

# Knowledge-Augmented Multimodal Clinical Rationale Generation for Disease Diagnosis with Small Language Models

Anonymous ACL submission

## Abstract

Interpretation is critical for disease diagnosis, but existing models struggle to balance predictive accuracy with human-understandable rationales. While large language models (LLMs) offer strong reasoning abilities, their clinical use is limited by high computational costs and restricted multimodal reasoning ability. Small language models (SLMs) are efficient but lack advanced reasoning for integrating multimodal medical data. In addition, both LLMs and SLMs lack of domain knowledge for trustworthy reasoning. Therefore, we propose ClinRaGen, enhancing SLMs by leveraging LLM-derived reasoning ability via rationale distillation and domain knowledge injection for trustworthy multimodal rationale generation. Key innovations include a sequential rationale distillation framework that equips SLMs with LLM-comparable multimodal reasoning abilities, and a knowledge-augmented attention mechanism that jointly unifies multimodal representation from time series and textual data in a same encoding space, enabling it naturally interpreted by SLMs while incorporating domain knowledge for reliable rationale generation. Experiments on real-world medical datasets show that ClinRaGen achieves state-of-the-art performance in disease diagnosis and rationale generation, demonstrating the effectiveness of combining LLM-driven reasoning with knowledge augmentation for improved interpretability.

## 1 Introduction

The widespread adoption of electronic health records (EHRs) has transformed deep learning applications in healthcare by providing diverse data modalities, including medical notes, laboratory (lab) test results, and clinical events. These multimodal inputs are crucial for disease diagnosis, mortality prediction, and drug discovery (Niu et al., 2024; Laghuvarapu et al., 2024). Large language models (LLMs) have recently demonstrated

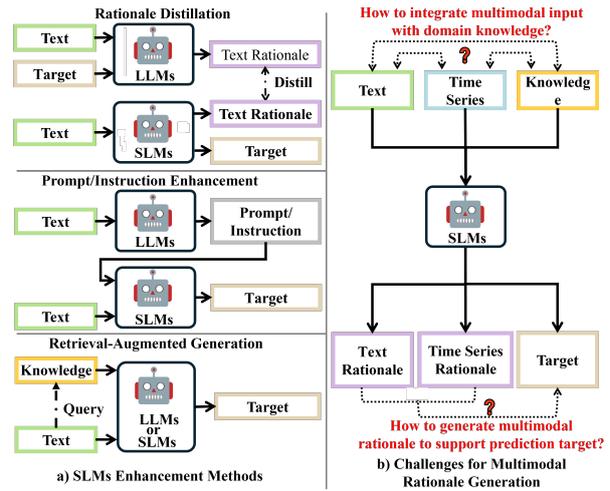


Figure 1: Existing SLM enhancement methods and challenges in multimodal rationale generation.

strong diagnostic performance and reasoning capabilities through techniques such as prompt learning and Chain-of-Thought (CoT) reasoning (Wei et al., 2022; Singhal et al., 2023; Chen et al., 2023). However, despite these advancements, LLMs face significant challenges in real-world clinical deployment due to high computational costs, the need for external domain-specific data integration, and difficulties in processing multimodal inputs—particularly numerical time-series lab test. More critically, LLMs lack the ability to generate clinically grounded multimodal rationales, limiting their interpretability in medical decision-making.

Small language models (SLMs) have emerged as a computationally efficient alternative, benefiting from recent advancements in rationale distillation, prompt learning, and retrieval-augmented generation (RAG) (Hsieh et al., 2023; Kang et al., 2024; Kwon et al., 2024). As shown in Figure 1a, these methods enable SLMs to inherit LLM-driven reasoning abilities, improve generalization through instruction-based adaptation, or leverage RAG for

064 more reliable outputs. However, as illustrated in  
065 Figure 1b, these approaches still suffer from two  
066 fundamental challenges. The first challenge is that  
067 they struggle to effectively integrate multimodal  
068 inputs with structured domain knowledge, as most  
069 methods focus on single-modality data (e.g., text-  
070 based rationales) rather than jointly processing text-  
071 ual and time series EHR data (Shi et al., 2024;  
072 Sohn et al., 2024). The second challenge is that  
073 they fail to provide coherent multimodal rationales  
074 that align with clinical decision-making, as ratio-  
075 nale generation often remains text-centric and lacks  
076 interpretability across different data modalities.

