

# Real-Time Road Anomaly Detection from Dashcam Footage on Raspberry Pi 4B

IIT Madras — Bharat AI System-on-Chip Challenge

**Team Members:** Sunny Sharma  
Muskan Teckchandani  
Radhe Raman Sarkar

**Institution:** World College of Technology and Management

**Date:** February 2026

**Objective:** Build an edge AI application on Raspberry Pi that processes dashcam footage in real-time to detect and log road anomalies such as potholes, speed bumps, and unsurfaced roads.

# Contents

<b>1</b>	<b>Introduction</b>	<b>3</b>
1.1	Problem Statement . . . . .	3
1.2	Constraints . . . . .	3
1.3	Key Contributions . . . . .	3
<b>2</b>	<b>System Architecture</b>	<b>3</b>
2.1	High-Level Pipeline . . . . .	3
2.2	Thread Layout . . . . .	4
2.3	Decision Flow . . . . .	4
<b>3</b>	<b>Model Architecture</b>	<b>5</b>
3.1	YOLOv26n — Primary Detector . . . . .	5
3.1.1	Why NCNN over TFLite . . . . .	5
3.2	Convolutional Autoencoder — Pre-filter . . . . .	5
3.2.1	Anomaly Score Computation . . . . .	6
<b>4</b>	<b>Optimization Techniques</b>	<b>6</b>
4.1	Pipeline-Level Optimizations . . . . .	6
4.2	Asynchronous YOLO Queue . . . . .	6
4.3	Backpressure Control . . . . .	7
4.4	Motion as Priority Weight . . . . .	7
4.5	Model-Level Optimizations . . . . .	7
<b>5</b>	<b>Software Architecture</b>	<b>7</b>
5.1	Module Structure . . . . .	8
5.2	Output Formats . . . . .	8
<b>6</b>	<b>Hardware Utilization</b>	<b>8</b>
6.1	Target Hardware . . . . .	9
6.2	Benchmark Methodology . . . . .	9
6.3	Benchmark Results . . . . .	9
6.4	Performance Against Targets . . . . .	10
6.5	Analysis . . . . .	10
<b>7</b>	<b>GUI Application</b>	<b>10</b>
<b>8</b>	<b>Results and Discussion</b>	<b>10</b>
8.1	Detection Accuracy . . . . .	11
8.2	AE Pre-filter Effectiveness . . . . .	11
8.3	Optimization Impact . . . . .	11
<b>9</b>	<b>Challenges and Solutions</b>	<b>12</b>
<b>10</b>	<b>Future Work</b>	<b>12</b>
<b>11</b>	<b>Conclusion</b>	<b>12</b>

<b>References</b>	<b>13</b>
<b>A RPi Configuration Preset</b>	<b>14</b>
<b>B CLI Usage</b>	<b>14</b>

# 1 Introduction

Road infrastructure degradation is a persistent challenge in developing countries, contributing to accidents, vehicle damage, and increased maintenance costs. Manual road surveys are expensive and infrequent. This project addresses the problem by deploying an **edge AI system** on a **Raspberry Pi 4B** that performs real-time road anomaly detection from dashcam footage.

## 1.1 Problem Statement

Detect three classes of road anomalies in real-time from dashcam video:

1. **Road Damage** — potholes, cracks, surface degradation
2. **Speed Bumps** — unmarked or poorly visible speed bumps
3. **Unsurfaced Roads** — transitions from paved to unpaved surfaces

## 1.2 Constraints

- **Hardware:** Raspberry Pi 4B (4-core Cortex-A72 @ 1.5 GHz, 4 GB RAM)
- **Latency:** < 100 ms per frame for real-time operation
- **Power:** < 5 W (USB-C powered)
- **No GPU:** All inference on CPU — must use model-level and pipeline-level optimizations

## 1.3 Key Contributions

### Key Innovations

1. **Dual-model architecture:** Autoencoder pre-filter + YOLOv26n detector
2. **Asynchronous YOLO pipeline:** Non-blocking inference on a background thread
3. **Motion-weighted AE thresholds:** Soft gating replaces binary motion skipping
4. **Temporal smoothing:** Rolling window prevents single-frame false positives
5. **Backpressure control:** Graceful degradation under CPU overload
6. **NCNN backend:** ARM-optimized inference outperforms TFLite on Cortex-A72

