporal feature learning for violence detection in surveillance videos. *Pattern Recognition*, 128, 108671.
- [13] Chen, X., Zhang, Y., Xu, Y., & Tan, W. (2022). TFNet: Temporal fusion network for violence detection with optical flow enhancement. *Computer Vision and Image Understanding*, 215, 103311.
- [14] Kim, J., Han, D., & Choi, J. (2022). Adaptive event detection in surveillance environments using multi-scale feature aggregation. *IEEE International Conference on Acoustics, Speech and Signal Processing*, 2837-2841.
- [15] Wang, Z., Zhou, Y., & Yu, M. (2022). Context-aware violent behavior recognition for public safety monitoring. *Multimedia Systems*, 28(4), 1145-1158.
- [16] Smith, R., Johnson, L., & Lee, T. (2022). Comparative analysis of object detection frameworks for real-time violence monitoring. *IEEE Transactions on Systems, Man, and Cybernetics: Systems*, 52(7), 4351-4363.
- [17] Garcia-Martinez, E., Lopez-Torres, M., & Fernandez-Caballero, A. (2023). Deep learning approaches for violence detection in social media video content. *Expert Systems with Applications*, 211, 118532.
- [18] Taylor, B., Richardson, K., & Cooper, M. (2023). Privacy-preserving violence detection using federated learning in edge environments. *IEEE Internet of Things Journal*, 10(4), 3192-3205.
- [19] Williams, S., Anderson, P., & Thompson, J. (2023). Benchmarking deep learning models for abusive content detection in video streaming platforms. *IEEE Conference on Computer Vision and Pattern Recognition Workshops*, 1123-1132.
- [20] Harris, D., Miller, A., & Brown, N. (2023). Lightweight models for violence detection on resource-constrained devices. *ACM Transactions on Sensor Networks*, 19(1), 1-24.
- [21] Jackson, R., Wilson, C., & Thomas, S. (2023). Transformer-based approaches for long-range violence detection in surveillance footage. *IEEE Transactions on Image Processing*, 32(5), 2185-2197.
- [22] Patel, V., Singh, A., & Kaur, G. (2023). Transfer learning strategies for violence detection across diverse domains. *Neural Computing and Applications*, 35(7), 4891-4904.
- [23] Chen, L., Zhang, R., & Pei, M. (2023). Uncertainty-aware fusion of detection models for robust violence identification. *IEEE International Conference on Multimedia and Expo*, 1-6.
- [24] Martinez, J., Rodriguez, P., & Garcia, D. (2023). Cross-domain adaptation techniques for violence detection in varied environmental conditions. *Computer Vision and Image Understanding*, 226, 103568.
- [25] Thompson, E., Wright, K., & Davis, M. (2023). Explainable AI approaches for violence detection systems. *AI Ethics*, 3(1), 87-101.
- [26] Lee, S., Park, J., & Kim, H. (2023). Contrastive learning for self-supervised violence detection in untrimmed videos. *IEEE/CVF Winter Conference on Applications of Computer Vision*, 1378-1387.

[27] Robertson, A., Phillips, J., & Evans, C. (2024). Multi-stream transformer architectures for fine-grained violence categorization. *Neural Networks*, 168, 542-557.

[28] Nguyen, T., Rahman, S., & Tran, D. (2024). Semi-supervised learning for violence detection with limited labeled data. *Pattern Recognition*, 138, 109426.

[29] Collins, M., Turner, R., & Baker, S. (2024). Real-time violence detection at the edge: Challenges and solutions. *ACM Transactions on Embedded Computing Systems*, 23(1), 1-25.

[30] Yamada, K., Peterson, L., & Garcia, M. (2024). Adaptive multi-resolution processing for efficient violence detection in streaming video. *IEEE Transactions on Circuits and Systems for Video Technology*, 34(2), 921-935.