077 To bring the best of both worlds, we propose Clin-  
078 RaGen, a knowledge-augmented framework for  
079 multimodal clinical rationale generation. Clin-  
080 RaGen enhances SLMs’ trustworthy mutlimodal  
081 reasoning capabilities from two aspects. First,  
082 it transfers LLM-derived reasoning to SLMs  
083 through a sequential rationale distillation paradigm.  
084 Second, unlike approaches that rely solely on  
085 LLM-generated rationales (Kwon et al., 2024) or  
086 resource-intensive RAG (Kang et al., 2024), we  
087 propose a knowledge-augmented attention mecha-  
088 nism that achieves dual functionality: Efficient in-  
089 tegration of external medical knowledge to enable  
090 multimodal rationale generation grounded in clin-  
091 ical validity, ensuring the production of clinically  
092 meaningful explanations; Unification of time-series  
093 and textual EHRs within a shared encoding space,  
094 thereby enhancing multimodal representation learn-  
095 ing and facilitates interpretable decision-making.

096 The main contributions of this paper are:

- 097 • We propose ClinRaGen, a multimodal frame-  
098 work that transfer LLM reasoning capabilities  
099 into SLMs for disease diagnosis and clinical  
100 rationale generation, achieving both accuracy  
101 and interpretability.
- 102 • We introduce a knowledge-augmented atten-  
103 tion mechanism that jointly encodes time-  
104 series EHRs into clinical textual representa-  
105 tions while injecting domain knowledge, sig-  
106 nificantly improving mutlimodal rationale re-  
107 liability and accuracy.
- 108 • State-of-the-art performance in disease di-  
109 agnosis and rationale generation, validated  
110 through extensive experiments on benchmark  
111 EHR datasets (Johnson et al., 2016, 2023).

## 2 Related Work 112

Recent advancements in large-scale high-quality 113  
datasets and computational resources have enabled 114  
significant progress in Natural Language Process- 115  
ing (NLP), with improved training methodologies 116  
fueling the development of LLMs (Touvron et al., 117  
2023; Achiam et al., 2023). In healthcare, LLMs 118  
have been applied to clinical question answering 119  
and diagnostic reasoning (Singhal et al., 2023; 120  
Yang et al., 2022). While effective in text-based 121  
tasks, these models struggle to generate clinically 122  
grounded multimodal rationales. Medical-specific 123  
LLMs (Chen et al., 2023; Zhang et al., 2023) miti- 124  
gate this issue through domain adaptation, but their 125  
high computational costs and reliance on large- 126  
scale training data limit scalability. 127

To improve efficiency, rationale distillation trans- 128  
fers LLM-derived reasoning ability to SLMs, reduc- 129  
ing computational overhead while preserving inter- 130  
pretability (Hsieh et al., 2023; Ho et al., 2023; Kang 131  
et al., 2024). Chain-of-thought prompting further 132  
enhances SLM reasoning capabilities (Wei et al., 133  
2022). However, most distillation approaches re- 134  
main text-centric and lack robust multimodal EHRs 135  
integration (Kang et al., 2024; Ho et al., 2023). 136  
RAG has been explored to improve rationale reli- 137  
ability by incorporating external knowledge, yet 138  
retrieval latency and adaptability remain key chal- 139  
lenges (Jiang et al., 2025). Despite these advance- 140  
ments, multimodal rationale generation remains an 141  
open challenge. Current models struggle to fuse 142  
textual, time-series, and structured medical knowl- 143  
edge into coherent clinical rationales. 144

