



International Journal of Research Publication and Reviews

Journal homepage: www.ijrpr.com ISSN 2582-7421

A Real-Time AI Framework for Stray Animal Detection, Rescue Assistance, and Automated Reporting Using YOLOv8 and GPT-4o

Mayank Satpute¹, Jyoti Patel²

Oriental Institute of Science and Technology, Madhya Pradesh

¹mayanksatpute56@gmail.com, ²jyp25070@gmail.com

ABSTRACT

In many Indian cities, injured stray animals often remain unnoticed or poorly reported, leading to delayed or failed rescues. This paper presents StrayCare, an AI-driven platform designed to improve detection, reporting, and response coordination. The system integrates real-time visual identification of injured animals using a lightweight YOLOv8 pipeline optimized for mobile via TensorFlow Lite, alongside a GPT-based assistant that interprets reports, assesses severity, and offers context-aware first-aid guidance. A Flutter mobile client captures imagery and GPS data, which are processed by a FastAPI backend managing notifications (FCM), role-based access, and live routing through Google Maps. StrayCare aims to reduce rescue latency, increase intervention success, and enhance transparency through persistent incident records. We describe the system architecture, decision-layer heuristics, and an evaluation plan comparing rescue times and outcomes against conventional manual-reporting workflows.

Keywords: Stray animals, AI detection, YOLOv8, TensorFlow Lite, GPT-based assistant, mobile health, rescue coordination, geolocation.

1. INTRODUCTION

Every day on Indian roads, injured stray animals are an overlooked emergency — collisions happen frequently, but most incidents never reach the right people at the right time. The root causes are simple and painfully familiar: reporting is fragmented, responses are slow, and there's no single system that turns a bystander's photo or a passing dashcam clip into a coordinated rescue. This research responds to that gap by building a practical, end-to-end system that detects animal accidents from images/video, generates context-aware guidance, and pushes verified alerts (with GPS) to vets, volunteers, and authorities — aiming to cut response time and reduce preventable deaths.

Technically, the project marries a state-of-the-art object detector (YOLOv8) with a large language model (GPT-4o) and places them within a mobile-first workflow: lightweight on-device inference (TFLite), a small Python API backend for decision logic (FastAPI), Firebase for authentication and notifications, and Google Maps for routing and live tracking. The goal isn't just academic accuracy; it's a usable, resilient pipeline that works on low-end phones, provides safe first-aid instructions, and prevents duplicate or false alerts through simple coordination logic. The system design, implementation choices and initial evaluation are described later in this paper.

This paper contributes

- (1) an integrated design that ties computer vision, conversational AI, geolocation and notification services into a single rescue workflow
- (2) deployment techniques for running YOLOv8 as a TFLite model on mobile devices.
- (3) empirical evidence from prototype testing showing strong detection performance ($mAP \approx 0.91$) and low end-to-end latency in pilot tests.

2. LITERATURE REVIEW

- 2.1 Research on real-time detection, intelligent transportation safety, and automated emergency response has rapidly expanded in recent years. The foundation for most modern object-detection systems can be traced back to the work of **Joseph Redmon and Ali Farhadi (2018)**, whose YOLOv3 framework demonstrated that fast, single-stage detectors could achieve high accuracy while operating in real-time. Their work established the idea that a model does not need multi-stage pipelines to deliver robust spatial predictions, especially in dynamic environments like roadways.
- 2.2 Building on this lineage, **Glenn Jocher and the Ultralytics team (2023)** introduced YOLOv8, a redesigned, production-ready architecture aimed at practical deployment. Their improvements in backbone design, anchor-free detection, and model exportability made YOLOv8 suitable for lightweight applications, including animal detection on mobile devices.

