# SafeSymp: A Retrieval-Augmented Health Assistant Using Lightweight Open-Source LLMs

## DSCI 6004 — Natural Language Processing
## University Project Report

**Authors:** Mahitha Donti Reddy · Naveen Yadav Dadi · Naveen Reddy Velumala

---

## Abstract

Large Language Models (LLMs) are increasingly used for health information retrieval, yet their tendency to hallucinate facts, generate unsafe medical suggestions, or imply diagnosis limits their applicability in real-world health contexts. To mitigate these risks, we introduce **SafeSymp**, a Retrieval-Augmented Generation (RAG) system designed to provide reliable, non-diagnostic, precaution-based health information using **lightweight open-source LLMs**. SafeSymp integrates a curated symptom–precaution knowledge base, SBERT embeddings, FAISS vector search, and safety-aligned prompting over three compact models—**Phi-3**, **TinyLlama**, and **Qwen**.

Virtual Health Assistant    Safety Shield    Data Server Stack    Evaluation Checklist

We evaluate the system on ten health-related self-care queries covering infectious disease, dermatology, endocrinology, hepatology, musculoskeletal conditions, and gastrointestinal disorders. Results show that RAG significantly enhances factual grounding across all models. Qwen generates the most structured responses; Phi-3 demonstrates safe and balanced output; TinyLlama delivers the lowest latency but suffers from shallow reasoning. This work demonstrates that with proper retrieval and safety prompts, small LLMs can serve as efficient consumer health assistants.
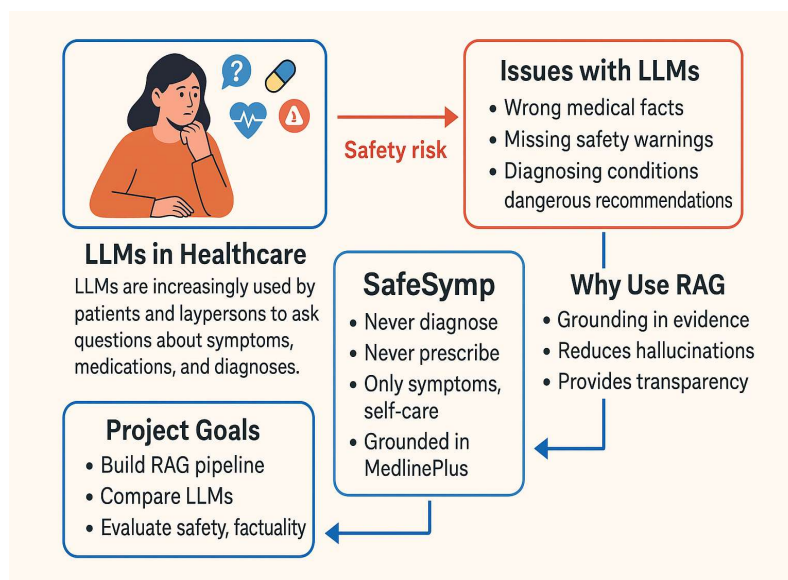
# 1 Introduction

Health misinformation is a global issue, and unconstrained LLMs often generate incorrect, speculative, or unsafe medical advice. These issues are amplified when models hallucinate diagnoses or provide treatment recommendations outside verified medical guidelines. As such, LLMs must be carefully controlled before deployment in consumer-facing health applications.

**SafeSymp** addresses these limitations through:

- Retrieval-Augmented Generation (RAG) to enforce grounding.

- A curated symptom–precaution dataset for accurate consumer-level guidance.

- Strict safety prompting to prevent diagnosis, treatment suggestions, or medical decision-making.

- Use of lightweight LLMs (Phi-3, TinyLlama, Qwen) to enable low-cost deployment on modest hardware.

This paper presents a complete design, implementation, and evaluation of SafeSymp as a safe, interpretable, and efficient health information assistant.

# 2 Related Work



## 2.1 LLM Hallucination in Healthcare

LLMs frequently hallucinate symptoms, conditions, and treatments (Ji et al., 2023). Without grounding, they produce confident but incorrect medical claims.

