

SafeTriage: A Hybrid AI Clinical Escalation Assistant

1. Problem Statement

Large Language Models (LLMs) show strong reasoning ability in medical triage tasks. However, pure LLM-based systems present a critical risk: under-triage. Missing a high-risk case such as chest pain, stroke symptoms, or unconsciousness can have severe consequences.

In community health settings, especially in low-resource environments, frontline workers require fast, reliable, and safety-focused triage support. The system must prioritize patient safety over conversational fluency.

SafeTriage addresses this by combining LLM reasoning with deterministic safety guardrails to eliminate missed high-risk cases.

2. Motivation

Most AI triage systems rely solely on probabilistic outputs. While flexible, they can fail in edge cases or adversarial inputs.

In real-world healthcare:

- False negatives (missing emergencies) are unacceptable.
- Over-escalation is safer than under-escalation.
- Systems must operate offline or with limited connectivity.

SafeTriage is built with a safety-first philosophy: zero tolerance for missed high-risk cases.

3. Solution Overview

SafeTriage is a hybrid architecture composed of two layers:

1. LLM Clinical Reasoning Layer
2. Deterministic Safety Override Layer

The LLM generates structured JSON outputs including:

- Summary
- Risk level (Low / Moderate / High)
- Red flags

- Recommended action
- Missing information

The safety layer then validates and escalates decisions based on predefined clinical rules.

4. System Architecture

Patient Input

↓

LLM Structured Output

↓

Deterministic Safety Rules

↓

Final Escalation Decision

The deterministic layer overrides model predictions in cases such as:

- Elderly patients with chest discomfort
- Severe hypertension (≥ 180 systolic)
- Severe hypotension (< 90 systolic)
- Tachycardia (> 130 bpm)
- Critical neurological symptoms
- Unconsciousness

This ensures safety constraints are always enforced.

5. Safety Layer Design

The safety override layer is rule-based and interpretable. It performs:

- Vital sign extraction
- Symptom keyword detection
- Age-based escalation logic
- Emergency symptom enforcement

If a high-risk condition is detected, the system escalates to "High" risk regardless of LLM output.

This hybrid approach prevents model hallucination or under-confidence from compromising patient safety.

6. Evaluation Methodology

SafeTriage was evaluated using:

- 25 structured clinical test cases
- Adversarial stress test cases
- Quantitative performance metrics

Evaluation metrics included:

- Accuracy
 - Confusion matrix
 - False negative count
 - False positive count
 - High-risk recall
-

7. Results

Structured Evaluation (25 cases):

- Accuracy: 100%
- False Negatives (High-risk missed): 0
- False Positives (Over-escalation): 1

Stress Testing:

- Maintained 100% high-risk recall
- Demonstrated conservative escalation bias

SafeTriage achieved zero missed high-risk cases across evaluation.

This validates the effectiveness of the hybrid safety architecture.

8. Limitations

- Rule-based logic may over-escalate in borderline cases.

- Clinical validation with real patient data is required.
 - Performance depends on input quality and completeness.
-

9. Future Work

- Integrate real-world clinical datasets
 - Add explainability dashboards
 - Optimize inference latency
 - Deploy lightweight offline models for rural environments
-

10. Ethical Considerations

SafeTriage is not a medical diagnosis tool. It is designed as a decision-support assistant for trained health workers.

The system prioritizes safety and transparency. All escalation rules are deterministic and interpretable.

Patient data privacy and secure deployment must be ensured in production environments.

11. Conclusion

SafeTriage demonstrates that combining LLM reasoning with deterministic safety guardrails can significantly reduce the risk of under-triage.

By eliminating false negatives in structured testing and maintaining conservative escalation in stress scenarios, SafeTriage offers a robust and safety-focused AI triage framework suitable for community healthcare support.