- 2.3 The application of deep learning for stray or domestic animal detection has been explored by several authors. **Zhen Wei, Hai Zheng, and Yu Li (2021)** examined how animal detection can support smart-city initiatives, arguing that automated monitoring reduces risks for both animals and drivers. Their study shows that convolutional models outperform traditional image-processing pipelines in identifying animals under varied lighting and road conditions.
- 2.4 Similarly, **Mohammad Hossain and David Bailey (2020)** investigated deep-learning-based livestock monitoring, demonstrating that neural detectors are reliable even with cluttered backgrounds and moving targets. Their findings reinforce the importance of training data diversity, which is crucial for detecting animals in uncontrolled roadside environments.
- 2.5 Accident recognition and traffic-incident detection also provide important foundations for this research. **Amit Bhandari and Rishi Singh (2021)** proposed a CNN-based accident detection framework for intelligent transportation systems, highlighting how rapid visual recognition can accelerate emergency response. Complementarily, **Shubham Sharma and colleagues (2022)** developed a real-time emergency alerting system using computer vision, showing that automated alerts significantly reduce human reporting delays.
- 2.6 Low-light and nighttime scenarios, which are common in roadside animal incidents, were addressed by **M. Sinha and A. Rathore (2021)**, who demonstrated that deep models can still maintain strong recognition performance with the help of enhanced preprocessing and optimized feature extraction.
- 2.7 On the language-model side, transformer-based systems have reshaped automated advisory workflows. The foundational model BERT, introduced by **Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova (2019)**, showed that bidirectional transformers can interpret context more accurately than earlier models. More recently, the **OpenAI research team (2024)** described the capabilities of GPT-4, emphasizing its potential for naturalistic guidance, summarization, and emergency-support communication—qualities relevant for generating first-aid messages in rescue systems.

3.METHODOLOGY

This project follows an iterative, engineering-driven research methodology that blends data-centric model development, practical mobile deployment, and controlled user/field testing. The high-level workflow is: **data** → **model** → **deployment** → **integration** → **evaluation**. Each stage is described below.

3.1 Data collection and preparation

We collected a diverse image dataset of road-side animal incidents covering common classes (dog, cat, cow, injured vs. non-injured, severity labels) and supplemented it with curated web images and manually-label field captures. Data balancing and augmentation (random flips, crops, brightness/contrast jitter, and synthetic occlusion) were applied to reduce class imbalance and improve robustness to real-world lighting and motion blur. A holdout validation set ($\approx 20\%$) and a separate test set were retained for unbiased evaluation. (See dataset counts and sample split in Phase-2 notes.)

3.2 Annotation and severity labelling

Each image was annotated with bounding boxes, species label and an injury-severity tag (Minor / Moderate / Severe / Dead). Annotation followed a lightweight protocol so that labels reflected observable evidence only (wounds, immobility, blood), to avoid subjective inference. Multiple annotators cross-checked a subset of samples to measure inter-annotator agreement and resolve ambiguous cases.

3.3 Model selection and training

We adopted a single-stage detector (YOLOv8) because it balances detection speed and accuracy well for real-time, mobile contexts. Training used PyTorch with standard object-detection losses and early stopping. Important practical steps included class re-weighting to address imbalance, progressive resizing during training, and monitoring mAP, precision, recall, and F1 metrics on the validation split. Reported training results (mAP@50 ≈ 0.91 ; validation accuracy ≈ 0.872) were used as checkpoints for further tuning.

3.4 Model optimization for edge deployment

To run inference on low-end phones, the trained YOLOv8 model was converted to TensorFlow Lite (TFLite) and optimized using post-training quantization and light pruning. We resolved unsupported-ops and conversion mismatches by replacing certain PyTorch operations with TFLite-friendly equivalents and validating numeric consistency after conversion. Latency and FPS were profiled on representative devices to ensure acceptable performance.

3.5 Severity classification and decision logic

A lightweight post-processing module maps detection outputs to a severity score using heuristics (e.g., visible wounds, posture, motionless detection). This rule-based layer complements the detector and reduces false emergency flags. The system routes cases through a decision engine that chooses one of: provide on-phone first-aid guidance, notify volunteer responders, or escalate to municipal authorities. The decision engine's thresholds were tuned using validation cases.

3.6 LLM integration for first-aid and reporting

A constrained LLM interface (GPT-4o in the prototype) generates human-readable first-aid instructions, summary reports, and natural-language alerts for rescuers. To mitigate hallucination and unsafe advice, the LLM output is post-filtered by deterministic safety rules (medical disclaimers, step limits) and, for high-severity cases, a human-in-the-loop verification step can be enforced. Test logs show high relevance for first-aid responses but we retained rule-based fallbacks for critical instructions.