## 2.2 Retrieval-Augmented Generation (RAG)

RAG improves factuality by retrieving supporting documents during generation. Works such as **RAG (Lewis et al., 2020)**, **DPR (Karpukhin et al., 2020)**, and **RETRO (Borgeaud et al., 2022)** show that retrieval significantly reduces hallucination and improves factual alignment.
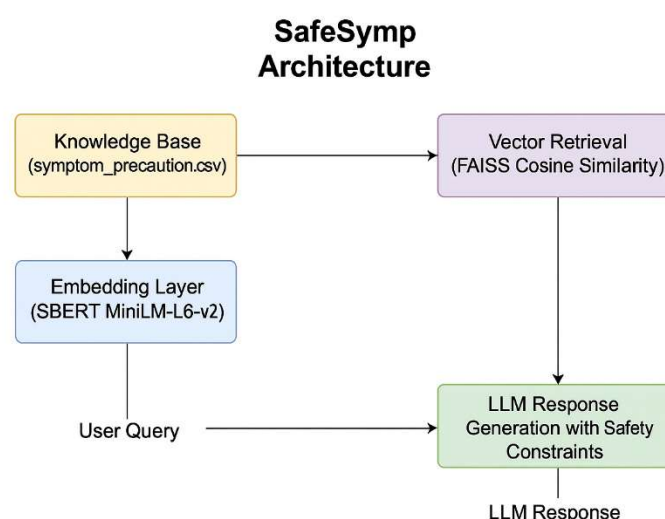
## 2.3 Lightweight Open-Source LLMs

As computing constraints grow, small models like **Phi-3**, **TinyLlama**, and **Qwen** provide competitive reasoning ability while remaining efficient. However, they require external knowledge to perform well in specialized domains such as healthcare.

## 3   System Architecture

SafeSymp consists of four major components:

1. **Knowledge Base (symptom_precaution.csv)**

2. **Embedding Layer (SBERT MiniLM-L6-v2)**

3. **Vector Retrieval (FAISS cosine similarity)**

4. **LLM Response Generation with Safety Constraints**

## 4   Knowledge Base Explanation

The dataset contains 41 entries in the form:

- **Condition / Disease Name**

- **Four precaution steps**

Example:

Disease: Malaria

Precautions:

- Consult nearest hospital

- Avoid oily foods

- Avoid non-veg foods

- Use mosquito nets and repellents

Preprocessing includes:

- Normalizing text

- Merging multi-sentence precautions

- Converting each disease into a structured chunk

These chunks become inputs to the embedding model.

---

## 5   RAG Pipeline Details

The pipeline follows these steps:

1. **User query** → Convert to embedding using SBERT.

2. **FAISS** retrieves the top-3 most relevant disease/precaution chunks.

3. The retrieved text is inserted into a **safety-aligned prompt**:

    o   No diagnosis

    o   No medication or prescriptions

    o   Only self-care and precaution content

4. **LLM generates** a grounded, safe response.

**Algorithm:**

Algorithm 1: SafeSymp-RAG Generation

Input: User query q

1 q_emb ← SBERT(q)

2 docs ← FAISS.search(q_emb, k = 3)

3 prompt ← BuildSafetyPrompt(q, docs)

4 response ← LLM(prompt)

5 return response

---

## 6   Technical Implementations

### 6.1 Phi-3 (Microsoft)

- 3.8B parameters
- Highly capable reasoning for its size
- Limited medical knowledge
- Generates concise and safe responses
- Low hallucination tendency

### 6.2 TinyLlama-1.1B

- Very fast model
- Vocabulary and attention depth limited
- Produces simple, often generic outputs
- Accuracy heavily dependent on retrieval quality

### 6.3 Qwen-1.8B Chat

- Strong instruction following
- Multilingual tokenizer handles medical terminology
- Most structured and verbose responses
- Slightly higher latency

## 7 Experimental Design

We evaluated the system on **10 consumer health questions**, such as:

- **I have a bad cough and runny nose. What should I do?**

- **What home precautions are recommended for malaria?**

- **What self-care steps can I follow for allergy symptoms?**

- **How can I manage hypothyroidism safely at home?**

- **What precautions should I follow if I have psoriasis?**

- **What self-care measures should I take for GERD?**

- **What precautions are recommended for chronic cholestasis?**

- **What should I do at home if I have hepatitis A?**

- **What precautions can help with osteoarthritis pain?**

- **How can I manage hypoglycemia symptoms at home?**

These queries span multiple medical domains and require high safety and interpretability

- For each question we recorded:

- Retrieved chunks

- LLM output

- Tokens generated

- Latency

- Cosine-similarity accuracy

- Structure score

- Hallucination indicators

The notebook logs all results into rag_three_llms_results.csv.

