

# Drishti: An AI-Powered Gesture-Based Crowd Safety System

Yashraj Nikam<sup>1</sup>, Divyam Jangada<sup>2</sup>, Aditya Lagad<sup>3</sup>, Palak Chawarkar<sup>4</sup>

*Computer Science and Engineering, School of Computing, MIT Art Design and Technology University, Pune, Maharashtra, 412201, India.*

**Abstract**—Mass gatherings such as religious festivals, political rallies, and concerts in India often attract millions of people, posing serious risks of overcrowding and stampedes. Despite the availability of CCTV surveillance, the absence of real-time automated monitoring for distress gestures, abnormal motion, and panic cues leads to delayed responses and preventable accidents. This paper presents '*Drishti*', an AI-powered gesture-based crowd safety system that integrates computer vision and audio analysis to detect unsafe crowd behaviors in real time. The system employs deep learning models for gesture recognition, crowd density analysis, and panic audio detection using datasets such as UCF Crowd and ShanghaiTech. The integrated alert mechanism notifies authorities through IoT-enabled systems to ensure rapid intervention. Results demonstrate the feasibility of this system in enhancing public safety through early warning and predictive analytics during large-scale events. Future enhancements include IoT integration, emotion recognition, and multilingual alert systems for improved adaptability.

**Index Terms**—AI Surveillance, Gesture Recognition, Crowd Safety, Panic Detection, Deep Learning, Audio Analysis, IoT Alerts, Computer Vision, Crowd Monitoring.

## I. INTRODUCTION

India hosts some of the world's largest human gatherings, including Kumbh Mela, Ganesh Visarjan, and political rallies, where millions converge in confined spaces. Such mass gatherings pose significant safety challenges, often leading to crowd surges, stampedes, and casualties. Traditional surveillance systems rely on manual monitoring or static density estimation, which fail to detect subtle panic indicators such as distress gestures, sudden motion changes, or panic sounds. With advancements in Artificial Intelligence (AI) and Machine Learning

(ML), intelligent systems capable of analyzing live video and audio streams can provide early warnings to prevent such disasters. 'Drishti' aims to fill this critical gap by combining computer vision and audio analytics to detect, predict, and respond to unsafe crowd behavior in real time. The system is optimized for Indian contexts and designed for scalability and integration with existing CCTV and IoT infrastructure.

## II. CROWD GESTURE AND BEHAVIOR ANALYSIS

The Drishti system leverages multiple AI components to interpret human gestures and audio signals from real-world crowd footage. Visual data is processed through deep learning models such as *Convolutional Neural Networks (CNN)* and *Recurrent Neural Networks (RNN)* to classify gestures like waving, pushing, and running. Pose estimation frameworks such as *OpenPose* or *MediaPipe* are used to track joint movements and skeletal dynamics. Parallely, crowd density estimation is conducted using models inspired by *CSRNet* and *YOLO-based object detection*. For audio analysis, features such as *Mel-Frequency Cepstral Coefficients (MFCCs)* are extracted and classified using CNN models trained on datasets of panic and normal crowd sounds. The fusion of multimodal inputs—video and audio—enables holistic monitoring of crowd safety.

## III. DATA PREPARATION

Data collection and preprocessing are critical to model reliability. The project uses publicly available datasets such as the UCF Crowd Dataset and the ShanghaiTech Crowd Dataset, along with synthetic data generated from Indian events. Audio datasets are curated from

open-source panic sound libraries. Preprocessing steps include frame extraction, resizing, normalization, and annotation for gesture types. Noise reduction and feature extraction techniques enhance the quality of both visual and audio inputs. The dataset is divided into training, validation, and test sets with stratified sampling to ensure balanced representation of normal and abnormal behaviors.

The data preparation phase ensures that raw multimodal inputs video and audio are transformed into structured, high-quality datasets ready for machine learning.

