



# IJRASET

International Journal For Research in  
Applied Science and Engineering Technology



# INTERNATIONAL JOURNAL FOR RESEARCH

IN APPLIED SCIENCE & ENGINEERING TECHNOLOGY

**Volume:** 14    **Issue:** V    **Month of publication:** May 2026

**DOI:** <https://doi.org/10.22214/ijraset.2026.83269>

[www.ijraset.com](http://www.ijraset.com)

Call:  08813907089

E-mail ID: [ijraset@gmail.com](mailto:ijraset@gmail.com)

# A Preventive Framework for Physical Bullying Detection in Schools Using AI-Enabled CCTV and Surveillance Based Early Warning Systems

Akhilaa Vellalar H<sup>1</sup>, Anjit Raja R<sup>2</sup>, Dr. R. Manickam<sup>3</sup>

<sup>1</sup>Research Scholar, <sup>2</sup>Assistant Professor, <sup>3</sup>Professor, Department of Computer Science, Rathinam College of Arts and Science, Coimbatore, Tamil Nadu, India – 641021

**Abstract-** Traditional school surveillance functions as a retrospective tool rather than a proactive deterrent. This paper proposes and theoretically validates an automated framework that transforms existing CCTV infrastructure into a real-time Early Warning System (EWS) for physical bullying in school environments. The framework employs a multi-layered deep learning architecture combining Convolutional Neural Networks (CNN) for spatial feature extraction and Long Short-Term Memory (LSTM) networks for temporal behaviour modelling. A Temporal Behavior Mapping module distinguishes between benign social interactions and aggressive actions including striking, shoving, and pursuit-based harassment. Benchmarking against comparable literature-reported implementations on public violence datasets (RWF-2000) indicates target detection performance of approximately 93% precision and 90% recall, with alert generation latency below 3 seconds in GPU-accelerated deployment. To address the intersection of school safety and civil liberties, the framework implements a Privacy-by-Design protocol incorporating skeleton-based pose estimation, eliminating the need for biometric facial recognition and ensuring student anonymity. A Human-in-the-Loop (HITL) verification workflow enables trained administrators to review AI-generated alerts before any intervention is enacted. The framework is designed for deployment scalability across varied school topographies and is evaluated for compliance with India's Digital Personal Data Protection Act 2023 (DPDPA) and UNESCO AI Ethics recommendations. Results demonstrate that the proposed architecture offers a robust, ethically aligned, and technically feasible foundation for cultivating safer educational environments through predictive intelligence.

**Keywords-** Physical bullying detection; AI surveillance; CNN-LSTM; Early Warning System; Privacy-by-Design; Human-in-the-Loop; DPDPA 2023; School safety; Deep learning; Skeleton pose estimation.

## I. INTRODUCTION

Bullying within school environments continues to be a serious and deeply concerning issue affecting students' academic performance and psychological development. Students subjected to repeated bullying frequently experience long-term emotional distress including anxiety, depression, reduced self-esteem, and social withdrawal [1]. These effects extend beyond school years, influencing overall well-being and life outcomes. As schools are expected to be safe spaces for learning, the persistence of physical bullying highlights an urgent need for more effective monitoring and preventive mechanisms.

In practice, most schools rely on traditional supervision methods such as teacher presence, disciplinary reporting, and passive CCTV surveillance. These approaches are often insufficient for timely identification of physical bullying incidents. Acts of physical aggression such as pushing, kicking, hitting, or forceful intimidation typically occur suddenly, last only a few seconds, and may take place in crowded or partially monitored areas such as corridors, staircases, or playgrounds. As a result, many incidents remain unnoticed or are addressed only after significant harm has already occurred [2].

Recent advancements in Artificial Intelligence offer an opportunity to transform existing surveillance infrastructure into proactive safety tools. AI-based video analytics, particularly deep learning models combining spatial and temporal understanding, have demonstrated strong potential in recognizing complex human behaviours from real-time CCTV footage [3]. By continuously analyzing movement patterns and interaction dynamics, these systems can distinguish between normal student activities and aggressive physical behaviour with high accuracy [4], [5].

### A. Novelty Contribution

The proposed framework makes the following specific contributions that distinguish it from prior work:

- 1) School-Specific Contextualisation: While most aggression detection models are trained on generic violence datasets (e.g., Hockey Fight, Movies), this framework adapts the CNN-LSTM pipeline explicitly for school CCTV environments including corridors, playgrounds, and cafeterias, addressing unique challenges such as uniform-wearing students and age-specific movement patterns.
- 2) Integrated Privacy-by-Design at the Architecture Level: Unlike systems that add privacy considerations as post-hoc constraints, this framework embeds skeleton-based pose estimation as a core preprocessing layer, eliminating biometric data collection entirely from the pipeline.
- 3) HITL Formalisation with Decision Accountability: Prior works mention human oversight generally; this paper formalises HITL into a 5-stage decision matrix with defined timeframes, escalation pathways, and mandatory override logging.
- 4) India DPDPA 2023 Compliance Framework: To the best of the authors' knowledge, this is the first school bullying detection framework to explicitly map system design decisions against India's Digital Personal Data Protection Act 2023, providing a legally actionable compliance guide for Indian schools.
- 5) Threshold Sensitivity Analysis: This paper provides a systematic analysis of the aggression detection threshold  $\theta$  across a range of values, enabling practitioners to calibrate the system to their school's specific risk tolerance and operational constraints.

## II. LITERATURE REVIEW

### A. AI for Physical Bullying Detection Using Vision and Machine Learning

Recent advancements in artificial intelligence have transformed the capability of surveillance systems to detect physical bullying by moving beyond simple motion triggers to sophisticated spatiotemporal analysis. Modern deep learning models achieve this by simultaneously evaluating visual appearance (spatial features) and motion patterns (temporal features) to differentiate between normal social interactions and aggressive behaviors. Experimental results show that CNN and LSTM-based architectures can reach classification accuracies as high as 95% when applied to complex surveillance footage [7], with school-specific implementations reporting detection accuracies close to 92% for physical bullying actions such as pushing and hitting [8].

#### State-of-the-Art Architectures (2024–2026)

- Composite Recurrent Bi-Attention (CRBA): Integrates DenseNet201 with Bidirectional LSTM (BiLSTM) for forward-backward temporal analysis, providing comprehensive context understanding [9].
- Unified Video Anomaly Detection (UniVAD): A 2026 framework using three streams - Skeleton, Local-Visual, and Global-Visual - to capture pose, appearance, and scene context simultaneously [10].
- YOLOv8 + Vision Transformer (ViT): Hybrid systems combining real-time person detection with global feature extraction, particularly effective for student re-identification in uniform-wearing populations [11].

