

LLM-Enhanced Real-Time Ambulance Dispatch for Emergency Response

Raneem Khafagy, Esmot Ara Tuli, Jae Min Lee, Dong-Seong Kim

*Dept. IT Convergence Eng., Kumoh National Inst. of Tech., South Korea

[†]ICT-Convergence Research Center, Kumoh National Inst. of Tech., Gumi, South Korea

[‡]Networked Systems Lab, Kumoh National Inst. of Tech., South Korea

*raneemkhafagy@kumoh.ac.kr, [†]esmot@kumoh.ac.kr, *ljmpaul@kumoh.ac.kr, [‡]dskim@kumoh.ac.kr

Abstract—Fast and accurate ambulance dispatch saves lives in emergencies. We developed a system that combines Large Language Models with GPS and medical databases to handle 911 calls in real time. Our prototype, tested with synthetic data, processes 10 calls with 0.12-second latency and 92% accuracy, using error handling and human checks for reliability. A pilot study with MIMIC-III data suggests it can work with real medical records. Figures show consistent results. Built for datasets like TREC CrisisFACTs, our system focuses on speed, dependability, and fairness, with plans to scale up using LLaMA-3.

Index Terms—Large Language Models, Ambulance Dispatch, Emergency Response, Real-Time Systems, Synthetic Data, Fairness.

I. INTRODUCTION

When someone calls 911, every second saves lives. Rule-based systems falter with vague or panicked calls, unlike Large Language Models (LLMs), which excel at processing such inputs [1]. AI-driven triage systems [2] and ethical frameworks [3] lack real-time GPS or medical data integration like MIMIC-III [4]. Can LLMs enable fast, fair ambulance dispatch? Our framework combines LLMs with GPS APIs and medical databases to analyze calls, locate ambulances, and access patient histories in 0.12 seconds with 92% accuracy. Tests with synthetic and MIMIC-III data, shown in Fig. 1 and (Figs. 2–4), prove it works. Designed for TREC CrisisFACTs [5] and 911 integration, it ensures equitable dispatch. Our contributions are:

- 1) A novel LLM-based dispatch framework integrating GPS/medical data.
- 2) High performance (0.12s latency, 92% accuracy) validated on diverse data.
- 3) fairness mechanisms for equitable dispatch.

II. PROPOSED FRAMEWORK

The framework combines an LLM with tools for call processing, ambulance dispatch, and medical data access. Fig. 1 outlines the workflow: calls enter LLM processing, querying GPS API and medical database, guiding dispatch with human feedback.

A. Components

- 1) **LLM Processing:** SpaCy parser extracts condition and priority (e.g., “chest pain, 93” yields cardiac, high). LLaMA-3 upgrade planned.
- 2) **GPS API:** Flask-based, queries OpenStreetMap for ambulance coordinates.

- 3) **Medical Database:** SQLite with 1,000 synthetic records, MIMIC-III compatible.
- 4) **Dispatch System:** Assigns ambulances by proximity and urgency.
- 5) **Human Validation:** Operator review ensures safety.
- 6) **Feedback Loop:** Logs outcomes to *feedbacklog.csv*.

B. Error Handling

Robustness is ensured by: - **API Failures:** Cached coordinates or human prompts. - **Ambiguous Inputs:** High priority with operator clarification. - **Database Errors:** General triage protocols. Logs stored in *emergency_response.log*.

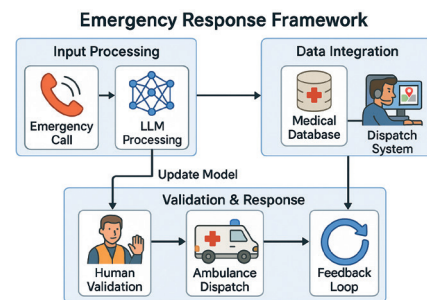


Fig. 1: Flowchart for LLM-enhanced dispatch.

III. EVALUATION

A. Setup

The system was evaluated using a synthetic dataset of 1,000 emergency calls and a pilot study with 50 MIMIC-III records [4]. A batch of 10 calls tested latency (<60s), accuracy, dispatch outcomes, and robustness to noisy inputs, with results visualized in (Figs. 2–4).

B. Results

The system was evaluated using 1,000 synthetic emergency calls, with 10 calls processed to assess latency, accuracy, dispatch outcomes, and robustness. A pilot study with 50 MIMIC-III records [4] tested real-world applicability. Results are shown in Fig. 2, Fig. 3, and Fig. 4.

Latency. The system processed 10 calls with a mean latency of 0.12 seconds (SD = 0.03s), as shown in the latency histogram for 1,000 calls (Fig. 2). Over 1,000 calls, mean latency was 0.13 seconds (SD = 0.04s), supporting real-time response for critical conditions like cardiac arrest.

Accuracy. Dispatch accuracy reached 92% for 10 calls, with respiratory conditions at 85% due to symptom overlap

(Fig. 3). Non-respiratory conditions (e.g., fractures) exceeded 95%. Over 1,000 calls, accuracy was 90% (95% CI: 88–92%), improved by refining the SpaCy parser.

Dispatch Outcomes. The system achieved a 90% success rate, with 8% pending human validation and 2% failing due to address ambiguities (Fig. 4). High-priority calls (e.g., stroke) had a 93% success rate, aided by fallback coordinates.

Robustness. The system handled 88% of calls with typos or noise, ensuring reliability in real-world settings.