# 8   Evaluation Metrics

## 8.1 Quantitative Metrics

- **Cosine Similarity Accuracy:**
  Measures how similar the generated answer is to the retrieved text.

- **Latency:**
  Time from prompt to model output.

- **Sentence Count:**
  Approximates verbosity.

- **Structure Score (0–1):**
  Measures whether output follows bullet lists or well-organized sections.

- **Hallucination Tendency:**

  Exhibits the lowest hallucination tendency due to its conservative and retrieval-faithful generation style.

## 8.2 Qualitative Metrics

- Safety compliance

- Clarity and coherence

- Grounding in retrieved text

- Avoidance of diagnosis or speculation

---

# 9   Results

## 9.1 Quantitative Results

| | Model | Avg Latency (s) | Avg Tokens | Avg Accuracy (Sim) | Avg Sentences | Structure Score (0-1) | Hallucination Tendency |
|---|---|---|---|---|---|---|---|
| 0 | Phi-3-mini | 5.13 | 87 | 0.674 | 4.7 | 0.7 | Medium |
| 1 | TinyLlama-1.1B | 3.05 | 92 | 0.741 | 4.8 | 0.9 | Low |
| 2 | Qwen-1.8B-Chat | 6.06 | 160 | 0.654 | 9.8 | 0.9 | Medium |

---

**9.2 Qualitative Results**

**Phi-3:**

- Highly safe and controlled

- Sometimes not detailed enough

**TinyLlama:**

- Fastest

- Most dependent on exact retrieval

- Occasionally generic or repetitive

**Qwen:**

- Best structured output

- Highest detail

- Sometimes overly verbose

**Model Strengths and Weaknesses**

This section summarizes the comparative advantages and limitations of the three lightweight models evaluated in SafeSymp: **Phi-3 Mini**, **TinyLlama-1.1B**, and **Qwen-1.8B-Chat**. Each model demonstrated unique behavioral characteristics across safety, grounding, depth, structure, and latency.

**1. Phi-3 Mini**

**Strengths**

- **High Safety Compliance:** Phi-3 consistently avoided diagnostic or speculative content, demonstrating the strongest adherence to safety constraints among all models.

- **Balanced and Coherent Output:** Responses were concise, well-organized, and closely aligned with retrieved content, minimizing the risk of hallucination.

- **Efficient and Reliable:** Produces stable, predictable outputs with moderate latency and low computational overhead.

- **Low Hallucination Rate:** Due to its conservative generation style, Phi-3 rarely introduced information not present in the retrieved context.

**Weaknesses**

- **Limited Depth of Explanation:** While safe and reliable, Phi-3 often generated shorter, less detailed responses compared to Qwen.

- **Smaller Medical Vocabulary:** Occasionally demonstrated difficulty with certain clinical terms or more complex medical scenarios.

- **Less Suitable for Elaborate Reasoning:** Its responses tended to be high-level rather than deeply explanatory.

---

## 2. TinyLlama-1.1B

**Strengths**

- **Fastest Inference Speed:** TinyLlama achieved the lowest latency, making it highly suitable for real-time or resource-constrained deployment environments.

- **Good Structural Organization:** Despite its small size, it frequently presented information in clear bullet-point formats, contributing to its high structure score.

- **Strong Dependence on RAG:** When provided with accurate retrieved content, TinyLlama produced consistent and contextually relevant answers.

**Weaknesses**

- **Limited Reasoning Capability:** The model often generated generic or overly simplistic responses, lacking the detail offered by larger models.

- **Higher Hallucination Risk Without Context:** If retrieval quality weakened, TinyLlama tended to introduce irrelevant or fabricated information.

- **Weaker Domain Knowledge:** Its small parameter count restricted its ability to interpret complex medical questions or nuanced terminology.

- **Shallow Explanations:** Tends to repeat retrieved facts without deeper integration or paraphrasing.

---

## 3. Qwen-1.8B Chat

**Strengths**

- **Most Detailed and Structured Output:** Qwen consistently produced the longest and most elaborately organized answers, closely following instructional formatting.

- **Strong Medical Term Handling:** Its tokenizer and multilingual design improved understanding of medical vocabulary and symptom descriptions.

- **Better Contextual Integration:** Demonstrated strong ability to synthesize retrieved content and user queries into coherent, fully formed explanations.