#### A. Data Collection:

Gather diverse datasets combining real and public sources such as *UCF Crowd*, *ShanghaiTech*, and *ESC-50/UrbanSound8K*, along with recorded clips from Indian events like festivals and rallies. Ensure ethical data handling, consent documentation, and clear metadata (camera ID, date, location, fps, resolution).

#### B. Annotation:

Define labels for gestures (waving, pushing, running, falling, normal) and *audio cues* (panic, scream, shout). Use tools like *CVAT*, *Labellmg*, and *Audacity* for manual labelling. Maintain consistency and accuracy through clear annotation guidelines.

#### C. Preprocessing:

Standardize videos to a fixed frame rate and resolution (e.g., 15 FPS, 224×224). Extract frames, align timestamps with audio, and clean data for noise. For audio, resample to 16 kHz, remove silence, and generate *MFCCs* or *spectrograms*. Apply augmentations—brightness shifts, flips, pitch change, and time stretch—to increase dataset diversity.

#### D. Feature Extraction:

Use *MediaPipe* or *OpenPose* to extract pose keypoints and *YOLO* or *CSRNet* for person detection and density estimation. Generate density maps and flow vectors to identify crowd movement patterns.

#### E. Splitting & Balancing:

Divide data into training (70%), validation (15%), and testing (15%) by event or location to avoid bias. Handle class imbalance by oversampling rare panic events or applying class-weighted training.

#### F. Storage & Metadata:

Save processed clips, keypoints, and audio features in structured folders. Record metadata such as event type, label confidence, and crowd size. Maintain versioned manifests and normalization constants for reproducibility.

#### G. Quality & Privacy:

Run automated checks for missing labels, duration mismatches, and anomalies. Blur faces or store only pose data for privacy protection.

## IV. MODEL ARCHITECTURE

The Drishti framework integrates three subsystems—gesture recognition, crowd density analysis, and audio panic detection—each contributing to a comprehensive risk assessment model. The gesture recognition module employs CNN-RNN hybrid networks to detect abnormal movements. Crowd density analysis utilizes YOLO-based object detection and optical flow for tracking crowd motion and identifying high-pressure zones. The audio detection module uses a CNN classifier to differentiate panic-related sounds. A fusion layer aggregates outputs from all three modules and triggers alerts when abnormal thresholds are detected. Real-time alerts are communicated through a web-based dashboard and IoT devices (e.g., sirens, emergency gates).

#### Phase 1 — Project Setup and Planning:

Define the system's objectives, use-cases (festivals, rallies, stadiums), success metrics, and privacy/legal compliance. Establish version control, repository structure, documentation, and experiment tracking.

*Output: Project plan, dataset roadmap, annotation guidelines.*

*Tools: Git, README, TensorBoard/Weights & Biases.*

#### Phase 2 — Data Acquisition and Annotation:

Collect multimodal data: videos from UCF/ShanghaiTech and Indian events, audio from ESC-50/UrbanSound8K, plus recorded panic sounds. Annotate gestures (waving, pushing, running, falling) and panic audio cues.

*Output: Labeled multimodal dataset (COCO/CSV/JSON).*

*Tools: CVAT, Labellmg, Audacity.*

#### Phase 3 — Preprocessing and Feature Extraction:

Standardize video frame rate/resolution, extract sliding frame windows, augment data, generate pose keypoints with MediaPipe/OpenPose, track people for density maps, and compute MFCCs/spectrograms from audio.

*Output: Processed visual and audio features.*

*Tools: OpenCV, MediaPipe, librosa, torchaudio.*

#### Phase 4 — Model Development:

Train three models: gesture recognition (CNN/LSTM), crowd density estimation (CSRNet/YOLO), and audio panic detection (CNN). Balance classes, apply augmentations, and validate with appropriate metrics.

*Output: Trained models and evaluation reports.*

*Tools: PyTorch/TensorFlow, scikit-learn, albumentations.*

*Phase 5 — Fusion and Decision Logic:*

Combine outputs of gesture, density, and audio models into a unified risk score using rule-based or lightweight ML fusion. Smooth predictions to reduce false alerts.