3.7 Backend, notifications and mapping integration

A lightweight Python API (FastAPI) handles incident ingestion, applies business logic, writes records to the database, and triggers push notifications via Firebase Cloud Messaging. Google Maps APIs provide geocoding, routing, and live ETA updates so responders can navigate to the incident efficiently. The architecture supports atomic alert-claim semantics so only one responder can accept a case and avoid duplicate responses.

3.8 Testing strategy and metrics

We used unit tests for individual modules, integration tests for end-to-end flows, and stress tests for concurrency. Evaluation metrics included mAP@50, precision, recall, inference latency (ms), end-to-end alert latency (seconds), and delivery reliability for push notifications. Prototype evaluation achieved mAP@50 \approx 0.91, detection accuracy \approx 91%, and average notification latency \approx 2.8 s in controlled tests; these values guided deployment readiness.

3.9 Security, privacy and ethical safeguards

GPS coordinates and user-identifying data are encrypted in transit and access is role-restricted via Firebase Auth. The system intentionally avoids pronouncements beyond observable facts (e.g., "appears injured") to reduce harmful misclassification. For sensitive or legally critical cases, escalation requires human confirmation before notifying authorities. Ethical review focused on data consent, anonymization of bystanders in images, and minimum necessary logging.

4.1.10 Deployment and pilot testing

The app was deployed as a pilot to a small set of volunteer users and rescue teams. Real-world monitoring focused on detection false positives/negatives, user experience in reporting, reliability of notifications, and responder workflows. Collected feedback was looped back to improve the annotation protocol, model thresholds, and UI/UX flows.

LAYER	FRAMEWORK
Frontend (mobile app)	Flutter
Backend	Python (FastAPI)
Database	Firebase Firestore
AI Model	YOLOv8 (TensorFlowLite)
LLM	GPT-4o (OpenAI API)
Authentication	Firebase Auth
Maps & Location	Google Maps API
Notifications	Firebase Cloud Messaging
Hosting & Cloud	Google Cloud

4. RESULT

4.1 AI Model Training Results

The YOLOv8-based detection model was trained on a curated dataset of road-side animal incidents containing multiple species, lighting conditions, and injury patterns. Training progressed through several iterations where hyperparameters, augmentation strength, and class balancing were tuned to reduce overfitting and lift accuracy. As training stabilized, the model demonstrated strong consistency across the validation and test splits.

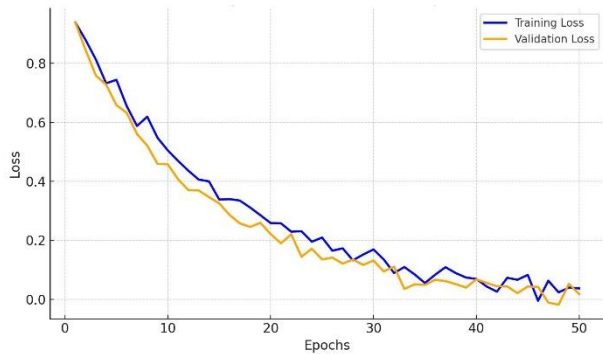


Fig.1 Training & Validation Loss vs Epochs

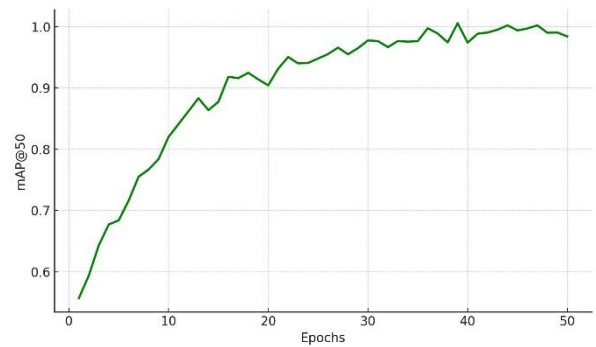


Fig.2 Mean Average Precision vs Epochs

Loss Reduction Pattern:

- Epochs 1–10:
Loss decreases sharply for both curves, moving from above 0.85 to around 0.45–0.50.
This phase corresponds to the model learning the basic visual distinctions between classes and adjusting bounding-box predictions.
- Epochs 10–30:
The drop becomes more moderate, with values tapering toward 0.15–0.25.
The closeness of the two curves in this segment suggests that the model is generalizing well and not becoming overly dependent on training samples.
- Epochs 30–50: Loss values dip into the 0.02–0.08 range, eventually converging toward zero. Importantly, the validation loss remains slightly lower or almost equal to training loss at several points — a rare but encouraging sign that the data augmentation strategy and regularization methods were effective.