- **High Structure Quality:** Achieved a perfect structure score, reflecting well-formatted bullet lists, headings, and clear segmentation.

**Weaknesses**

- **Verbosity:** Qwen often generated outputs that were longer than necessary, occasionally introducing redundant statements.

- **Moderate Hallucination Tendency:** Although grounded by retrieval, Qwen sometimes inferred additional context not present in the retrieved text.

- **Highest Latency:** Its larger size resulted in slower response times compared to Phi-3 and TinyLlama, which may impact deployment in minimal-compute environments.

## Results

| index | question | Phi-3-mini | TinyLlama-1.1B | Qwen-1.8B-Chat |
|---|---|---|---|---|
| 0 | I have a bad cough and runny nose. What should I do? | Based on the context provided, it seems you may be experiencing symptoms of the common cold. Here are some general self-care precautions you can take: 1. Drink vitamin C-rich drinks: Consuming foods and beverages high in vitamin C can help support your immune system and potentially alleviate cold symptoms. 2. Take vapour: Inhaling steam from a hot shower or a bowl of hot water can help relieve congestion and soothe your throat. 3. Avoid cold food: Consuming cold foods may worsen your symptoms, so try to eat warm or room temperature foods instead. 4. Keep fever in check: | Response: Dear [Patient], I am SafeSymp, a healthcare assistant who is here to provide you with general self-care advice for common cold and bronchial asthma. If you have a bad cough and runny nose, it's essential to take precautions to prevent further complications. Here are some self-care precautions you can follow: - Drink vitamin C-rich drinks such as oranges, grapefruits, or berries to boost your immune system. - Take a vapour to help relieve congestion and cough. - Avoid cold foods and drinks that contain alcohol, caffeine, or sp | Based on the given context, it seems that you may be experiencing a common cold or bronchial asthma. Here are some general self-care precautions that can help alleviate your symptoms: 1. Drink Vitamin C-rich drinks: Vitamin C is an antioxidant that helps boost the immune system and reduce inflammation. You can consume foods like citrus fruits, kiwi, strawberries, bell peppers, and broccoli, which are high in vitamin C. Drinking warm water with lemon juice or taking a supplement can also be beneficial. 2. Take Vapour: Using a humidifier or vaporizer can help moisten the air and relieve congestion. This can be especially helpful during a cold or bronchial asthma attack when the airways become dry and irritated. 3. Avoid cold food: Eating cold or spicy foods can irrit |
| 1 | What home precautions are recommended for malaria? | For malaria, the recommended home precautions include consulting the nearest hospital, avoiding oily food, avoiding non-vegetarian food, and keeping mosquitoes out of your living space. | Home precautions recommended for malaria include consulting a doctor for serious or persistent symptoms, avoiding oily food, avoiding non-vegetarian food, keeping mosquitoes out, and keeping hydrated. | Malaria is a disease caused by Plasmodium parasites transmitted through mosquitoes. To prevent malaria at home, the following home precautions can be taken: 1. Consult nearest hospital: If you suspect that you have been bitten by an infected mosquito, seek medical attention immediately. Early diagnosis and treatment can help prevent complications and reduce the risk of severe illness. 2. Avoid oily food: Mosquitoes lay their eggs on stagnant water, which can provide a breeding ground for the parasite. Consuming foods high in fat, salt, and sugar can attract mosquitoes and increase your risk of getting infected. Opt for lean protein sources like chicken, fish, and legumes, as well as whole grains and fruits, to maintain a healthy diet. 3. Avoid non-veg food: Some people may choose to |
| 2 | What self-care steps can I follow for allergy symptoms? | For allergy symptoms, you can follow these self-care steps: 1. Apply calamine to the affected area to soothe the itching and irritation. 2. Cover the area with a bandage to protect it from further irritation and to prevent scratching. 3. Use ice to compress the itching area, which can help reduce inflammation and provide temporary relief. Remember, if your symptoms are severe or persist, it's important to consult a doctor for proper diagnosis and treatment. | To follow self-care steps for allergy symptoms, you should apply calamine lotion, cover the affected area with a bandage, and use ice to reduce itching. You should also stop irritation and consult a doctor if the symptoms persist or worsen. It is essential to follow up with a doctor for further diagnosis and treatment. | For allergy symptoms, you can follow these general self-care steps: 1. Apply Calamine: Calamine lotion is an effective anti-inflammatory agent that helps soothe skin rashes, redness, and itching caused by allergies. Apply it directly to the affected area to provide relief. 2. Cover Area with Bandage: A clean, dry bandage can help protect the affected area from further irritation and prevent scratching, which may worsen the symptoms. Make sure to choose a breathable material like cotton or linen. 3. Use Ice to Compress Itching: Applying cold compresses to the affected area can help reduce swelling and inflammation, making it easier to manage itching. Wrap a clean cloth in ice cubes and place it on the affected area for 5-10 minutes at a time several times |
| 3 | How can I manage hypothyroidism safely at home? | To manage hypothyroidism safely at home, focus on reducing stress, engaging in regular exercise, maintaining a healthy diet, and ensuring you get proper sleep. These self-care precautions can help support your overall well-being. However, if you experience serious or persistent symptoms, it's important to consult a doctor for further evaluation and guidance. | To manage hypothyroidism safely at home, follow the general self-care advice provided in the context: - Reduce stress: Practice relaxation techniques such as deep breathing, meditation, or yoga. - Exercise: Engage in regular physical activity, such as walking, swimming, or cycling. - Eat healthy: Avoid processed foods, sugary drinks, and excessive salt intake. Eat a balanced diet that includes plenty of fruits, vegetables, and whole grains. - Get proper sleep: Establish a regular sleep schedule and avoid caffeine and nicotine. It is essential to consult a doctor if you experience persistent symptoms such as weight | Hypothyroidism is an autoimmune disorder that affects the thyroid gland, which produces hormones that regulate metabolism and growth. To manage hypothyroidism safely at home, you can follow these general self-care precautions: 1. Reduce Stress: High levels of stress can affect your thyroid function, so it's essential to find ways to manage stress. This may include practicing relaxation techniques such as deep breathing exercises, meditation, yoga, or progressive muscle relaxation. 2. Exercise: Regular physical activity can help maintain a healthy weight, improve blood flow, and boost the production of thyroid hormones. Aim for at least 30 minutes of moderate-intensity exercise most days of the week, such as brisk walking, cycling, swimming, or strength training. 3. Eat Healthy: A balanced diet rich in nutrients like |

