

Human Activity Recognition

By: Aditi Srivastava

Overview

Human Activity Recognition (HAR) involves identifying physical movements of humans such as walking, running, sitting, jumping, and waving. This project uses pose estimation and deep learning to classify actions in real-time using webcam input.

Problem Statement

The goal is to detect and classify human activities such as walking, jogging, running, hand-waving, boxing, and clapping from video frames. Traditional CNN approaches struggle because they analyze single images and fail to understand motion.

Activities:

- Walking
- Running
- Boxing
- Hand Waving, Jogging, Clapping



Why CNN Fails

CNN processes frames independently and cannot capture sequential motion. Even averaging predictions across frames produces inaccurate classification for actions requiring temporal understanding.

CNN Limitation:

Single-frame analysis → No motion understanding

Training Pipeline

1. Capture video frames from webcam.
2. Extract body keypoints using MediaPipe Pose.
3. Store extracted keypoints into CSV files.
4. Build dataset containing sequences of poses.
5. Train an LSTM deep learning model.
6. Save the trained model for inference.

Training Pipeline Diagram

Video -> Pose Extraction -> CSV Dataset





CSV -> LSTM Training -> Saved Model

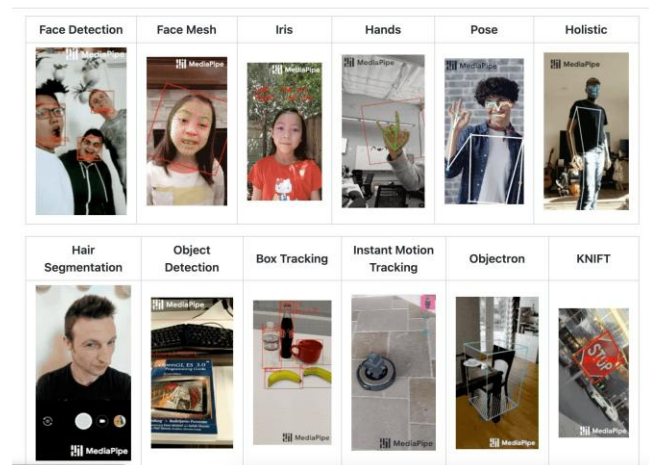
Real-time ready model

MediaPipe

Live ML anywhere

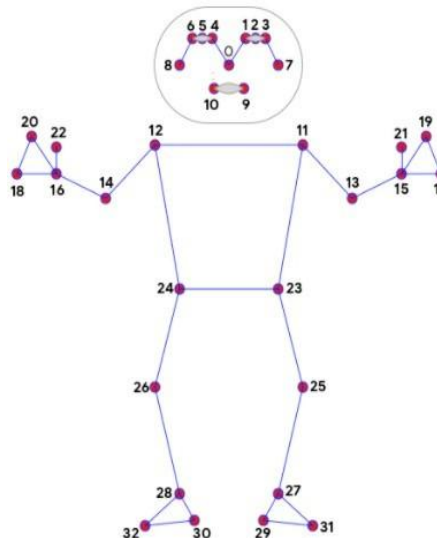
MediaPipe offers cross-platform, customizable ML solutions for live and streaming media.

	
<i>End-to-End acceleration: Built-in fast ML inference and processing accelerated even on common hardware</i>	<i>Build once, deploy anywhere: Unified solution works across Android, iOS, desktop/cloud, web and IoT</i>
	
<i>Ready-to-use solutions: Cutting-edge ML solutions demonstrating full power of the framework</i>	<i>Free and open source: Framework and solutions both under Apache 2.0, fully extensible and customizable</i>



MediaPipe Pose

MediaPipe Pose provides 33 body landmarks with real-time performance. It extracts skeleton keypoints that represent human posture and motion essential for understanding activities.



- | | |
|--------------------|----------------------|
| 0. nose | 17. left_pinky |
| 1. left_eye_inner | 18. right_pinky |
| 2. left_eye | 19. left_index |
| 3. left_eye_outer | 20. right_index |
| 4. right_eye_inner | 21. left_thumb |
| 5. right_eye | 22. right_thumb |
| 6. right_eye_outer | 23. left_hip |
| 7. left_ear | 24. right_hip |
| 8. right_ear | 25. left_knee |
| 9. mouth_left | 26. right_knee |
| 10. mouth_right | 27. left_ankle |
| 11. left_shoulder | 28. right_ankle |
| 12. right_shoulder | 29. left_heel |
| 13. left_elbow | 30. right_heel |
| 14. right_elbow | 31. left_foot_index |
| 15. left_wrist | 32. right_foot_index |
| 16. right_wrist | |

Why LSTM?

LSTM (Long Short-Term Memory) networks are excellent for sequential data. Human activities rely on how movement changes across frames. LSTM learns motion patterns over time, making it ideal for HAR.

Why LSTM Works:

- Captures sequential motion
- Remembers long-term patterns
- More accurate than CNN for actions

Inference Pipeline

1. Read frame from webcam
2. Extract pose landmarks using MediaPipe
3. Structure keypoints into a sequence
4. Feed the sequence into the trained LSTM
5. Predict action and show label on screen

Inference Pipeline

Webcam Frame → Pose Extractor

Pose Landmarks → LSTM Model

Prediction → Display on Screen

Demo Output

The system displays the detected skeleton and predicted activity (e.g., Walking, Running, Jumping) in real time. This makes it suitable for monitoring, safety, fitness, and robotics applications.