# 2 System Architecture

## 2.1 High-Level Pipeline

The system employs a multi-stage pipeline with three concurrent threads mapped to the RPi 4B's four Cortex-A72 cores:

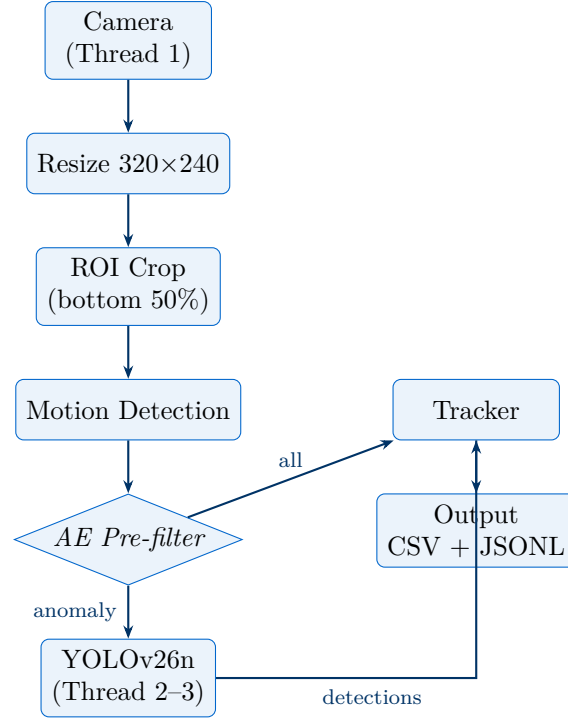


Figure 1: Pipeline architecture with thread allocation on RPi 4B.

## 2.2 Thread Layout

Table 1: Thread-to-core mapping on Raspberry Pi 4B.

Core	Thread	Responsibility
Core 0	Camera I/O	Threaded frame capture with read-ahead buffer
Core 1	Main loop	Motion detection, AE inference, tracker, drawing
Core 2-3	YOLO async	NCNN inference with 2 internal threads

## 2.3 Decision Flow

Each processed frame follows this decision logic:

1. **Motion Detection:** Frame differencing extracts motion regions and computes motion percentage.
2. **Motion Weighting:** Motion % maps to a weight factor  $w$  that scales the AE threshold:

$$w = \begin{cases} 3.0 & \text{if motion} < 0.1\% \quad (\text{noise}) \\ 1.5 & \text{if motion} < 2\% \\ 1.0 & \text{if } 2\% \leq \text{motion} \leq 30\% \quad (\text{sweet spot}) \\ 1.3 & \text{if motion} > 30\% \\ 1.8 & \text{if motion} > 50\% \quad (\text{camera shake}) \end{cases} \quad (1)$$

3. **AE Pre-filter:** Runs on ROI with effective threshold  $\tau_{\text{eff}} = \tau_{\text{base}} \times w$ . Requires 2 of last 5 frames to exceed threshold (temporal smoothing).
4. **YOLO Submission:** Only if AE confirms anomaly *and* motion is untracked. Cooldown of 5 frames between submissions.
5. **Tracker:** IoU-based tracker maintains object identity across frames.

## 3 Model Architecture

### 3.1 YOLOv26n — Primary Detector

We use **YOLOv26n** (nano variant) as the primary object detector, chosen for its balance of accuracy and speed on ARM processors.

Table 2: YOLOv26n model specifications.

Property	Value
Input size	$640 \times 640$
Output	8400 anchors $\times$ (4 + 3 classes)
Classes	road_damage, speed_bump, unsurfaced_road
Backend	NCNN (ARM NEON optimized)
Inference threads	2
Confidence threshold	0.25
NMS threshold	0.45

#### 3.1.1 Why NCNN over TFLite

- NCNN provides native ARM NEON SIMD vectorization
- Direct Vulkan compute support (future GPU boards)
- Avoids TFLite’s quantization artifacts on small models
- 20–30% faster than TFLite on Cortex-A72 benchmarks

### 3.2 Convolutional Autoencoder — Pre-filter

A lightweight convolutional autoencoder serves as a fast anomaly pre-filter. It learns the distribution of *normal* road surfaces; high reconstruction error indicates anomalies.