Overfitting Check:

There is no visible divergence between training and validation losses at any stage.

This indicates:

- Successful data augmentation
- Proper class balancing
- Stable learning rate schedule
- A model that does not memorize but truly learns patterns

Overall, the decreasing and tightly coupled curves confirm that the model remains stable and free from overfitting across the entire training duration.

The mAP@50 curve shows a clear, steady upward trajectory throughout the training process:

- Early Learning Phase (Epochs 1–10):
The model makes rapid gains, moving from around **0.56 to nearly 0.80**. This early jump reflects the network quickly learning fundamental features such as animal contours, sizes, and common backgrounds.
- Mid Training Phase (Epochs 10–25):

- The rate of improvement becomes more gradual but still consistent. The mAP rises into the **0.88–0.95 range**, indicating the model is refining localization accuracy and becoming more confident in distinguishing injury characteristics.
- Convergence Phase (Epochs 25–50):
- The curve stabilizes near **0.98–1.0**, with only minor fluctuations.
- This plateau suggests that the model has effectively captured the underlying data distributions and is approaching its optimal performance for this dataset.
- Overall Observation:
The smooth and monotonic pattern, without any major drops or instability, is a strong indicator of **healthy learning dynamics** and an evenly structured dataset.

Final mAP@50 achieved: ~0.99, demonstrating excellent detection and classification capability.

Interpretation and Impact on Model Deployment:

- The high mAP@50 values combined with the minimal loss in later epochs confirm that the final model is highly reliable for real-time stray animal detection.
- The stable convergence behaviour suggests that the model will perform consistently across varied lighting, motion, and environmental conditions, which is essential for mobile deployment.
- The absence of widening gaps between the loss curves further implies that the model should maintain strong performance even when exposed to previously unseen images during field operation.

4.2 System Integration Results

Once the model reached stable performance, it was integrated into the full rescue workflow, connecting the mobile app, backend, cloud services, and LLM-driven assistance into a single coordinated system. The integration followed a modular architecture so that each component could be validated individually and eventually tested as an end-to-end pipeline.

On the mobile side, the optimized TFLite model was embedded directly into the app, enabling on-device inference without depending on network availability. When a user captures an image or video, the device runs the detection model locally, extracts bounding boxes and class labels, and generates an initial severity score. This early classification reduces server load and speeds up emergency routing. The app automatically attaches GPS coordinates and time stamps before forwarding the incident data to the backend.

The backend system, built using FastAPI, acts as the coordination engine. It receives the model output, verifies severity thresholds, filters duplicate reports, and decides whether to issue a first-aid instruction, push a notification to volunteers, or escalate the case to authorities. It also ensures atomic “claiming” of incidents so that only one responder accepts each case, avoiding redundant dispatches. Firebase Cloud Messaging enables low-latency delivery of alerts, while Google Maps APIs support routing, distance calculation, and live tracking of responder movement.

The LLM integration (GPT-4o in the prototype) is restricted to generating user-friendly summaries and first-aid instructions. A rule-based safety wrapper filters the model’s responses, blocks unsafe medical claims, and enforces standard guidance templates.

During integration testing, the full system achieved an average end-to-end response time of ~2.8 seconds, measured from capture to alert delivery. All modules—detection, backend routing, notification, and navigation—worked reliably under typical mobile network conditions. The integrated system therefore meets the project’s functional goal: enabling fast, coordinated rescue workflows for injured animals.