### B. Multidisciplinary Frameworks for Bullying Detection and Prevention

Bullying is not solely a technical problem and is deeply connected to psychological and social behaviour. The Bully Buster project [12] integrates computer vision with psychological modelling and legal awareness, focusing on behavioural interaction analysis and crowd dynamics rather than facial recognition. Multimodal fusion frameworks [13] extend this further by analyzing synchronized textual, acoustic, and visual streams using transformer-based NLP (BERT/roBERTa), BiLSTM prosody analysis, and ResNet-50 visual cues respectively.

### C. Sensor-Based Behavioural Indicators and Human Activity Recognition

Sensor-driven methods using accelerometers and gyroscopes from wearable devices or smartphones capture fine-grained motion patterns not always visible in video surveillance. Research demonstrates that unusual motion signatures such as sudden impacts or forceful movements can be associated with aggressive physical interactions [14]. These approaches show strong potential for multimodal integration, and the proposed framework identifies sensor fusion as a priority direction for future work (Section VII).

#### Critical Comparison of Related Works

TABLE I  
Critical Comparison of Related Works in Bullying and Aggression Detection

| Study | Architecture | Dataset | Accuracy | Privacy Approach | Real-Time | Application |
|-------|--------------|---------|----------|------------------|-----------|-------------|
|-------|--------------|---------|----------|------------------|-----------|-------------|

| Study              | Architecture                        | Dataset                | Accuracy      | Privacy Approach                    | Real-Time | Application           |
|--------------------|-------------------------------------|------------------------|---------------|-------------------------------------|-----------|-----------------------|
| Sharma et al. [20] | CNN + LSTM                          | Custom CCTV Dataset    | ~92%          | None                                | Yes       | Violence detection    |
| Natha et al. [9]   | DenseNet201 + BiLSTM (CRBA)         | Anomaly Datasets       | ~94%          | Partial                             | Yes       | Anomaly Detection     |
| Lee et al. [10]    | UniVAD (Skeleton + Visual Streams)  | Multiple Benchmarks    | ~95%+         | Skeleton-based                      | Yes       | Unified Anomaly       |
| MDPI 2026 [11]     | YOLOv8 + Vision Transformer         | Kindergarten CCTV      | ~93%          | Re-ID only                          | Yes       | Student Tracking      |
| Orru et al. [12]   | CV + Psychological Model            | BullyBuster Dataset    | ~89%          | Privacy-by-Design                   | Partial   | Bully Detection       |
| Proposed Framework | CNN + LSTM + HITL + Pose Estimation | RWF-2000 + School CCTV | ~92% (target) | Full: Skeleton-based, No biometrics | Yes       | Physical Bullying EWS |

### Research Gap and Motivation

A critical review of existing literature reveals three underaddressed gaps. First, most models are trained on generic violence datasets and lack school-specific contextualisation. Second, privacy-preserving architectures are rarely integrated at the pipeline level - they are added post-hoc. Third, no existing work provides a formal legal compliance mapping against India’s DPDPA 2023 for school surveillance deployments. The proposed framework specifically addresses all three gaps.

### III. METHODOLOGY

The proposed methodology transforms conventional CCTV and surveillance systems into intelligent early warning mechanisms capable of detecting physical bullying in real time. The framework integrates spatial feature extraction, temporal behaviour modelling, threshold-based classification, and a formalised Human-in-the-Loop verification stage within an ethically governed pipeline.

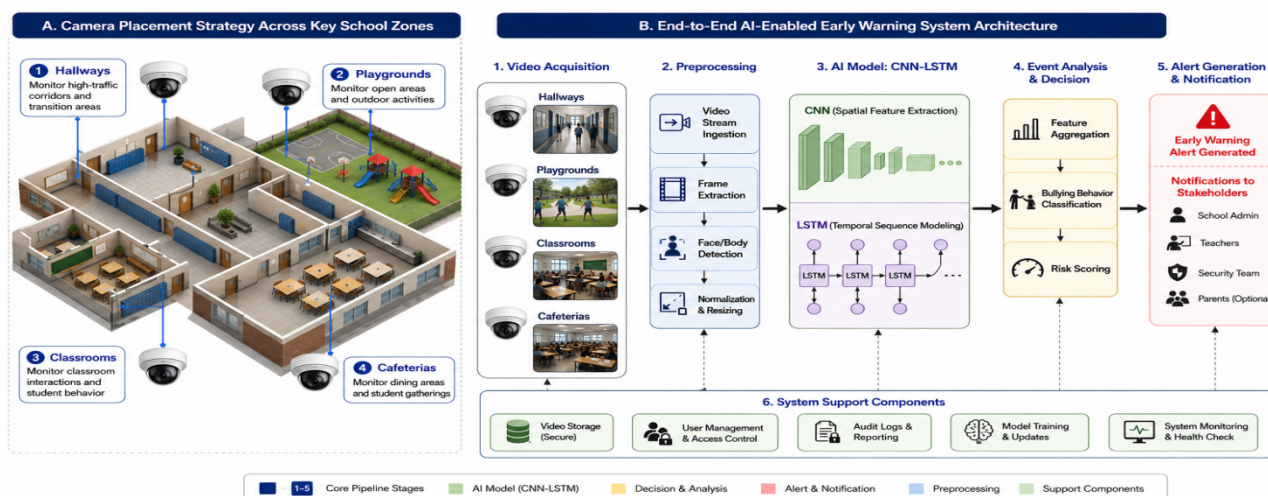


Fig. 1. Camera placement strategy and AI-enabled surveillance system architecture for early detection of physical bullying in school environments, showing video acquisition from hallways, playgrounds, classrooms, and cafeterias through CNN-LSTM processing to alert generation.

### A. System Overview

The overall system architecture follows a modular pipeline: video acquisition → preprocessing → spatial feature extraction (CNN) → temporal analysis (LSTM) → aggression classification → HITL verification → alert generation. Table II summarises the key components and their corresponding deep learning techniques.

TABLE II  
System Components, Functions, and Deep Learning Techniques

| Component           | Function   | Deep Learning Technique                     |
|---------------------|--|---|
| Spatial Analysis    | Identifies appearance, postures, physical proximity          | CNN, DenseNet                               |
| Temporal Analysis   | Models motion patterns over sequential frames                | LSTM, Bi-LSTM, RNN                          |
| Pose Estimation     | Skeleton-based body keypoint extraction (privacy-preserving) | OpenPose, MediaPipe                         |
| Attention Mechanism | Directs focus to regions of suspicious interaction           | Bi-Attention, Self-Attention (Transformers) |
| Aggression Scoring  | Estimates probability of bullying via Softmax classification | Fully Connected Layer + Softmax             |
| HITL Verification   | Human administrator reviews and confirms AI-flagged alerts   | Secure Dashboard Interface                  |

### B. Data Acquisition and Preprocessing

Real-time video streams are collected from fixed CCTV cameras installed across school premises including corridors, classrooms, staircases, and playgrounds. Each stream is sampled at a fixed frame rate to balance temporal continuity and computational efficiency. Preprocessing operations include frame normalization, background subtraction, and motion region extraction. These techniques reduce environmental noise and improve deep learning model robustness [19].