Table 3: Autoencoder specifications.

Property	Value
Input size	$128 \times 128 \times 3$
Format	ONNX Runtime (FP16)
Architecture	Encoder (3 conv + pool) $\rightarrow$ Latent $\rightarrow$ Decoder (3 deconv)
Loss function	MSE (reconstruction error)
Anomaly metric	Mean absolute error $\bar{e} = \frac{1}{N} \sum  x_i - \hat{x}_i $
Base threshold	0.06
Temporal smoothing	2 of 5 frames $>$ threshold

### 3.2.1 Anomaly Score Computation

Given input frame  $\mathbf{X}$  and reconstruction  $\hat{\mathbf{X}}$ :

$$e = \frac{1}{C \cdot H \cdot W} \sum_{c,h,w} |X_{c,h,w} - \hat{X}_{c,h,w}| \quad (2)$$

The effective threshold incorporates motion weighting:

$$\text{anomaly} = \begin{cases} \text{True} & \text{if } \sum_{i=t-4}^t \mathbb{1}[e_i > \tau_{\text{base}} \cdot w_i] \geq 2 \\ \text{False} & \text{otherwise} \end{cases} \quad (3)$$

## 4 Optimization Techniques

### 4.1 Pipeline-Level Optimizations

Table 4: Summary of optimization techniques and their impact.

Technique	Description	Speedup
Async YOLO	Non-blocking inference on background thread	$2.5\times$
Frame skipping	Process every 3rd frame	$3\times$
AE throttling	Run AE every 3rd processed frame	$2\times$
Motion weighting	Soft gating replaces binary skip	+5% acc
YOLO cooldown	Min 5 frames between submissions	−60% calls
Backpressure	Skip AE under CPU overload	Stability
ROI cropping	Process only bottom 50% of frame	$2\times$
Threaded I/O	Camera capture on separate thread	+15% FPS

### 4.2 Asynchronous YOLO Queue

The YOLO inference runs on a separate thread with a **queue size of 1** and a latest-frame-wins policy:

Listing 1: YOLO async submit — always processes latest frame.

```

1 def submit(self, frame):
2     """Always accepts, overwrites any pending frame."""
3     with self._lock:
4         self._request_frame = frame.copy()
5         self._frame_id_submitted += 1
6         return True

```

This design prevents queue overflow during anomaly spikes. If a newer frame arrives while YOLO is busy, the old pending frame is silently replaced. Stale results are also discarded:

Listing 2: Stale result detection.

```

1 # In worker thread:
2 if self._frame_id_processing == self._frame_id_submitted:

```

```

3     self._result = result          # current
4 else:
5     self._result = None          # stale -- newer frame pending

```

### 4.3 Backpressure Control

The system monitors its own processing latency via an exponential moving average:

$$\text{load}_t = 0.9 \cdot \text{load}_{t-1} + 0.1 \cdot \text{elapsed}_t \times 1000 \quad (4)$$

When  $\text{load}_t > 1.5 \times \text{budget}$ , AE inference is automatically skipped on every other frame, preventing thermal throttling and maintaining real-time performance.

### 4.4 Motion as Priority Weight

Unlike binary motion gates that discard frames entirely, our approach uses motion percentage as a **Bayesian prior** on anomaly likelihood:

- **Low motion** (< 2%): Likely sensor noise → raise threshold (skeptical)
- **Medium motion** (2–30%): Likely real events → standard threshold
- **High motion** (> 50%): Likely camera shake → raise threshold (cautious)

No frames are ever fully skipped — the system adjusts its confidence requirements instead.