## 3 Methodology 145

We introduce ClinRaGen, a knowledge-augmented 146  
framework designed to enhance disease diagnosis 147  
and clinical rationale generation in SLMs by in- 148  
tegrating LLM-derived reasoning and structured 149  
domain knowledge. ClinRaGen bridges the gap be- 150  
tween large-scale medical knowledge and efficient 151  
multimodal reasoning, enabling the generation of 152  
two types of rationales: 1). medical note-based 153  
rationales ( $R^m$ ) and 2). lab test-based rationales 154  
( $R^l$ ) from medical notes ( $M$ ), time-series lab test 155  
results ( $T$ ), and disease-specific knowledge ( $K$ ). 156

ClinRaGen consists of two key components: 157  
Knowledge Retrieval and LLM-Guided Rationale 158  
Generation (Section 3.1), which collects domain 159  
knowledge and generates LLM-derived rationales 160

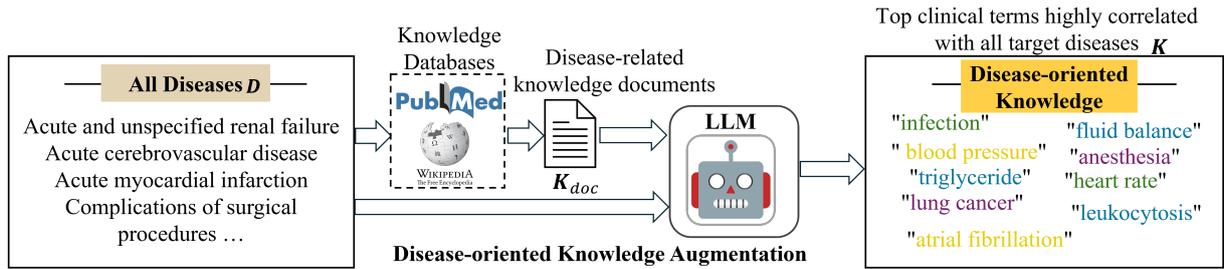


Figure 2: Knowledge augmentation in ClinRaGen. Given diagnosed diseases  $D$ , relevant descriptions  $K_{doc}$  are retrieved and processed by an LLM to extract key clinical terms  $K$ , enhancing multimodal rationale generation.

as distillation data for subsequent model training, and Knowledge-augmented Attention with Sequential Multimodal Rationale Distillation (Section 3.2), which progressively integrates structured knowledge to enhance multimodal reasoning in SLMs.

### 3.1 Knowledge Retrieval and LLM-Guided Rationale Generation

This step focuses on gathering domain knowledge and leveraging LLMs to generate structured rationales. The generated rationales serve as distillation targets for training SLMs in later stages. This ensures that SLMs receive high-quality, structured reasoning data to develop robust multimodal reasoning capabilities.

#### 3.1.1 Collecting Domain-Specific Medical Knowledge

LLMs encode extensive medical knowledge but are computationally expensive and impractical for direct deployment. Meanwhile, SLMs such as Flan-T5 and Flan-PaLM (Chung et al., 2024) are computationally efficient but lack sufficient domain-specific expertise to perform complex medical reasoning (Kang et al., 2024; Ho et al., 2023). To bridge this gap, ClinRaGen retrieves relevant medical knowledge from external sources and structures it for integration into SLM training.

As shown in Figure 2, ClinRaGen collects disease-related documents  $K_{doc}$  from PubMed<sup>1</sup> and Wikipedia<sup>2</sup>, extracting key medical terms using LLM-based processing to construct a structured knowledge base  $K$ :

$$K = \operatorname{argmax}_{K'} P_{LLM}(K' | D, K_{doc}). \quad (1)$$

This structured knowledge base is not used directly by the SLM during inference but instead supports

rationale generation in the next step. The retrieval and extraction process iterates until a stable set of key medical terms is obtained.

#### 3.1.2 Generating Rationales for Distillation

ClinRaGen employs LLMs to generate structured rationales that serve as distillation targets for SLM training. Unlike direct knowledge retrieval (Kang et al., 2024; Jiang et al., 2025), this step synthesizes structured explanations that explicitly link medical knowledge with clinical decision-making, enabling SLMs to internalize complex reasoning patterns during later training stages. To construct high-quality rationale data, we collaborated with clinicians to curate representative EHR samples and formulate corresponding gold-standard rationales  $O$ . These rationales guide the LLM in generating structured explanations, ensuring that the distilled knowledge supports multimodal reasoning. To improve LLM comprehension of numerical lab test data, we applied anomaly detection (Vinutha et al., 2018) and designed structured prompts that convert numerical values into interpretable textual explanations  $T^*$  (see Appendix A.1).