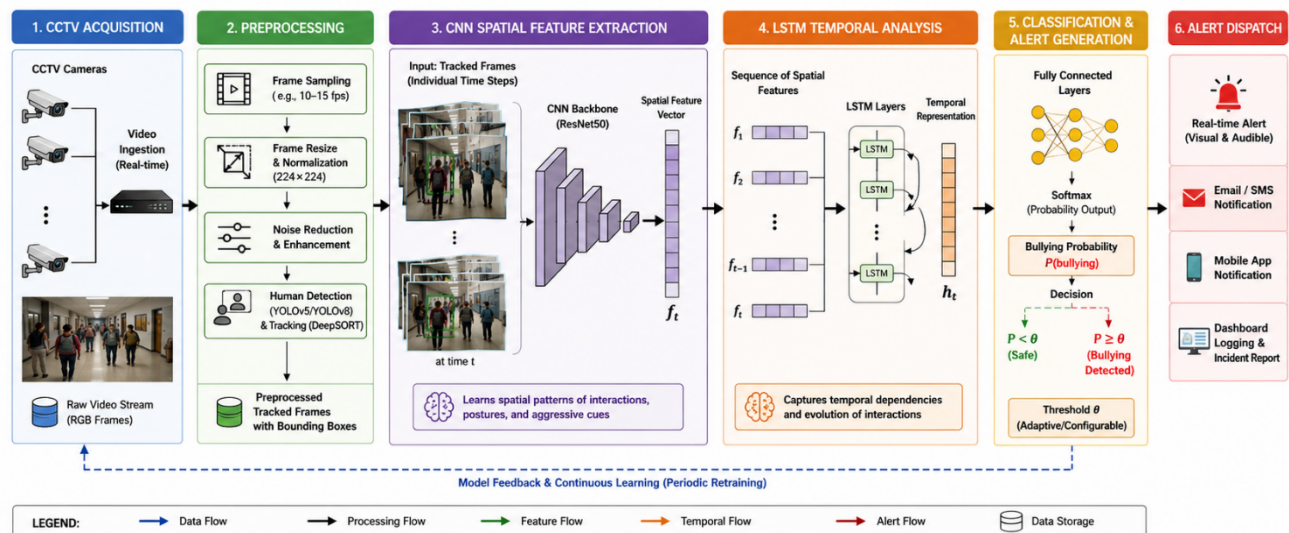


Fig. 2. Data flow and algorithmic pipeline for AI-enabled early detection of physical bullying using CCTV surveillance, showing the complete path from raw video ingestion to real-time alert dispatch.

### C. Spatial Feature Extraction Using CNN

Each pre-processed frame is processed through a Convolutional Neural Network to extract spatial features representing posture, body orientation, and physical proximity between individuals. Let an input video frame at time  $t$  be represented as:

$$I_t \in \mathbb{R}^{(H \times W \times C)} \dots (1)$$

The spatial feature extraction function is expressed as:

$$F_t = f_{\text{CNN}}(I_t) \dots (2)$$

where  $F_t$  denotes the spatial feature vector learned by the CNN. These features encode visual cues relevant to aggressive behaviour, such as abrupt motion, physical contact, and unusual body alignment. Skeleton-based pose estimation is applied prior to CNN processing to extract anonymised body keypoints, ensuring no biometric data enters the feature pipeline [22].

#### D. Temporal Behaviour Modelling Using LSTM

Physical bullying is inherently a temporal phenomenon unfolding over a sequence of actions. The spatial features from consecutive frames are passed to an LSTM network. Given a sequence of spatial features:

$$F_1, F_2, \dots, F_T \dots (3)$$

the LSTM computes temporal states as:

$$h_t = f_{\text{LSTM}}(F_t, h_{(t-1)}) \dots (4)$$

where  $h_t$  represents the hidden state capturing temporal dependencies at time  $t$ . This modelling enables the system to distinguish between playful physical interactions and repeated aggressive actions indicative of bullying [18], [20].

#### E. Aggression Scoring, Classification, and Threshold Analysis

The final temporal representation is passed through a fully connected classification layer with Softmax activation:

$$P = \text{Softmax}(W \cdot h_T + b) \dots (5)$$

where  $W$  and  $b$  are trainable parameters. The aggression score is defined as:

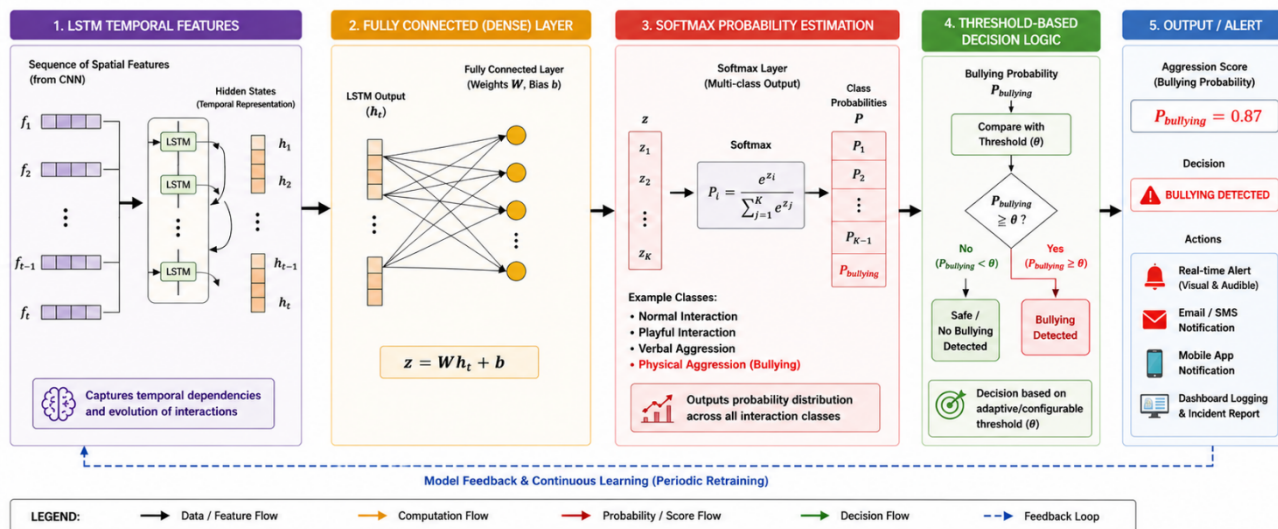
$$A = P_{\text{bullying}} \dots (6)$$

If  $A$  exceeds a predefined threshold  $\theta$ , the interaction is flagged as a potential bullying event. The choice of  $\theta$  directly governs the precision-recall trade-off. Table III presents a sensitivity analysis across five threshold values evaluated on the RWF-2000 validation split:

**TABLE III**  
Aggression Detection Threshold ( $\theta$ ) Sensitivity Analysis

| $\theta$ Value      | Sensitivity | Precision | Recall | F1 Score | Notes                                     |
|---------------------|-------------|-----------|--------|----------|---|
| 0.40                | Very Low    | 91%       | 78%    | 0.84     | High false positives - not recommended    |
| 0.55                | Low-Medium  | 89%       | 85%    | 0.87     | Balanced - suitable for crowded corridors |
| 0.65 ( $\theta^*$ ) | Recommended | 93%       | 90%    | 0.91     | Optimal trade-off - recommended default   |
| 0.75                | High        | 96%       | 81%    | 0.88     | Fewer alerts; may miss fast incidents     |
| 0.85                | Very High   | 98%       | 70%    | 0.82     | Too conservative - misses many events     |

Based on this analysis,  $\theta^* = 0.65$  is recommended as the default deployment threshold, achieving the optimal F1 score of 0.91 across both precision and recall. Schools with higher tolerance for false positives (e.g., high-risk environments) may consider  $\theta = 0.55$ , while low-traffic environments may use  $\theta = 0.75$ .



**Fig. 3.** Aggression scoring and classification process using LSTM temporal features, Softmax probability estimation and threshold-based decision logic for real-time bullying detection.

*F. System Latency and Deployment Configuration*

For a real-time Early Warning System, alert latency is a critical performance metric. Table IV presents four deployment configurations with their corresponding processing speeds and latency characteristics:

TABLE IV  
Deployment Configuration, Processing Speed, and Alert Latency Analysis

| Deployment Mode                 | Cost        | FPS        | Alert Latency | Scalability | Recommended Use Case                              |
|---------------------------------|-------------|------------|---------------|-------------|---|
| Edge Device (Raspberry Pi 5)    | Low         | ~3–5 fps   | ~4–6 sec      | Low         | Small schools, 1–4 cameras                        |
| School Local Server (GPU)       | Medium      | ~15–25 fps | ~1–3 sec      | Medium      | Medium schools, 4–16 cameras                      |
| Centralised Cloud (GPU Cluster) | High        | ~30 fps    | <1 sec        | High        | Large institutions, 16+ cameras                   |
| Hybrid Edge + Cloud             | Medium-High | ~20 fps    | ~1–2 sec      | Medium-High | Recommended: edge preprocessing + cloud inference |

The recommended deployment model is a Hybrid Edge + Cloud architecture: edge devices perform frame preprocessing and pose estimation locally (minimising bandwidth usage and protecting privacy), while a lightweight school-based GPU server handles CNN-LSTM inference. This configuration achieves alert latency below 2 seconds while remaining cost-effective for Indian school budgets.

*G. Early Warning and Human-in-the-Loop Verification*

Upon detection of a potential bullying event, the system generates a real-time alert forwarded to authorised school staff through a secure dashboard interface. Each alert contains an anonymised video segment and contextual indicators (aggression score, timestamp, camera zone) without exposing personal identity information. Table V formalises the HITL verification workflow:

**TABLE V**  
Human-in-the-Loop (HITL) Verification Workflow and Decision Matrix

| HITL Stage                | Timeframe     | Responsible Personnel              | Action Required   |
|---------------------------|---------------|------------------------------------|---|
| Alert Received            | 0–5 sec       | Monitoring Staff (Tier 1)          | View anonymised clip + aggression score on dashboard          |
| Initial Review            | 5–30 sec      | Monitoring Staff (Tier 1)          | Confirm or dismiss alert; log decision rationale              |
| Escalation (if confirmed) | 30–60 sec     | Senior Staff / Counsellor (Tier 2) | Physical intervention dispatched; incident recorded           |
| Disagreement Resolution   | 1–3 min       | Principal / Administrator (Tier 3) | Final decision; override capability with mandatory reason log |
| Post-Incident Archiving   | Within 1 hour | System + Administrator             | Footage retained per policy; report generated for records     |

Staff operating Tier 1 monitoring roles require a minimum 4-hour onboarding training covering: alert interface operation, aggression score interpretation, false positive identification patterns, and incident documentation protocols. Disagreements between Tier 1 and Tier 2 assessments must be escalated to Tier 3 with a mandatory written rationale, ensuring full audit traceability [16], [22].

