

GROUP 7

Team Members:

- Amal Ajay A - KTE22EC013
- Joshwin Binu - KTE22EC038
- Judin Joe Mathew - KTE22EC039
- Justin Varghese - KTE22EC040

Project Guide:

 Asst. Prof. Shinoj K Sukumaran

Introduction: The Challenge

- **Kerala's Vulnerability:** Prone to floods, landslides, and other natural disasters.
- **Critical Infrastructure Failure:** Disasters often knock out power, internet, and mobile networks.
- **The Consequence:** Isolated communities are cut off from vital, life-saving information precisely when they need it most.

Our Solution: The Offline Crisis Assistant

Why *Offline*?

Cloud-Based Assistants (e.g., Siri, Alexa)






- ✗ **Require Internet**
- ✗ **Fail when networks are down**
- ✗ **Not tailored for disaster medicine**

Our Crisis Assistant

- ✓ **Fully Offline Operation**
- ✓ **Rugged and Portable**
- ✓ **Specialized for First Aid & Survival**

A rugged, portable device powered by a specialized offline LLM to provide immediate, expert-level emergency guidance without any connectivity.




Objectives

-  To design and develop a rugged, portable emergency assistant device.
-  To optimize and run a fine-tuned offline LLM on embedded hardware (Raspberry Pi).
-  To support intuitive voice and text input for context-aware guidance.
-  To provide reliable information on first aid, survival, and emergency protocols.
-  To function entirely without internet dependency.

Abstract

In disaster-stricken areas where internet connectivity is often unavailable, access to timely and reliable emergency information is critical. This project presents the design and development of a rugged, portable device powered by a fine-tuned Large Language Model (LLM) that operates entirely offline. The system provides context-aware assistance on first aid, CPR, wound care, and survival techniques via text or speech input. Trained on verified content from the WHO and Red Cross, the LLM ensures reliable responses. Built on a Raspberry Pi with a microphone, speaker, and battery pack, this project demonstrates how embedded AI can strengthen disaster preparedness and response without relying on cloud services.

Literature Review

	Paper & Key Points	Relevance to Our Project
	Zheng et al. (2024) <i>A Review on Edge Large Language Models</i> Techniques for running LLMs on low-power devices.	Provides the foundational methodology for model optimization (quantization, pruning) crucial for our Raspberry Pi deployment.
	Basit et al. (2024) <i>MedAide: On-Premise Medical Assistance</i> Offline medical chatbot using LoRA fine-tuning.	Directly validates our core concept of an offline, reliable medical assistant, informing our fine-tuning strategy.
	Goecks & Waytowich (2023) <i>DisasterResponseGPT</i> Generates disaster response plans from scenarios.	Informs the design of our system's logic for structured emergency guidance and planning in chaotic environments.

Phase 1 Work Plan (July 8 – September 30)

Stage	Dates	Key Tasks
1. Ideation	Jul 8 – 21	Brainstorming, topic finalization, preliminary literature survey on edge AI and disaster tech.
2. Scoping	Jul 22 – Aug 4	Finalize objectives, draft hypothesis, initial block-level design.
3. Planning	Aug 5 – 18	Complete system architecture, assign team roles (HW/SW), list requirements.
4. Feasibility	Aug 19 – Sep 1	Preliminary tests: model size vs. hardware capacity, component testing.
5. Refinement	Sep 2 – 22	Iterate on design based on test results, compile initial analysis.
6. Reporting	Sep 23 – 30	Prepare and review Phase 1 final report and presentation.

Current Progress



Output Obtained

```
jane ~ 01:51 ) cd
jane ~ 01:51 ) cd Documents/SRE-RAG
jane ~/SRE-RAG 1 main ? 01:51 ) cd knowledgebase
jane ~/SRE-RAG/knowledgebase 1 main ? v3.13.3 01:52 ) python3 rag_assistant\_.py
▲ Rich not installed. Using basic formatting.
▲ OFFLINE CRISIS ASSISTANT - DEMO VERSION
Simulating LLM-powered emergency guidance system
No internet connection required!

Choose mode:
1. Demo Mode (for presentation)
2. Interactive Mode (for testing)
3. Show System Status

Select mode (1, 2, or 3): 1
```

```
■ STEP-BY-STEP PROCEDURE:
1. SCENE SAFETY - Ensure area is safe for you and victim
2. UNIVERSAL PRECAUTIONS - Wear gloves or use barrier if available
3. EXPOSE THE WOUND - Remove/cut clothing to see injury clearly
4. DIRECT PRESSURE - Apply firm pressure with clean cloth/gauze
5. MAINTAIN PRESSURE - Do not lift to check if bleeding stopped
6. ELEVATE IF POSSIBLE - Raise injured area above heart level
7. PRESSURE POINTS - If bleeding continues, apply pressure to arterial points
8. SECURE BANDAGE - Apply pressure bandage, check circulation below wound
9. TREAT FOR SHOCK - Keep victim warm, lying down, legs elevated
10. MONITOR VITALS - Check breathing, pulse, consciousness level

▲ CRITICAL WARNINGS:
▲ NEVER remove embedded objects (knives, glass, etc.)
▲ Do NOT use tourniquets unless trained (life/limb situations only)
▲ If blood soaks through bandage, add more layers - don't remove
▲ Watch for signs of shock: pale, cold, weak pulse, confusion

📖 SEEK IMMEDIATE MEDICAL HELP IF:
• Bleeding won't stop after 10 minutes of direct pressure
```

```
📖 OFFLINE CRISIS ASSISTANT - SYSTEM STATUS
=====
System Mode: OFFLINE OPERATIONAL
Knowledge Base: LOADED
Emergency Procedures: 6 PROCEDURES READY
Response Time: <3 SECONDS
Network Dependency: NONE (FULLY OFFLINE)
=====

■ AVAILABLE EMERGENCY PROCEDURES:
-----
• Severe Bleeding Control - HIGH PRIORITY
• Cardiopulmonary Resuscitation (CPR) - CRITICAL PRIORITY
• Thermal Burns Treatment - HIGH PRIORITY
• Airway Obstruction (Choking) - CRITICAL PRIORITY
• Bone Fractures and Sprains - MODERATE PRIORITY
• Medical Shock Treatment - CRITICAL PRIORITY
=====

● DEMO: Processing Sample Emergency Queries
```

```
=====
📖 EMERGENCY RESPONSE: SEVERE BLEEDING CONTROL
=====
▲ URGENCY LEVEL: HIGH
★ QUERY: Someone is bleeding heavily from their arm
● TIME: 2025-09-30 01:56:26
● CONFIDENCE: 6.7%
● SOURCE: WHO Emergency Care Guidelines 2023

▲ REQUIRED SUPPLIES:
• Clean cloth or sterile gauze pads
• Pressure bandages or elastic wrap
• Medical gloves (if available)
• Scissors to cut clothing
• Blanket for shock treatment

■ STEP-BY-STEP PROCEDURE:
1. SCENE SAFETY - Ensure area is safe for you and victim
2. UNIVERSAL PRECAUTIONS - Wear gloves or use barrier if available
3. EXPOSE THE WOUND - Remove/cut clothing to see injury clearly
4. DIRECT PRESSURE - Apply firm pressure with clean cloth/gauze
5. MAINTAIN PRESSURE - Do not lift to check if bleeding stopped
6. ELEVATE IF POSSIBLE - Raise injured area above heart level
7. PRESSURE POINTS - If bleeding continues, apply pressure to arterial points
8. SECURE BANDAGE - Apply pressure bandage, check circulation below wound
9. TREAT FOR SHOCK - Keep victim warm, lying down, legs elevated
10. MONITOR VITALS - Check breathing, pulse, consciousness level
```