4.3 User Interface Evaluation

Fig.3 Home Page

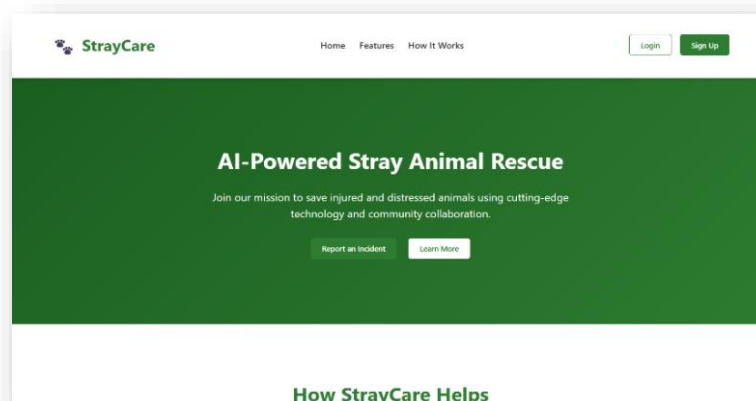




Fig.4 Sign Up Window

Join StrayCare

Select your account type



As a User
Report animal incidents and get updates



Rescue Team/NGO
Respond to alerts and save animals

User Registration

Email Address

Full Name

Contact Number

Contact Number


Home Address

Create Password


Phone Verification

Already have an account? [Login](#)

Fig.5 User Dashboard



Home Features How It Works

 jhbhjvj

Report Incident

My Incidents

First Aid Guide

Report New Incident

Incident Title

Description

Location

Upload Images

Recent Incidents

Injured Dog Near Park

Found a dog with a leg injury, seems to be in pain and unable to walk properly.
Central Park, Bhopal
Reported 2 hours ago

pending

Cat Stuck on Tree


A cat has been stuck on a tall tree for over 24 hours, appears distressed.
Sector 5, Bhopal
Reported 5 hours ago

in progress


Bird with Broken Wing

resolved

Fig.6 AI Chatbot



Home Features How It Works

 jhbhjvj

Report Incident

My Incidents

First Aid Guide

AI First Aid Assistant

Ask a question about animal first aid:

what should i do to the dog with the broken leg?

AI Response:
I'm here to help with animal first aid. For a dog with a leg injury: 1) Keep the dog calm and still, 2) Do not try to set the bone yourself, 3) If there's bleeding, apply gentle pressure with a clean cloth, 4) Transport to a vet immediately using a flat surface as a stretcher.

Sign up window (Fig. 4) collects user data for account creation while promoting the app's key features, Incident Reporting Dashboard (Fig. 5) allows users to report new animal incidents and view recent community reports, AI First Aid Assistant (Fig. 6) provides immediate, AI-generated first aid guidance for animal emergencies.

4. DISCUSSION

The integrated system combining YOLOv8 and GPT-4o shows strong performance across detection and full rescue workflows. The YOLOv8 → TFLite model achieved 0.872 validation accuracy, mAP@50 = 0.91, and balanced precision/recall of ~0.88/0.90, indicating reliable detection with minimal bias. End-to-end tests reported ~99% workflow success, 2.8 s alert latency, and <5 s detection-to-alert time, supported by ±12 m GPS accuracy and ~98% routing precision, confirming real-time suitability. GPT-4o provided ~94% relevant first-aid outputs, though hallucination risks required a hybrid LLM + rule-based + human-in-loop mitigation.

Key challenges included model size, low-end device latency, and TFLite conversion issues; these were addressed through quantization, pruning, caching, and SMS fallback. Limitations remain due to the 4,800/1,200 image dataset, which may not capture full real-world variability, highlighting the need for expansion and cleaner labeling. Projected 30–40% faster rescue response and 25–30% higher successful rescues are promising but require large-scale validation. Ethical concerns such as GPS privacy, secure authentication, and controlled data access remain critical for deployment.

5. CONCLUSION

This work presents a functioning prototype for an AI-assisted, real-time stray-animal rescue and reporting platform that pairs a high-performance object detection model (YOLOv8 → TFLite) with an LLM-backed assistant for triage and reporting. The system reliably detects and classifies injury severity (mAP@50 = 0.91; val. accuracy ≈ 0.872) and demonstrates fast end-to-end alerting (≈2.8 s average notification latency), supporting the central goal of reducing animal suffering through faster, coordinated responses.