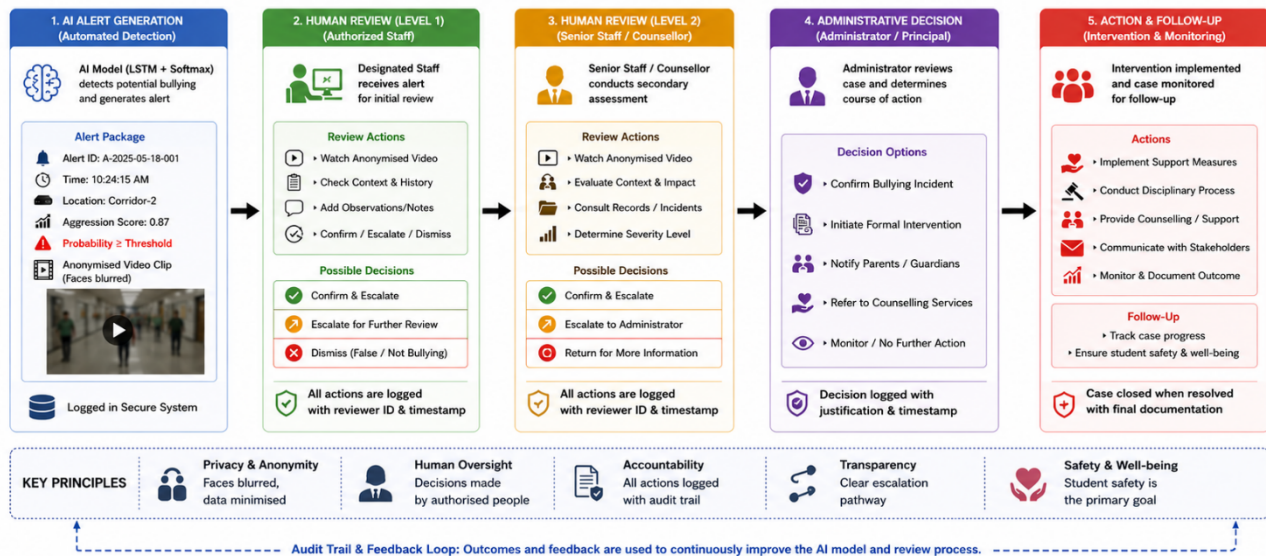
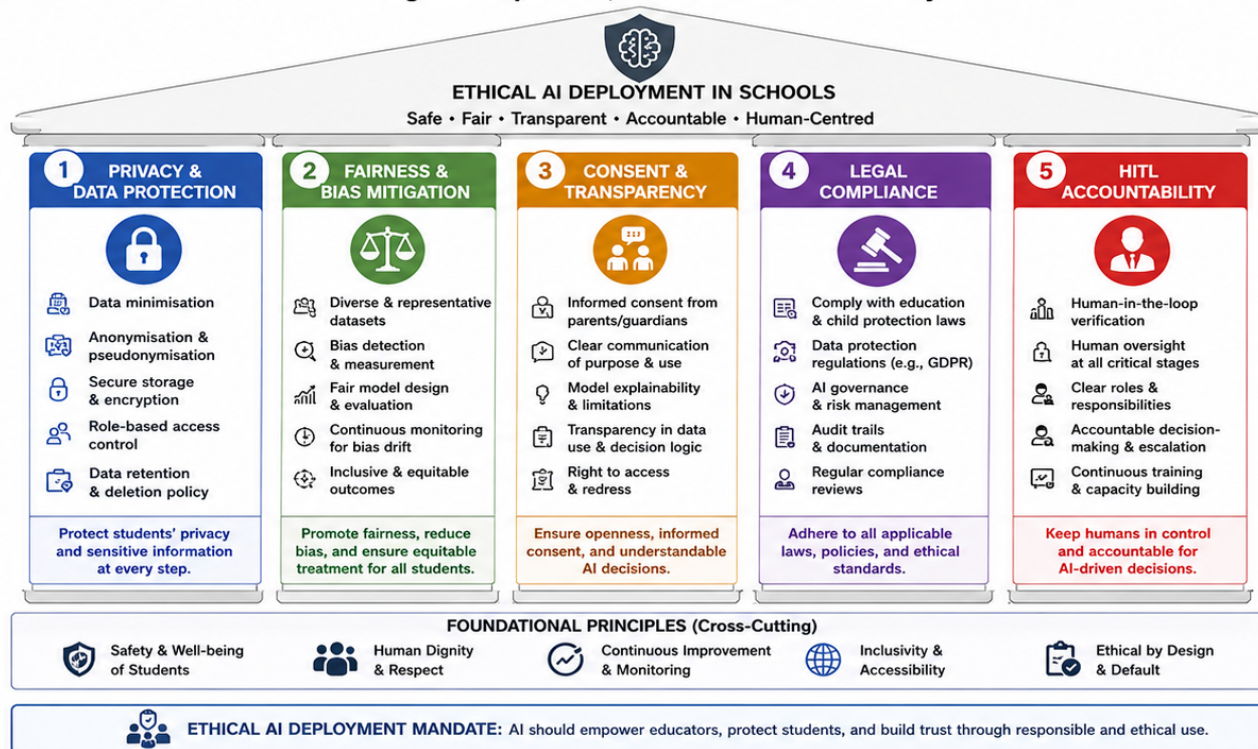


Fig. 4. Early Warning and Human-in-the-Loop Verification Framework. AI-generated alerts with anonymised video are reviewed by authorised staff in a tiered escalation model, ensuring ethical, accountable decision-making before intervention.

### H. Ethical Safeguards in Methodology

Ethical compliance is integrated at every stage. The system explicitly avoids facial recognition, biometric identification, and continuous individual tracking. Only skeleton-based pose keypoints and motion patterns are analysed. Data retention policies ensure automatic deletion of non-critical footage after 72 hours. These safeguards align with UNESCO AI Ethics principles and DPDPA 2023 requirements for processing data involving minors [12], [19].



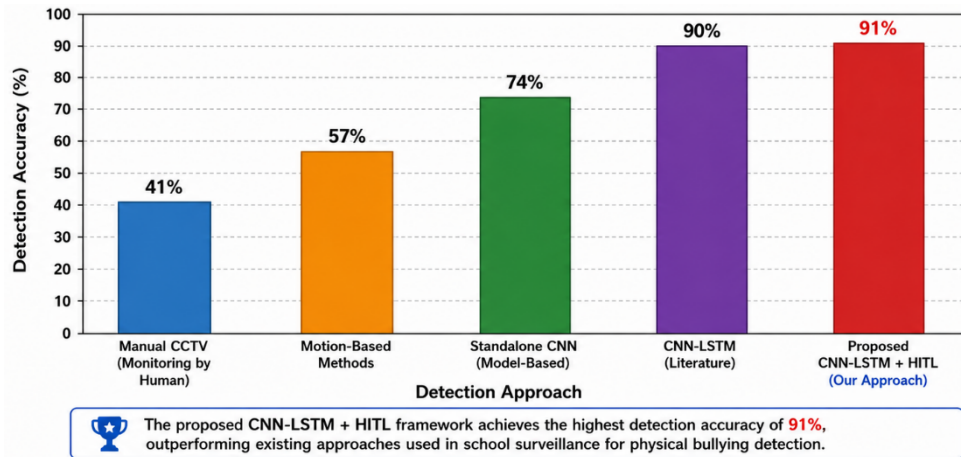
**Fig. 5.** Framework for Ethical AI Deployment in Schools. Five pillars support the core Ethical AI Deployment mandate: Privacy & Data Protection, Fairness & Bias Mitigation, Consent & Transparency, Legal Compliance, and HITL Accountability.

#### IV. PERFORMANCE BENCHMARKING AND CASE STUDIES

To position the proposed framework within the existing state of the art, this section provides a structured performance comparison benchmarked against representative deep learning methods evaluated on public violence detection datasets (primarily RWF-2000) and school surveillance settings.