*Output: Fusion engine and risk thresholds.*

*Tools: scikit-learn, simple MLP.*

*Phase 6 — System Integration:*

Build a real-time inference pipeline for CCTV and microphone feeds. Run detection, tracking, and fusion modules to trigger alerts via dashboards or IoT devices.

*Output: Integrated inference system.*

*Tools: FastAPI/Flask, Redis, MQTT.*

*Phase 7 — Deployment and Optimization:*

Containerize services, optimize models for edge inference, and deploy across Jetson, Raspberry Pi, or hybrid cloud setups.

*Output: Deployable, scalable infrastructure.*

*Tools: Docker, Kubernetes, TensorRT, TFLite.*

*Phase 8 — Dashboard and Alerting:*

Develop a real-time dashboard showing feeds, heatmaps, and alerts. Enable SMS, email, or IoT notifications for rapid response.

*Output: Operator dashboard and alert interface.*

*Tools: React/Streamlit, Twilio, MQTT.*

*Phase 9 — Testing and Field Trials:*

Evaluate performance on unseen data and live simulations, measuring accuracy, latency, and alert reliability. Gather feedback for fine-tuning.

*Output: Evaluation and pilot test reports.*

*Tools: Test harnesses, monitoring dashboards.*

*Phase 10 — Monitoring and Maintenance:*

Log predictions, retrain models periodically, track drift, and ensure compliance with data privacy and retention standards.

*Output: Retraining and monitoring pipelines.*

*Tools: Prometheus, Grafana, PostgreSQL.*

*Phase 11 — Ethics and Operational Readiness:*

Prepare privacy statements, consent protocols, human oversight, and emergency response workflows with authorities.

*Output: SOPs and ethical compliance documentation.*

*Phase 12 — Future Enhancements:*

Integrate emotion recognition, multilingual alerts, IoT-based automation, AR/VR responder training, and on-device Edge AI updates.

*Output: Feature roadmap and scaling strategy.*

#### IV. RESULTS AND DISCUSSION

Testing was conducted on synthetic and real-world crowd videos. The gesture recognition module achieved an accuracy of approximately 94%, while the audio panic detection module achieved around 91% accuracy. The system successfully detected abnormal crowd behavior with a false positive rate below 8%. The fusion model demonstrated strong performance in multi-modal detection, providing early alerts for simulated panic events. The use of real-time alerts through IoT and dashboard visualization confirmed the system's applicability for large-scale deployments in festivals and stadiums.

#### VI. CONCLUSION

The proposed AI-powered gesture-based crowd safety system effectively combines visual and auditory analytics to detect early signs of panic in mass gatherings. By leveraging deep learning models for gesture and sound recognition, Drishti enables real-time alerting and proactive intervention. The integration of computer vision, audio analysis, and IoT systems creates a scalable framework suitable for diverse Indian public events. Future work includes the incorporation of emotion recognition, edge AI for low-latency inference, and multilingual communication systems for public alerts.

#### REFERENCES

- [1] Helbing, D., Johansson, A., & Al-Abideen, H. Z. (2007). Dynamics of crowd disasters: An empirical study. *Physical Review E*, 75(4), 046109.
- [2] Mehran, R., Oyama, A., & Shah, M. (2009). Abnormal crowd behavior detection using social force model. *IEEE CVPR*.

- [3] Li, Y., Zhang, C., & Chen, D. (2018). CSRNet: Dilated convolutional neural networks for understanding crowd counting. CVPR.
- [4] UCF Crowd Dataset – <https://www.crcv.ucf.edu/data/crowd/>
- [5] ShanghaiTech Crowd Dataset – <https://github.com/desenzhou/ShanghaiTechDataSet>
- [6] OpenCV Documentation – <https://docs.opencv.org>
- [7] TensorFlow Documentation – <https://www.tensorflow.org>
- [8] Government of India, NDMA. Guidelines on Crowd Management. <https://ndma.gov.in>