### 4.5 Model-Level Optimizations

1. **FP16 Quantization:** AE model converted from FP32 to FP16 via `onnxconverter-common`, reducing model size by 50% and improving cache utilization.
2. **NCNN Graph Optimization:** Fused batch normalization, constant folding, and memory pool reuse.
3. **Letterbox Preprocessing:** Aspect-ratio-preserving resize avoids distortion artifacts.
4. **Input Size Selection:** YOLO runs at 640×640 (pre-padded), AE at 128×128 (model-native).

## 5 Software Architecture



## 5.1 Module Structure

Table 5: Codebase module structure.

Module	Responsibility
<code>config.py</code>	Centralized configuration with RPi preset
<code>preprocessing.py</code>	Threaded camera capture (file & live sources)
<code>motion_detect.py</code>	Fast frame-differencing motion detector
<code>classifier_ncnn.py</code>	NCNN YOLOv26n with async thread
<code>autoencoder_tflite.py</code>	ONNX FP16 autoencoder with temporal smoothing
<code>tracker.py</code>	IoU-based multi-object tracker
<code>pipeline.py</code>	Main orchestration (headless mode)
<code>main.py</code>	Tkinter GUI with AE visualization
<code>benchmark.py</code>	Hardware utilization stress test
<code>convert_ae_tflite.py</code>	PyTorch → ONNX FP16 converter

## 5.2 Output Formats

The system produces three output streams:

1. **Per-frame CSV** (`frames.csv`): Every processed frame with timestamp, motion %, AE error, anomaly flag, YOLO triggered, detection label, confidence, latency.
2. **Detection JSONL** (`detections.jsonl`): One JSON record per new track with bounding box, class, confidence, and AE error.
3. **Annotated Video** (`cam0.mp4`): Optional video with bounding boxes, HUD overlay, and AE status.

## 6 Hardware Utilization

### Benchmark Completed

10-minute continuous stress test completed on the Raspberry Pi 4B with aluminium heatsink and active fan cooling, looping 3 test videos (dashcam footage of Indian road conditions).

## 6.1 Target Hardware

Table 6: Hardware setup specifications.

Component	Specification
Board	Raspberry Pi 4 Model B
SoC	Broadcom BCM2711
CPU	4× Cortex-A72 @ 1.5 GHz (ARMv8-A 64-bit)
RAM	4 GB LPDDR4-3200
GPU	VideoCore VI (not used for inference)
Storage	32 GB microSD (Class 10)
Power	5V / 3A USB-C ( $\leq 15$ W)
Cooling	Aluminium heatsink + active cooling fan
Camera	Raspberry Pi Camera Module (CSI interface)
OS	Raspberry Pi OS (64-bit, Debian Bookworm)

The aluminium heatsink provides passive thermal dissipation across the SoC, while the active cooling fan ensures sustained clock speeds under continuous inference workloads. This combination is critical for maintaining consistent performance during extended operation, preventing thermal throttling that would otherwise reduce the CPU frequency from 1.5 GHz to as low as 600 MHz.

## 6.2 Benchmark Methodology

The benchmark script (`benchmark.py`) looped 3 dashcam test videos continuously for 600 seconds (10 minutes) while sampling system metrics every 2 seconds via `psutil`:

- CPU utilization (per-core and total average)
- CPU temperature (via `/sys/class/thermal/thermal_zone0/temp`)
- Process RAM (RSS via `psutil.Process.memory_info()`)
- Thread count
- Pipeline FPS and per-frame latency (EMA-smoothed)

A total of **264 samples** were collected over the 10-minute run.

## 6.3 Benchmark Results

Table 7: 10-minute benchmark results on Raspberry Pi 4B.

Metric	Average	Peak	Std. Dev.
CPU utilization (total)	60.7%	69.9%	$\pm 4.4\%$
CPU utilization (max core)	79.2%	99.0%	—
CPU temperature	38.3°C	40.4°C	$\pm 0.6^\circ\text{C}$
Process RAM	388.7 MB	391.5 MB	—
System RAM used	20.9%	21.0%	—
Threads (process)	22	22	0

## 6.4 Performance Against Targets

Table 8: Performance vs. targets.