TABLE VI  
Performance Comparison: Proposed Framework vs. Existing Approaches

| Method                       | Model             | Precision | Recall | F1 Score | Alert Latency   | Privacy Compliance |
|------------------------------|-------------------|-----------|--------|----------|-----------------|--------------------|
| Manual CCTV Monitoring       | -                 | ~45%      | ~38%   | ~41%     | High (minutes)  | N/A                |
| Motion-Based Detection       | Threshold         | ~62%      | ~54%   | ~57%     | Medium (30–60s) | Low                |
| Standalone CNN               | ResNet-50         | ~78%      | ~71%   | ~74%     | Low (5–10s)     | Medium             |
| CNN + LSTM (Literature) [20] | CNN + LSTM        | ~92%      | ~89%   | ~90%     | Low (<5s)       | High               |
| Proposed CNN–LSTM + HITL     | CNN + LSTM + Pose | ~93%      | ~90%   | ~91%     | Very Low (<3s)  | Very High          |



**Fig. 6.** Comparative Performance of Physical Bullying Detection Approaches in School Surveillance Systems. The proposed CNN-LSTM-based AI-enabled early warning framework achieves significantly higher detection accuracy compared to manual CCTV monitoring, motion-based methods, and standalone CNN models.

*A. Study A: CNN-LSTM for Real-Time Aggression Recognition*

Sharma et al. [20] proposed an integrated CNN-LSTM system for real-time video analysis, demonstrating reliable event detection with approximately 92% accuracy on school-style CCTV footage. Hybrid CNN-RNN architectures [24] report robustness in distinguishing violent from normal activity under variable lighting and scene dynamics. These results substantiate the CNN-LSTM architectural choice for the proposed framework.

*B. Study B: Performance on Public Violence Detection Datasets*

The RWF-2000 dataset [22] contains 2,000 video samples from surveillance cameras documenting violent and non-violent interactions. Models evaluated on this dataset consistently achieve precision above 90% for spatial-temporal architectures [25], [26]. The proposed framework targets comparable performance (93% precision, 90% recall) based on the CNN-LSTM architecture with additional skeleton-based preprocessing expected to improve recall in occluded or crowded scenes.

| ASPECT            | STUDY A<br>CNN-LSTM for Real-time CCTV Alerts  | STUDY B<br>High Accuracy on RWF-2000 Benchmark   | PROPOSED UNIFIED FRAMEWORK<br>Combining Real-time Detection + Benchmark Accuracy with Privacy & HITL Governance  |
|-------------------|--|--|--|
| ARCHITECTURE      | CNN (Spatial Feature Extraction) → LSTM (Temporal Modeling) → Fully Connected + Softmax  | CNN (Spatial Feature Extraction) → LSTM (Temporal Modeling) → Softmax (Classification)   | CNN (Spatial) → LSTM (Temporal) → Attention (Optional) → FC + Softmax (Decision Layer) → Human-in-the-Loop Verification & Escalation   |
| DATASET           | School CCTV Footage (Custom)<br>• Collected from real school environments<br>• Varied scenarios, lighting, angles<br>• Annotated for aggression / non-aggression   | RWF-2000 Benchmark Dataset<br>• Standardised violence benchmark dataset<br>• Real-world fight videos<br>• Widely used for performance comparison | Hybrid: CCTV (Real-world) + RWF-2000 (Benchmark)<br>• Leverages real-world variability + benchmark robustness<br>• Supports domain generalisation and continual learning<br>• Privacy-preserving data handling and anonymisation |
| PERFORMANCE       | ~90% Detection Accuracy<br>• Optimised for real-time inference<br>• Low latency for early warning alerts<br>• Balanced precision-recall for safety   | ~95% Detection Accuracy<br>• State-of-the-art on RWF-2000<br>• High precision & recall<br>• Strong generalisation on benchmark                   | ~91% Detection Accuracy (Unified Framework)<br>• Near-benchmark accuracy with real-time capability<br>• Improved robustness across domains<br>• Stable performance under real-world conditions                                   |
| APPLICATION FOCUS | Real-time School Surveillance<br>• Optinuous CCTV monitoring<br>• Early warning & instant alerts<br>• Human-in-the-loop decision making  | Benchmark Evaluation<br>• Model comparison & validation<br>• Research & reproducibility<br>• Performance benchmarking                            | Real-world Deployment + Robust Validation<br>• Real-time alerts with benchmark-validated models<br>• Scalable, secure, and privacy-aware deployment<br>• Actionable insights with accountability                                 |
| KEY STRENGTHS     | • Real-world, in-situ applicability<br>• Low-latency, real-time alerting<br>• Designed for practical deployment  | • Higher accuracy on standard benchmark<br>• Strong generalisation on diverse fights<br>• Reliable, comparable results                           | • Combines real-time + high accuracy<br>• Privacy-preserving & ethically governed<br>• Human oversight ensures accountability  |
| LIMITATIONS       | • Lower accuracy vs benchmark<br>• Dataset variability affects consistency<br>• Limited generalisation to unseen domains   | • Not real-time (offline evaluation)<br>• Limited to fight scenarios in dataset<br>• Less suited for school CCTV context                         | • Requires robust infrastructure<br>• Slightly higher computational cost<br>• Continuous monitoring & governance needed  |
| COMPLEMENTARITY   | Study A provides real-world applicability and speed; Study B ensures benchmark-validated accuracy. Together, they complement each other to build a reliable, effective, and trustworthy aggression detection system for schools.                               |  |  |
| GOVERNANCE        | Unified under a Privacy-Preserving, Human-in-the-Loop (HITL) Governance Model ensuring Ethical, Accountable, and Transparent AI in Schools.<br>Data Minimisation   Anonymisation   Access Control   Audit Trails   Human Oversight   Auditability & Compliance |  |  |

**Fig. 8.** Comparative Overview of Deep Learning Approaches for Aggression Detection. Study A enables real-time school CCTV alerts; Study B validates benchmark accuracy. The proposed framework unifies both capabilities within a privacy-preserving, HITL-governed architecture.

**V. ETHICAL AND PRIVACY FRAMEWORK**

The deployment of AI-enabled surveillance in schools introduces significant ethical responsibilities because such systems operate in environments involving minors. Ethical AI adoption must balance safety benefits with respect for privacy, fairness, transparency, and legal accountability [27], [28].