Figure 3 illustrates the multimodal rationale generation process. ClinRaGen first generates rationales ( $R^m$ ) based on medical notes:

$$R^m = \operatorname{argmax}_{R'} P_{LLM}(R' | M, D, O). \quad (2)$$

Then, lab test-based rationales ( $R^l$ ) are generated using insights from both medical notes, time series anomalies, and the generate note-based rationales:

$$R^l = \operatorname{argmax}_{R'} P_{LLM}(R' | M, T^*, D, O, R^m). \quad (3)$$

These LLM-generated rationales form the foundation of the subsequent distillation process (detailed in Section 3.2) and enable SLMs to learn structured, multimodal reasoning efficiently. For further

<sup>1</sup><https://pubmed.ncbi.nlm.nih.gov/>

<sup>2</sup><https://www.wikipedia.org/>

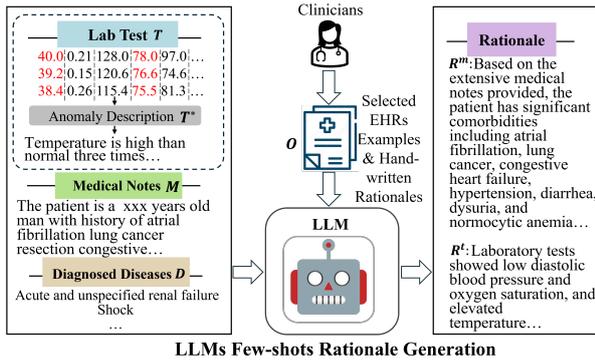


Figure 3: LLM-based clinical rationale generation. Medical notes ( $M$ ), lab test results ( $T$  and  $T^*$ ), diagnosis ( $D$ ), and clinicians provide examples ( $O$ ) are used to produce medical note-based ( $R^m$ ) and lab test-based ( $R^t$ ) rationales.

details on data processing and prompt engineering, refer to Appendix A.2.

### 3.2 Multimodal Rationale Distillation

Figure 4a presents the ClinRaGen framework, which comprises a Time Series Encoder for processing numerical lab test data, a Knowledge-Augmented Attention Module for integrating structured domain knowledge, and a SLM for generating disease diagnoses and structured clinical multimodal rationales. The framework enables progressive multimodal reasoning by leveraging structured knowledge and sequential learning mechanisms.

As illustrated in Figure 4b, ClinRaGen employs a three-phase rationale distillation paradigm that systematically integrates textual, numerical, and structured domain knowledge. The first phase distills medical note-based rationales, allowing the SLM to develop a foundational understanding of textual clinical information. The second phase introduces knowledge-augmented attention, aligning numerical lab test with structured medical knowledge to for distilling lab test-based rationales. The final phase fully integrates textual and numerical inputs, enabling the SLM to generate clinically coherent multimodal rationales to support disease diagnosis.

#### 3.2.1 Phase 1: Rationale Distillation from Medical Notes

In the first phase, the SLM is trained exclusively on medical notes  $M$  to establish a foundational understanding of clinical reasoning. This stage enables the model to generate disease diagnoses  $D$  while also producing medical note-based rationales  $R^m$

and lab test-based rationales  $R^t$ . By learning to extract meaningful insights from structured textual data, the SLM develops its initial ability to infer clinical relationships. The model is trained using a language model generation objective:

$$\mathcal{L}_{note}(\theta) = \mathbb{E}[-\log P_{SLM_\theta}(D, R^m, R^t | M)], \quad (4)$$

where  $\theta$  represents the trainable parameters of the SLM. This phase not only enables the model to internalize explicit diagnostic reasoning from medical notes but also allows it to implicitly capture latent patterns associated with lab test results, laying the groundwork for multimodal integration in subsequent phases.