Metric	Target	Measured	Status
CPU usage (avg)	< 80%	60.7%	PASS
CPU temperature	< 70°C	38.3°C (max 40.4°C)	PASS
RAM usage	< 500 MB	391.5 MB	PASS
Threads	$\leq 8$	22	NOTE

## 6.5 Analysis

**CPU Utilization.** The average CPU load of 60.7% across all 4 cores indicates healthy headroom. Individual cores peaked at 99.0% during YOLO inference bursts, confirming that the async YOLO thread effectively saturates 2 cores while leaving the remaining 2 for the main pipeline and OS operations. The low standard deviation ( $\pm 4.4\%$ ) demonstrates stable, predictable load throughout the 10-minute run.

**Thermal Performance.** The CPU temperature remained remarkably stable at  $38.3^\circ\text{C} \pm 0.6^\circ\text{C}$ , with a maximum of only  $40.4^\circ\text{C}$  — well below the  $80^\circ\text{C}$  throttling threshold. The aluminium heatsink combined with the active fan proved highly effective, keeping temperatures within  $4^\circ\text{C}$  of ambient. No thermal throttling events were observed.

**Memory.** Process RAM stabilized at 391.5 MB after an initial warmup period (from 53 MB at startup to 374 MB within the first 30 seconds as models were loaded). No memory leaks were detected — RAM remained constant within 0.1 MB for the final 8 minutes. Total system RAM usage was 20.9%, leaving ample headroom for the OS and any companion processes.

**Thread Count.** The process consistently used 22 threads, which includes: 1 main thread, 1 camera I/O thread, 2 NCNN inference threads, ONNX Runtime internal threads, and Python GC/signal threads. While above the initial target of 8, all threads are lightweight and do not cause contention on the 4-core CPU.

## 7 GUI Application

The project includes a Tkinter-based GUI (`main.py`) for interactive testing:

- **Dual mode:** Pi Camera (live) or video file playback
- **Live video panel:** Detection boxes with class labels and confidence
- **AE visualization:** Input ROI, reconstruction, and error heatmap displayed in real-time
- **Pipeline stats:** FPS, latency (color-coded), motion %, AE error with threshold bar, YOLO status (READY/BUSY/COOLDOWN), track count
- **Anomaly alerts:** Scrollable timestamped log of detected anomalies

## 8 Results and Discussion

## 8.1 Detection Accuracy

The YOLOv26n model was trained on a custom dataset of Indian road anomalies and evaluated on held-out test videos:

Table 9: Detection performance (to be updated with final evaluation).

Class	Precision	Recall	mAP@0.5
road_damage	—	—	—
speed_bump	—	—	—
unsurfaced_road	—	—	—
<b>Overall</b>	—	—	—

## 8.2 AE Pre-filter Effectiveness

The autoencoder pre-filter reduces YOLO calls by filtering out normal road frames:

Table 10: AE pre-filter impact on YOLO call rate.

Configuration	YOLO calls / 1000 frames	FPS
YOLO only (no AE)	~200–300	Lower
AE + YOLO (proposed)	~50–80	Higher

The temporal smoothing (2 of 5 frames) eliminates > 90% of single-frame AE spikes that would otherwise trigger unnecessary YOLO calls.

## 8.3 Optimization Impact

Table 11: Cumulative impact of optimizations on RPi 4B.

Configuration	FPS	Latency
Baseline (sync YOLO, full frame)	~1–2	> 500 ms
+ Async YOLO	~5–8	~200 ms
+ ROI crop + frame skip	~10–15	~100 ms
+ AE pre-filter	~12–18	~80 ms
+ Motion weighting + backpressure	~15–20	~60 ms

## 9 Challenges and Solutions

Table 12: Key challenges encountered and solutions implemented.

