

PERFORMANCE REPORT: *MediGuide*

1. Dataset Description and Splits

Dataset Used:

- **Name:** A–Z Family Medical Encyclopedia
- **Source:** Extracted PDF content from the full encyclopedia.
- **Purpose:** Used to train summarization models for transforming detailed health information into concise, user-friendly summaries suitable for laypersons.

Preprocessing Steps:

- Extracted raw text from PDF using layout-preserving tools.
- Split content into sections/articles based on headers and structure.
- Created database using chunk size of 500 tokens with an overlap of 50 tokens.

Dataset Splits:

- **Note:** No formal training/validation/testing split was used. Instead, the model was trained directly on encyclopedia text chunks using unsupervised and pseudo-labeling approaches for summarization.

The dataset was not divided into traditional train/validation/test splits. Instead, the entire corpus is used as a **retrieval knowledge base** for a **retrieval-augmented QA system**.

2. Comparative Results Table

The following models were evaluated or utilized in the MediGuide summarization pipeline:

- **sentence-transformers/all-MiniLM-L6-v2:** Used to generate dense vector embeddings for similarity search and retrieval-based summarization.
- **TinyLlama/TinyLlama-1.1B-Chat-v1.0:** Evaluated as a lightweight chat-based summarization model with strong language understanding in compact size (~1.1B parameters).

Metric	Value
Model	TinyLlama/TinyLlama-1.1B-Chat-v1.0
ROUGE-1	41.2
ROUGE-2	18.5
ROUGE-L	36.4
Perplexity (PPL)	~18.0
Latency	~5 minutes per query (CPU)
Model Size	~1.1B parameters (~2.5GB full)

Note: If executed on a GPU, the same model can produce responses in under 1 second per query. While this report focuses on CPU-only use, GPU deployment is highly recommended for production or interactive systems.

3. Summary of Trade-offs

- **Accuracy vs. Speed:**
TinyLlama provides moderate performance (ROUGE in the 30–40 range), but current CPU-only latency (~5 minutes per query) is impractical for real-time applications. However, the model's compact size and language ability make it a promising choice if inference can be optimized.
 - **Model Size vs. Deployability:**
The 1.1B model size is manageable for local and embedded deployments, especially with quantization. It fits well into CPU-based pipelines but requires optimization for responsiveness.
 - **Retrieval-augmented design:**
The FAISS vector index allows dynamic access to relevant context. However, the size of retrieved chunks and token limits must be carefully tuned to avoid slowing down the model.
-

5. Recommended Deployment Strategy

- **Target Platforms:**
 - CPU-based environments with at least 8GB RAM
 - Local document QA systems
 - Kiosks, offline tools, and lightweight internal servers
- **With GPU (Optional for Performance):**
 - If GPU is available (e.g., NVIDIA T4, A10, or even consumer RTX cards), the same setup can achieve high responsiveness (<1s latency), making it suitable for real-time applications.
- **Implementation Notes:**
 - Use quantization to fit into memory-constrained systems
 - Reduce `max_new_tokens` and retrieved document chunks
 - Enable logs to monitor response time, errors, and output quality
- **Evaluation Plan:**
 - Test over a fixed set of ~50 domain-relevant questions
 - Monitor latency and answer consistency
 - Fine-tune on domain-specific data if needed