#### 3.2.2 Phase 2: Knowledge Injection and Time-Series Rationale Distillation

To enable the SLM to effectively interpret numerical lab test data and generate time-series-based rationales ( $R^t$ ) that support disease diagnosis ( $D$ ), we introduce a Knowledge-Augmented Attention (KA) Module. This mechanism integrates domain-specific medical knowledge into the reasoning process, enhancing the model’s ability to produce clinically coherent and robust multimodal rationales.

We first use a Time Series Encoder (TSE) to encode raw lab test values  $T$  into structured hidden embeddings  $T^e$ . To align domain knowledge with the SLM, we construct a domain-specific vocabulary  $V^k$  by filtering standard language vocabulary  $V$  based on structured medical knowledge  $K$ :

$$V^k = \{v_1, \dots, v_n \mid v_1 \in K, K \subseteq V\}. \quad (5)$$

A cross-attention mechanism is then applied to integrate knowledge-driven representations into the model. The lab test embeddings ( $T^e$ ) serve as the Query, while domain knowledge tokens ( $V^k$ ) act as the Key and Value:

$$\begin{aligned} H &= f_\phi(T, V^k), \\ &= \text{SoftMax}\left(\frac{(T^e W^q)(V^k \top W^k)}{\sqrt{d}}\right)(V^k W^v), \end{aligned} \quad (6)$$

where  $d$  is the hidden dimension of the SLM, and  $W^q, W^k, W^v$  are learnable attention weight matrices,  $f$  indicates the encoding function of the TSE and attention, and  $\phi$  represents the trainable parameters of  $f$ . The resulting knowledge-enhanced embeddings  $H$  are then fed into the frozen distilled SLM to refine its reasoning and generate lab test-based rationales ( $R^t$ ) and diagnosis ( $D$ ).