[illegible]

```

Ulama.context: freq_base = 1000.0
Ulama.context: freq_scale = 1
Ulama.context: CPU output buffer size = 0.22 MiB
Ulama.kv.cache: CPU KV buffer size = 1550.00 MiB
Ulama.kv.cache: size = 1530.00 MiB ( 400k slots, 32 rows, 1/1 seqs, K (f16): 768.00 MiB, V (f16): 768.00 MiB
Ulama.context: Flash attention was auto, set to enabled
Ulama.context: CPU context buffer size = 72.81 MiB
Ulama.context: graph nodes = 935
Ulama.context: graph compile = 1
common_init_from_params: added cudaDeviceProp limit bias = -1sf
common_init_from_params: added cudaDeviceProp limit bias = -1sf
common_init_from_params: setting dev_parallel_list.s to dev_size = 0sf
common_init_from_params: warning: do the model with no empty row - please wait ... (no-memory to display)
note: Ulama threshold: init, n_threads = 4
note: chat template is disabled, enabling conversation mode (disable it with no_conv)
*/* User-specified prompt will pre-start conversation, did you mean to set --system-prompt (-sys) instead?
note: chat template example:
[system]
You are a helpful assistant QwenD.
[user]
Hello QwenD.
[assistant]
There QwenD.
[user]
Hi QwenD.
[assistant]
system_info: n_threads = 4 (x.thread_count = 4) / 16 | CPU | SSE3 + 1 | SSE3 + 1 | AVX = 1 | AVX_VNNI = 1 | AVX2 = 1 | F16C = 1 | FMA = 1 | BMI2 = 1 | L1LARGE = 1 | OPMW = 1 | HPAK = 1
note: Interactive mode on.
sample_name: QWEN2.72B
tokenizer param:
repeat_last_n = 64, repeat_penalty = 1.000, frequency_penalty = 0.000, presence_penalty = 0.000
top_attenitor = 0.000, top_n = 1.750, top_attenitor_length = 2, top_attenitor_size = 400
top_k = 40, top_p = 0.950, min_p = 0.003, stop_probability = 0.003, stop_token_id = 0.100, typical_p = 1.000, top_n_logits = -1.000, temp = 0.800
top_n_logits = 0, min_tokens = 0.100, streamer = 0, streamer_size = 0.000
tokenizer chain: logits -> logits to penalties -> top-k -> top-p -> top-s -> min-p -> stop -> temp-temp -> dist
generator: k_sos = 400, n_batch = 2048, n_predict = 52, n_keep = 1
== Warning in interactive mode. ==
- Press Ctrl+C to interrupt at any time.
- Press Enter to return control to the AI.
- To return control without starting a new line, and your input with ' '.
- If you want to submit another line, and your input with '\n'.
- Not using system message. To change it, set a different value via --sys PROMPT
QwenD> What is the capital of France?QwenD>Paris
The capital of France is Paris. It is not only the country's largest city but also a global center for art, fashion, gastronomy, and

```

References



World Health Organization, *First Aid Guidelines*, 2023. <https://www.who.int>



International Federation of Red Cross, *First Aid Manual*, 2022. <https://www.ifrc.org>



Vicuna Team, "Vicuna: An Open LLM Dialogue Assistant," 2023. <https://lmsys.org>



ACM TIST / ArXiv, *LLMs for Emergency Response*, 2023. <https://arxiv.org/abs/2306.08956>



Basit et al., *MedAide*, 2024. <https://arxiv.org/abs/2403.00830>



Goecks & Waytowich, *DisasterResponseGPT*, 2023. <https://arxiv.org/abs/2306.17271>

Thank You

Questions?

✉ group-7@example.com github.com/group-7