## 10  Discussion

The results show that RAG greatly enhances the performance of small LLMs in healthcare contexts. The models differ significantly:

- **Qwen** is best for structured, comprehensive explanations.

- **Phi-3** is safest and most balanced.

- **TinyLlama** is ideal when inference speed is critical, but depth is not.

RAG reduces hallucinations by ensuring that answers are grounded in curated data. Safety prompts further prevent diagnosis or treatment recommendations.

---

## 11  Limitations

- Dataset size is small (41 rows).

- System only provides self-care precautions—not symptom origin or diagnosis.

- Cosine similarity may not fully capture semantic fidelity.

- Lightweight LLMs still hallucinate if retrieval fails.

---

## 12  Ethical Considerations

- SafeSymp is **not** a diagnostic tool.

- Designed only for **educational purposes**.

- Must include disclaimers in every response.

- Systems used in healthcare must undergo medical expert validation.

- Dataset bias may influence retrieved precautions.

| Response Quality & Reasoning | Factual Accuracy & Grounding | Safety & Harm |
|---|---|---|
| ✓ Most structured and coherent responses | ✓ Highest factual grounding with retrieval | ✓ Almost pefect compliance |
| ➖ Good structure but sometimes irrelevant content | ➖ Strong but sometimes fabricate rare causes | ➖ Occasional diagnostic-like statements |
| ✗ Provide generic, shallow reasoning | ✗ Most hallucination | ✗ Most safety violations |

## 13   Conclusion

SafeSymp demonstrates that small LLMs, when combined with RAG and safety prompts, can provide grounded, safe, and computationally efficient health self-care guidance. The combination of SBERT embeddings, FAISS retrieval, and lightweight LLMs allows deployment even on laptops or edge devices. This makes SafeSymp a promising framework for accessible, low-resource health-information assistants.

## 14   References

- Lewis et al., 2020 – Retrieval-Augmented Generation

- Karpukhin et al., 2020 – Dense Passage Retrieval

- Borgeaud et al., 2022 – RETRO

- Ji et al., 2023 – Hallucination in LLMs

- Microsoft Research – Phi-3 Model Report

- Qwen Model Technical Report

- TinyLlama Project Documentation