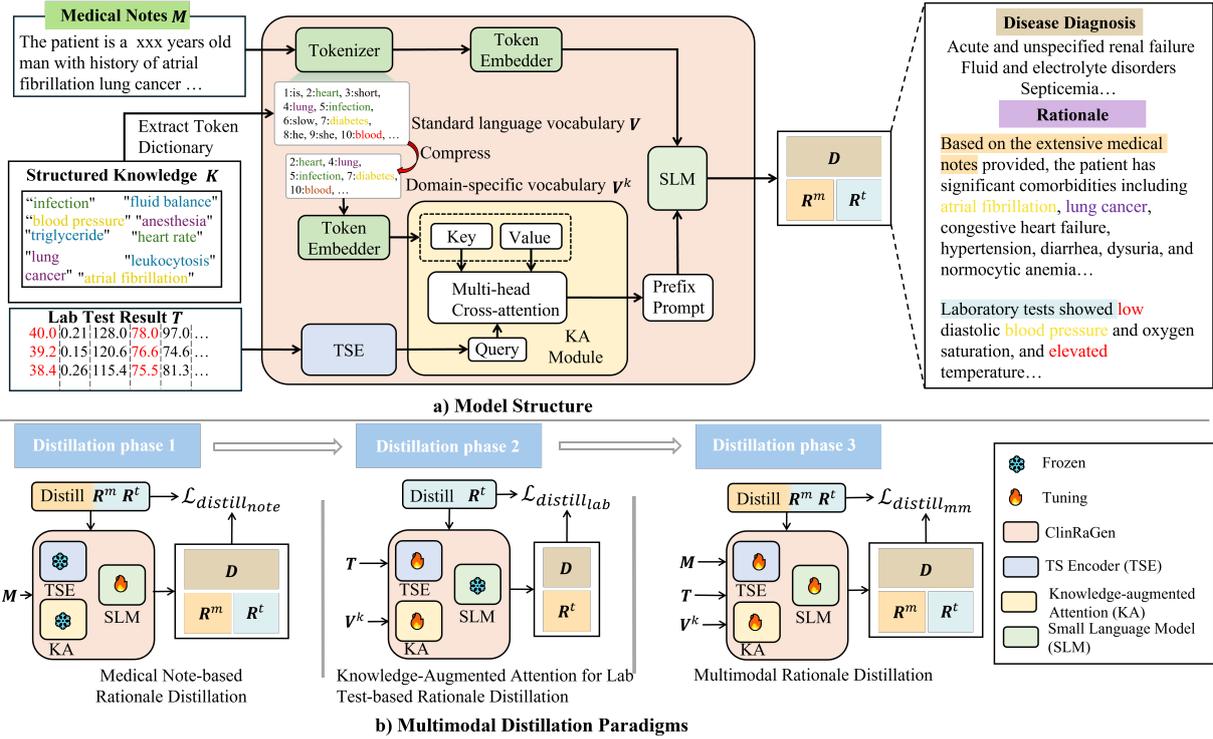


Figure 4: Overview of ClinRaGen. (a) Model structure comprising a time series encoder, knowledge-augmented attention module, and a SLM. (b) Three-phase rationale distillation: Medical Note-based Rationale Distillation, Knowledge-Augmented Attention for Lab Test-based Rationale Distillation, and Multimodal Rationale Distillation.

The model is trained using the following objective function:

$$\mathcal{L}_{lab}(\phi) = \mathbb{E}[-\log P_{SLM_\theta}(\mathbf{D}, \mathbf{R}^t | \mathbf{H})]. \quad (7)$$

This phase ensures that the SLM can naturally interpret lab test while effectively leveraging medical knowledge to enhance its reasoning capabilities.

### 3.2.3 Phase 3: Full Multimodal Rationale Distillation

In the final phase, ClinRaGen is trained to generate full multimodal clinical rationales by integrating medical notes, lab test, and structured domain knowledge. To ensure effective multimodal reasoning, lab test  $T$  is formatted as prefix prompts (Niu et al., 2024), allowing the SLM seamlessly incorporates it with textual EHRs.

During this stage, the model is optimized to generate both medical note-based rationales ( $R^m$ ) and lab test-based rationales ( $R^t$ ), ensuring that all available information contributes to clinically coherent and interpretable decision-making. The multimodal rationale distillation objective is formulated as follows:

$$\mathcal{L}_{mm}(\theta, \phi) = \mathbb{E}[-\log P_{SLM_\theta}(\mathbf{D}, \mathbf{R}^m, \mathbf{R}^t | \mathbf{M}, f_\phi(\mathbf{T}, \mathbf{V}^k))]. \quad (8)$$

The fine-tuning of all ClinRaGen components, ensuring that multimodal EHRs are effectively integrated, enhances diagnostic accuracy and produces modality-consistent rationales.

## 4 Experiments

### 4.1 Experimental Settings

**Dataset:** We evaluate ClinRaGen on two public EHR datasets: MIMIC-III (Johnson et al., 2016) (28,456 EHRs include medical notes and time series lab tests) and MIMIC-IV (Johnson et al., 2023) (28,900 EHRs). Both datasets use benchmark tools (Harutyunyan et al., 2019) for time series processing, with missing values filled by nearest available data. We target 25 disease phenotypes and follow a 4:1 training-to-testing split (Harutyunyan et al., 2019). Our model is available at github<sup>3</sup>.

**Baseline Methods:** To evaluate the effectiveness of ClinRaGen for disease diagnosis generation, we compared it with following baselines: Flan-T5 (Chung et al., 2024), PROMPTEHR (Wang and Sun, 2022), FROZEN (Tsimpoukelli et al., 2021), EHR-KnowGen (Niu et al., 2024), Clinical CoT