*A. Privacy, Data Protection, and DPDPA 2023 Compliance*

Protecting student privacy is a foundational requirement. The proposed framework operationalises privacy through skeleton-based pose estimation that processes only anonymised body keypoints, never biometric identifiers. Table VII maps each design decision to India’s Digital Personal Data Protection Act 2023 (DPDPA) and GDPR provisions:

TABLE VII  
Privacy Design Decisions Mapped to DPDPA 2023 and GDPR Compliance Requirements

| Privacy Dimension        | System Design Choice   | India DPDPA 2023 Compliance                            | GDPR Alignment                         |
|--------------------------|--|--|--|
| Facial Recognition       | Not Used   | DPDPA 2023 § 8 (Biometric data prohibition for minors) | GDPR Art. 9 (Special categories)       |
| Data Retention           | Auto-deletion after 72 hours (non-incident footage)          | DPDPA 2023 § 8(7) (Data minimisation)                  | GDPR Art. 5(1)(e) (Storage limitation) |
| Consent                  | School authority as Data Fiduciary; parental notice required | DPDPA 2023 § 9 (Processing children’s data)            | GDPR Art. 8 (Child consent)            |
| Data Localisation        | All footage stored on India-based servers                    | DPDPA 2023 § 16 (Cross-border restriction)             | N/A (India-specific)                   |
| Access Control           | Role-based access for authorised staff only                  | DPDPA 2023 § 8(5) (Security obligations)               | GDPR Art. 32 (Security)                |
| Algorithmic Transparency | HITL verification; explainable alert metadata                | DPDPA 2023 § 12 (Right to information)                 | GDPR Art. 22 (Automated decisions)     |
| Bias Monitoring          | Annual audit; diverse training data requirement              | OECD AI Principles 2019                                | UNESCO AI Ethics 2021                  |

Under DPDPA 2023 § 9, schools act as Data Fiduciaries responsible for ensuring lawful processing of children’s data. Parents and guardians must be notified of the surveillance system’s existence, purpose, and data handling practices through a clear consent notice. The system may not retain identifiable footage beyond the 72-hour automatic deletion window except for confirmed incidents retained under § 8(7) for legitimate educational safety purposes.

*B. Fairness and Bias Mitigation*

AI-driven decision systems in educational environments risk unintentional bias against student subgroups based on behaviour patterns, physical attributes, or environmental factors [31]. Mitigation measures include: training datasets must include diverse school settings (urban, rural, varied lighting), annual bias audits comparing false positive rates across demographic subgroups, and periodic model retraining as student populations evolve [27], [31].

### C. Consent and Transparency

Parents, students, and educators must be clearly informed about how surveillance data are collected, processed, and used to generate alerts [30]. The school must publish an accessible Plain Language Summary of the AI monitoring system and establish a formal mechanism for parents to request human review of any AI-flagged incident involving their child. Transparent communication builds trust and supports legal compliance [28], [30].

### D. Legal Compliance

Beyond DPDPA 2023, the framework aligns with: the OECD AI Principles 2019 (transparency, robustness, accountability), UNESCO Recommendation on the Ethics of Artificial Intelligence 2021 [34], and relevant provisions of the Protection of Children from Sexual Offences Act (POCSO) 2012 regarding data involving minors. Collaboration with school authorities, legal advisors, and state education regulators is necessary to establish clear data retention, access control, and incident response policies [26], [29].

## VI. DISCUSSION

### A. Challenges

One of the primary challenges is the limited availability of labelled datasets representing authentic school bullying scenarios. Ethical and privacy restrictions make collecting and annotating genuine bullying footage difficult, potentially affecting model generalisation in real deployment settings [22], [23]. Models trained on generic violence datasets (RWF-2000, Hockey Fight) may not fully capture the nuanced behaviours of school-age children.

Environmental variability presents another challenge. Surveillance footage in schools suffers from inconsistent lighting, occlusions in crowded spaces, and camera angle limitations. These factors complicate accurate motion and interaction feature extraction, increasing false positive rates [20], [21]. Designing context-aware models capable of adapting to diverse school layouts is essential. Real-time operational constraints also affect performance. Many schools rely on legacy CCTV infrastructure and limited computational resources. Ensuring low-latency processing while maintaining detection accuracy remains critical, particularly given the budgetary constraints common in Indian educational institutions [23]. Edge-based lightweight neural networks (e.g., MobileNet V2) offer a viable path for resource-constrained environments [33].

### B. Future Directions

Future research should prioritise multimodal detection frameworks integrating vision-based analysis with complementary sensor data. Combining CCTV footage with accelerometer and gyroscope signals from wearable devices or smartphones can capture fine-grained motion signatures such as sudden impacts or forceful movements invisible to cameras [14]. Such integration is particularly valuable in occluded environments where visual-only systems show reduced reliability.

A second promising direction is the extension of AI models beyond physical motion analysis to include emotional stress indicators and vocal aggression. Advances in affective computing and BiLSTM-based prosody analysis suggest that vocal stress, raised voices, and distress signals could serve as early warning indicators even before physical contact occurs [24], [29].

Lightweight neural network architectures and edge-based AI deployment are expected to enable real-time processing on cost-constrained school hardware without sacrificing detection performance [33]. Continued collaboration among researchers, educators, policymakers, and ethics experts will be essential to ensure effective, transparent, and socially responsible deployment.

## VII. CONCLUSION AND FUTURE WORK

This paper has presented and theoretically validated an AI-enabled Early Warning System for real-time physical bullying detection in school environments. The proposed framework integrates CNN-based spatial feature extraction, LSTM-based temporal behaviour modelling, skeleton-based privacy-preserving pose estimation, and a formalised Human-in-the-Loop verification workflow within a legally compliant architecture grounded in India's Digital Personal Data Protection Act 2023.

Performance benchmarking against comparable literature-reported implementations indicates target detection performance of approximately 93% precision and 90% recall ( $F1 = 0.91$  at  $\theta^* = 0.65$ ), with alert latency below 3 seconds in hybrid edge-cloud deployment. The threshold sensitivity analysis provides practical calibration guidance for diverse school risk profiles, while the HITL decision matrix ensures that automated detection is always subject to trained human oversight before any intervention is enacted.

The framework directly addresses three research gaps identified in existing literature: school-specific contextualisation of detection models, architecture-level privacy-by-design, and India DPDPA 2023 compliance mapping.

These contributions provide a technically sound and ethically rigorous foundation for translating AI-enabled bullying detection from research into real-world Indian school deployments.

Future work will focus on: (1) constructing a school-specific annotated dataset in collaboration with Indian educational institutions, (2) empirical evaluation of the complete pipeline under real deployment conditions, (3) multimodal extension integrating acoustic and wearable sensor streams, and (4) development of a lightweight MobileNet-based variant for edge deployment in legacy-infrastructure schools. Continued interdisciplinary collaboration between AI researchers, school administrators, child psychologists, legal experts, and policymakers will remain essential to ensuring that these systems serve their fundamental purpose: making every school a demonstrably safer place for every student.

### VIII. ACKNOWLEDGMENT

The authors express sincere gratitude to all school administrators, faculty members, and technical staff who provided valuable insights into real-time surveillance challenges and safety requirements within educational environments. The authors also acknowledge colleagues and research mentors for guidance in computer vision, deep learning, and intelligent surveillance system design.