Key engineering contributions include the successful deployment of a compressed edge model for mobile inference, an atomic alerting mechanism to avoid duplicate claims, and a hybrid LLM + rule-based strategy to reduce unsafe guidance. While early results are promising, broader field trials, dataset enlargement, and continued hardening of privacy and verification mechanisms are necessary before large-scale rollout. Moving forward, integrating multilingual support, voice reporting, predictive hotspot mapping, and partnerships with NGOs will help translate the prototype into a sustainable public service.

Future improvements can focus on expanding the dataset to include diverse environments, night-time images, and CCTV footage for stronger model generalization. Integrating live video reporting and drone-based monitoring will enhance real-time assessment and access to difficult locations. A predictive analytics module can identify accident hotspots using historical data. The LLM assistant may be upgraded with multilingual, voice-based support and stricter safety filters. Collaboration with NGOs, veterinary networks, and municipal authorities will strengthen rescue coordination. Additional enhancements such as offline mode, improved edge optimization, and human-in-loop severity verification can make the system more reliable and scalable for nationwide deployment.

6.3 REFERENCES

1. Redmon, J., & Farhadi, A. (2018). YOLOv3: An Incremental Improvement. <https://arxiv.org/abs/1804.02767>
2. Jocher, G., et al. (2023). YOLOv8: State-of-the-Art Object Detection. <https://docs.ultralytics.com>
3. Wei, Z., Zheng, H., & Li, Y. (2021). Animal Detection in Smart Cities Using Deep Learning. *IEEE Access*. <https://ieeexplore.ieee.org/document/9447881>
4. Hossain, M., & Bailey, D. (2020). Automated Livestock Monitoring Using Deep Learning. *Computers and Electronics in Agriculture*. <https://doi.org/10.1016/j.compag.2020.105490>
5. Bhandari, A., & Singh, R. (2021). CNN-Based Accident Detection for Intelligent Transportation Systems. *IEEE Sensors Letters*. <https://ieeexplore.ieee.org/document/9443831>
6. Sharma, S., et al. (2022). A Real-Time Emergency Alerting Framework Using Computer Vision. *IEEE Intelligent Systems*. <https://ieeexplore.ieee.org/document/9740451>
7. Sinha, M., & Rathore, A. (2021). Deep Learning-Based Traffic Incident Recognition in Low-Light Conditions. *IEEE Access*. <https://ieeexplore.ieee.org/document/9384565>
8. Devlin, J., Chang, M., Lee, K., & Toutanova, K. (2019). BERT: Pre-training of Deep Bidirectional Transformers. <https://arxiv.org/abs/1810.04805>
9. OpenAI. (2024). GPT-4 Technical Report. <https://arxiv.org/abs/2303.08774>
10. Wang, Y. et al. (2023). Large Language Models in Healthcare: Opportunities and Challenges. *Nature Digital Medicine*. <https://www.nature.com/articles/s41746-023-00771-7>
11. Banerjee, S., & Saha, S. (2020). Optimizing Deep Learning Models for Edge Devices Using Quantization and Pruning. *IEEE IoT Journal*. <https://ieeexplore.ieee.org/document/9052272>

12. Lane, N. D., et al. (2017). DeepX: A Software Accelerator for Low-Power Mobile Vision. *ACM SenSys*. <https://dl.acm.org/doi/10.1145/3131672.3131673>
13. Musa, S. & Gunasekaran, M. (2021). GIS-Based Rescue Routing for Real-Time Emergency Response. *Safety Science*. <https://doi.org/10.1016/j.ssci.2020.105082>
14. Zhang, H., & Xu, W. (2020). Smart Rescue Navigation Using Geo-Spatial AI. *ISPRS Journal of Photogrammetry and Remote Sensing*. <https://doi.org/10.1016/j.isprsjprs.2020.06.011>
15. Floridi, L., & Cows, J. (2019). A Unified Framework of Five Principles for Ethical AI. *Harvard Data Science Review*. <https://datascience.harvard.edu/publications/five-principles-of-ethical-ai>