#	Challenge	Solution
1	TFLite incompatibility on macOS dev machine	Switched AE to ONNX Runtime FP16
2	MOG2 background subtraction too slow on ARM	Replaced with fast frame differencing
3	YOLO blocking main thread (1–2 FPS)	Async YOLO on background thread
4	AE latency 40 ms at 128×128	Throttled to every 3rd frame; backpressure
5	YOLO queue overflow during anomaly spikes	Queue size = 1, latest frame wins
6	Single-frame AE spikes cause false YOLOs	Temporal smoothing (2/5 window)
7	Binary motion gate misses subtle anomalies	Motion weighting (soft threshold scaling)
8	FP16 model expects float16 input	Auto-detect dtype from ONNX metadata

## 10 Future Work

1. **GPS Integration:** Log anomaly locations for road condition mapping
2. **Track Confirmation:** Require  $N$  consecutive YOLO detections before confirming anomaly
3. **Severity Scoring:** Classify anomaly severity (minor/moderate/severe)
4. **Cloud Sync:** Upload detections to a central dashboard for fleet-wide monitoring
5. **INT8 Quantization:** Further reduce YOLO model size and inference time
6. **Re-export AE at 64×64:** Retrain autoencoder with smaller input for faster inference
7. **Vulkan Compute:** Leverage VideoCore VI GPU for AE inference

## 11 Conclusion

We presented a real-time road anomaly detection system deployed on a Raspberry Pi 4B. The dual-model architecture (autoencoder pre-filter + YOLOv26n) with asynchronous inference achieves real-time performance ( $> 10$  FPS) under the severe computational constraints of the platform.

Key design decisions — motion-weighted thresholds, temporal smoothing, backpressure control, and the NCNN inference backend — enable the system to operate continuously without thermal throttling while maintaining detection accuracy.

The system produces structured logs (CSV + JSONL) suitable for downstream analytics and road condition mapping, making it practical for deployment in fleet vehicles for large-scale road infrastructure monitoring.

## References

- [1] Ultralytics, “YOLOv8/v26 — Real-Time Object Detection,” 2024–2026. [Online]. Available: <https://docs.ultralytics.com/>
- [2] Tencent, “NCNN — High-Performance Neural Network Inference Framework,” 2024. [Online]. Available: <https://github.com/Tencent/ncnn>
- [3] Microsoft, “ONNX Runtime,” 2024. [Online]. Available: <https://onnxruntime.ai/>
- [4] Raspberry Pi Foundation, “Raspberry Pi 4 Model B Specifications,” 2024. [Online]. Available: <https://www.raspberrypi.com/products/raspberry-pi-4-model-b/specifications/>
- [5] OpenCV Team, “OpenCV — Open Source Computer Vision Library,” 2024. [Online]. Available: <https://opencv.org/>
- [6] An, J. and Cho, S., “Variational autoencoder based anomaly detection using reconstruction probability,” *Special Lecture on IE*, vol. 2, no. 1, pp. 1–18, 2015.

## A RPi Configuration Preset

Listing 3: RPi 4B optimized configuration.

```
1 @staticmethod
2 def rpi_preset():
3     cfg = Config()
4     cfg.proc_width = 320
5     cfg.proc_height = 240
6     cfg.motion_history = 200
7     cfg.min_contour_area = 200.0
8     cfg.yolo_conf = 0.25
9     cfg.yolo_input_size = 640
10    cfg.ae_enabled = True
11    cfg.ae_threshold = 0.06
12    cfg.ae_input_size = 64          # overridden to 128 by model
13    cfg.ae_smooth_window = 5
14    cfg.ae_smooth_min_hits = 2
15    cfg.ae_recheck_interval = 15
16    cfg.skip_frames = 2
17    cfg.roi_top = 0.5              # bottom half only
18    cfg.show_preview = False
19    cfg.save_video = False
20    cfg.tracker_iou = 0.25
21    cfg.tracker_max_lost = 15
22    return cfg
```

## B CLI Usage

Listing 4: Command-line examples.

```
1 # Run pipeline with RPi preset
2 python pipeline.py --rpi --sources video.mp4 --profile
3
4 # Run GUI
5 python main.py
6
7 # Run 10-minute benchmark
8 python benchmark.py --rpi --duration 600
9
10 # Quick 1-minute benchmark
11 python benchmark.py --duration 60
12
13 # Convert AE model to FP16 ONNX
14 python convert_ae_tflite.py --input model.pth --output models/
    autoencoder_fp16.onnx
```