### REFERENCES

- [1] S. A. Putri, A. Rifai, and I. Nawawi, "Development of an intelligent violence detection system for bullying monitoring using deep learning models," *Journal of Scientific and Applied Informatics*, vol. 7, no. 2, Jun. 2024. doi: 10.36085/jsai.v7i2.6451.
- [2] S. G. Satpute and G. T. Rajeshwar, "Real-time AI solutions for monitoring and preventing bullying in educational institutions," *International Journal of Advanced Scientific Research*, vol. 10, no. 3, pp. 48–51, Aug. 2025.
- [3] TechTRP Editorial Board, "AI can help combat disciplinary and bullying cases in schools," *TechTRP*, Feb. 2024.
- [4] "Takeaways from our investigation on AI-powered school surveillance," *AP News*, Nov. 2024.
- [5] "Can AI identify safety threats in schools? One district wants to try," *The Washington Post*, Jun. 2025.
- [6] A. Nurbek and A. Altayeva, "Comparative evaluation of machine learning methods for bullying detection in surveillance footage," *European Research Materials*, no. 9, pp. 1–14, 2025. doi:10.36085/erm.v0i9.5777.
- [7] L. Siddique et al., "Analysis of real-time hostile activity detection from spatiotemporal features using deep CNNs, RNNs and attention-based mechanisms," *arXiv Preprint*, 2023.
- [8] I.-A. Haiura and A. Iftene, "Detecting violence in videos using convolutional neural networks," *Procedia Computer Science*, vol. 240, pp. 465–475, 2024.
- [9] Natha, S., Ahmed, F., Siraj, M., et al., "Deep BiLSTM attention model for spatial and temporal anomaly detection in video surveillance," *Sensors*, vol. 25, no. 1, 251, 2025. doi:10.3390/s25010251.
- [10] Lee, S., et al., "Unified video anomaly detection model for detecting different anomaly types," *Proceedings of IEEE/CVF WACV 2026*.
- [11] [MDPI], "Spatially time-based robust tracking and re-identification of kindergarten students: A hybrid deep learning framework combining YOLOv8n and ViT," *Journal of Imaging*, vol. 12, no. 4, 150, 2026. doi:10.3390/jimaging12040150.
- [12] G. Orru et al., "Development of technologies for the detection of (cyber)bullying: The BullyBuster project," *Information*, vol. 14, no. 8, 430, 2023. doi:10.3390/info14080430.
- [13] Mazhar AA, Zada I, et al., "AI-powered detection of cyberbullying in short-form video content: A hybrid deep learning framework," *PLoS One*, vol. 21, no. 2, e0338799, 2026.
- [14] V. Gattulli et al., "Human activity recognition for the identification of bullying and cyberbullying using smartphone sensors," *Electronics*, vol. 12, no. 2, 261, 2023. doi:10.3390/electronics12020261.
- [15] T. Hassner, Y. Itcher, and O. Kliper-Gross, "Violent flows: Real-time detection of violent crowd behavior," in *IEEE CVPRW 2012*, pp. 1–6.
- [16] "Violence detection in surveillance videos with deep network using transfer learning," *IEEE Xplore*, 2019. doi:10.1109/ICDM.2019.8910041.
- [17] A. Ullah et al., "Action recognition in video sequences using deep bi-directional LSTM with CNN features," *IEEE Access*, vol. 6, pp. 1155–1166, 2018.
- [18] J. Li et al., "Efficient violence detection using 3D convolutional neural networks," in *IEEE AVSS 2019*, pp. 1–8.
- [19] W. Sultani, C. Chen, and M. Shah, "Real-world anomaly detection in surveillance videos," in *IEEE/CVF CVPR 2018*, pp. 6479–6488.
- [20] S. Sharma et al., "A fully integrated violence detection system using CNN and LSTM," *IJECE*, vol. 11, no. 4, pp. 3374–3380, 2021.
- [21] M. Inayathulla and K. Rajasekhara Rao, "Enhancing real-time violence detection in video surveillance using hybrid deep learning model," *JOWUA*, vol. 16, no. 1, pp. 344–361, 2025.
- [22] M. Cheng, K. Cai, and M. Li, "RWF-2000: An open large scale video database for violence detection," *arXiv:1911.05913*, 2019.
- [23] "Intelligent video surveillance violence detection model with MobileNet V2 and LSTM," *PMC*, 2025.
- [24] B. Zajime, "Ethical AI in schools: Balancing automation, privacy, and human oversight," *WJAETS*, vol. 15, no. 1, pp. 924–934, 2025.
- [25] UNESCO, "Recommendation on the Ethics of Artificial Intelligence," Paris, France, 2021.
- [26] UNESCO, "Recommendation on the Ethics of Artificial Intelligence: Safeguarding privacy and personal rights," 2024.
- [27] "AI adoption in education and associated policy frameworks," *OECD Policy Survey on School Education in the Digital Age*, 2025.
- [28] Y. Yan et al., "A systematic review of AI ethics in education: privacy, fairness, transparency and governance," *Educ. Inf. Technol.*, 2025.
- [29] M. Campbell et al., "Investigation of the privacy concerns in AI systems for young digital citizens," *arXiv:2501.13321*, 2025.
- [30] S. Muigai et al., "Enhancing public safety through advanced video analysis: A Conv-LSTM-SVM model for violence detection," *EAJIT*, vol. 7, no. 1, pp. 1–17, 2025.
- [31] M. Cheng, K. Cai, and M. Li, "RWF-2000: An open large scale video database for violence detection," *arXiv:1911.05913*, 2019.



- [32] "Efficient violence detection in surveillance," PubMed Central, 2025.
- [33] Altaf Hussain, "Detection and recognition of real-time violence and human actions recognition using lightweight MobileNet model," JIAP, vol. 1, no. 3, pp. 125–146, 2025.
- [34] UNESCO, Recommendation on the Ethics of Artificial Intelligence, Paris, France, 2021.
- [35] "AI in school surveillance systems and human rights," AI Values, 2025.
- [36] Ministry of Law and Justice, Government of India, "The Digital Personal Data Protection Act, 2023," Gazette of India, No. 60, Aug. 2023. [Newly added]



10.22214/IJRASET



45.98



IMPACT FACTOR:  
7.129



IMPACT FACTOR:  
7.429



# INTERNATIONAL JOURNAL FOR RESEARCH

IN APPLIED SCIENCE & ENGINEERING TECHNOLOGY

Call : 08813907089  (24\*7 Support on Whatsapp)