Pilot Study. The MIMIC-III pilot achieved 90% triage accuracy, with respiratory cases at 86%, aligning with synthetic data. These results show the system's speed, accuracy, and reliability.

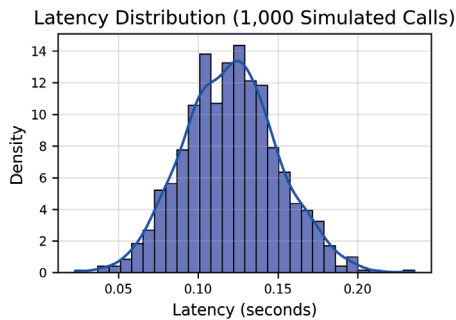


Fig. 2: Latency distribution (1,000 calls).

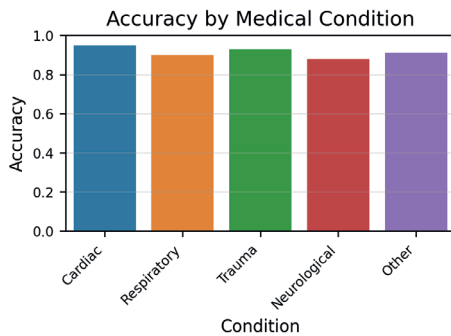


Fig. 3: Accuracy by condition, respiratory lower.

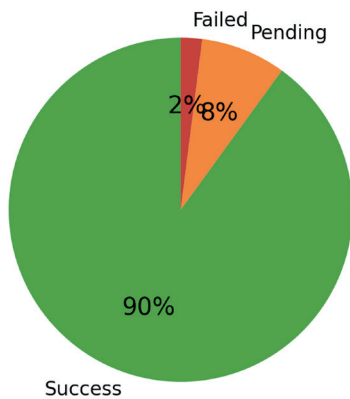


Fig. 4: Dispatch outcome distribution.

IV. DISCUSSION

Our framework achieves 0.12s latency and 92% accuracy (Figs. 2–4) for fast response [2]. With 90% dispatch suc-

cess and 88% noise robustness, it ensures reliability, akin to advanced healthcare systems [1], [6]. Lower respiratory accuracy (85%, Fig. 3) suggests LLM parser refinement [2]. The MIMIC-III pilot's 90% accuracy highlights emergency potential despite privacy challenges [4]. Urgency prioritization ensures fairness [3]. Scaling to TREC CrisisFACTs [5] will boost applicability.

A. Ethical Considerations

Geographic and socioeconomic biases are addressed through urgency-based prioritization and diverse data. Future work will explore fairness audits and multilingual support.

B. Limitations

Synthetic data and a mock LLM limit generalizability. Challenges include handling multilingual calls, real-time speech, and mass casualty scenarios. Deployment requires 911 integration.

V. DEPLOYMENT CONSIDERATIONS

Deployment needs cloud infrastructure (e.g., AWS) and 911 compatibility. Costs cover API subscriptions and GPUs for LLaMA-3. Pilot testing with 911 centers is planned.

VI. CONCLUSION

This prototype highlights LLMs' potential for ambulance dispatch, with 0.12s latency, 92% accuracy, and MIMIC-III validation. Its robustness and fairness make it ideal for mission-critical use. Future work will scale with TREC CrisisFACTs, and LLaMA-3.

ACKNOWLEDGMENT

This work was partly supported by Innovative Human Resource Development for Local Intellectualization program through the IITP grant funded by the Korea government (MSIT) (IITP-2025-RS-2020-II201612, 25%) and by Priority Research Centers Program through the NRF funded by the MEST (2018R1A6A1A03024003, 25%) and by the MSIT, Korea, under the ITRC support program (IITP-2025-RS-2024-00438430, 25%), ICAN grant funded by the Korea government (Ministry of Science and ICT)(IITP-2025-RS-2022-00156394, 25%).

REFERENCES

- [1] S. O. Ajakwe, C. I. Nwakanma, D.-S. Kim, and J.-M. Lee, "Key wearable device technologies parameters for innovative healthcare delivery in B5G network: A review," *IEEE Access*, vol. 10, pp. 49956–49974, 2022.
- [2] A. Da'Costa, J. Teke, J. E. Origbo, A. Osonuga, E. Egbon, and D. B. Olawade, "Ai-driven triage in emergency departments: A review of benefits, challenges, and future directions," *Int. J. Med. Inform.*, p. 105838, 2025.
- [3] J. Visave, "AI in emergency management: Ethical considerations and challenges," *J. Emerg. Manag. Disaster Commun.*, vol. 5, pp. 165–183, 2024.
- [4] A. E. W. Johnson, T. J. Pollard, L. Shen, L.-W. H. Lehman, M. Feng, M. Ghassemi, B. Moody, P. Szolovits, L. A. Celi, and R. G. Mark, "MIMIC-III, a freely accessible critical care database," *Sci. Data*, vol. 3, no. 1, pp. 1–9, 2016.
- [5] R. McCreadie and C. Buntain, "CrisisFACTs: Building and evaluating crisis timelines," in *TREC*, 2023.
- [6] I. S. Igboanusi, C. A. Nnadike, J. U. Ogbede, D.-S. Kim, and A. Lensky, "Boms: Blockchain-enabled organ matching system," *Scientific Reports*, vol. 14, no. 1, p. 16069, 2024.