<sup>3</sup><https://anonymous.4open.science/r/ClinRaGen-6C9D/>

Models	Size	Modality		Micro			Macro		
		Lab	Note	Precision	Recall	F1	Precision	Recall	F1
<b>MIMIC-III</b>									
Flan-T5	60M		✓	0.5812 <sub>(0.11)</sub>	0.6623 <sub>(0.07)</sub>	0.6203 <sub>(0.05)</sub>	0.5656 <sub>(0.10)</sub>	0.6247 <sub>(0.08)</sub>	0.5887 <sub>(0.07)</sub>
PROMPTEHR	75.2M		✓	0.5929 <sub>(0.11)</sub>	0.6553 <sub>(0.07)</sub>	0.6224 <sub>(0.02)</sub>	0.5744 <sub>(0.10)</sub>	0.6287 <sub>(0.06)</sub>	0.5910 <sub>(0.03)</sub>
LLaMA-ft	7B	✓	✓	0.6142 <sub>(0.21)</sub>	0.6598 <sub>(0.15)</sub>	0.6364 <sub>(0.04)</sub>	0.6108 <sub>(0.15)</sub>	0.6164 <sub>(0.13)</sub>	0.6055 <sub>(0.04)</sub>
FROZEN	265M	✓	✓	0.6102 <sub>(0.18)</sub>	0.6401 <sub>(0.16)</sub>	0.6231 <sub>(0.03)</sub>	0.5976 <sub>(0.16)</sub>	0.6001 <sub>(0.17)</sub>	0.5915 <sub>(0.03)</sub>
EHR-KnowGen	77M	✓	✓	0.6001 <sub>(0.03)</sub>	0.6551 <sub>(0.02)</sub>	0.6262 <sub>(0.01)</sub>	0.5834 <sub>(0.04)</sub>	0.6181 <sub>(0.03)</sub>	0.5944 <sub>(0.01)</sub>
Clinical CoT									
-w/o TSE	60M		✓	0.6115 <sub>(0.03)</sub>	0.6402 <sub>(0.04)</sub>	0.6311 <sub>(0.03)</sub>	0.6024 <sub>(0.04)</sub>	0.5989 <sub>(0.06)</sub>	0.5969 <sub>(0.03)</sub>
-w/ TSE	85M	✓	✓	0.5967 <sub>(0.05)</sub>	0.6607 <sub>(0.06)</sub>	0.6328 <sub>(0.03)</sub>	0.5924 <sub>(0.06)</sub>	0.6092 <sub>(0.07)</sub>	0.5975 <sub>(0.05)</sub>
LLM Zero-shot									
-LLaMA	7B	✓	✓	0.1227 <sub>(0.08)</sub>	0.0421 <sub>(0.06)</sub>	0.0627 <sub>(0.06)</sub>	0.0392 <sub>(0.06)</sub>	0.0622 <sub>(0.06)</sub>	0.0438 <sub>(0.05)</sub>
-ChatGPT	175B	✓	✓	0.4474 <sub>(0.07)</sub>	0.1405 <sub>(0.05)</sub>	0.2139 <sub>(0.05)</sub>	0.4883 <sub>(0.08)</sub>	0.1872 <sub>(0.05)</sub>	0.2188 <sub>(0.04)</sub>
ClinRaGen	87M	✓	✓	0.6104 <sub>(0.02)</sub>	0.6751 <sub>(0.02)</sub>	<u>0.6410</u> <sub>(0.01)</sub>	0.5991 <sub>(0.03)</sub>	0.6311 <sub>(0.04)</sub>	<u>0.6113</u> <sub>(0.02)</sub>
ClinRaGen*	793M	✓	✓	0.6047 <sub>(0.03)</sub>	0.6875 <sub>(0.03)</sub>	<b>0.6501</b> <sub>(0.02)</sub>	0.5943 <sub>(0.04)</sub>	0.6531 <sub>(0.03)</sub>	<b>0.6196</b> <sub>(0.03)</sub>
<b>MIMIC-IV</b>									
Flan-T5	60M		✓	0.6624 <sub>(0.05)</sub>	0.6953 <sub>(0.02)</sub>	0.6792 <sub>(0.04)</sub>	0.6428 <sub>(0.06)</sub>	0.6601 <sub>(0.05)</sub>	0.6479 <sub>(0.04)</sub>
PROMPTEHR	75.2M		✓	0.6524 <sub>(0.07)</sub>	0.7031 <sub>(0.06)</sub>	0.6802 <sub>(0.02)</sub>	0.6353 <sub>(0.05)</sub>	0.6702 <sub>(0.07)</sub>	0.6501 <sub>(0.03)</sub>
LLaMA-ft	7B	✓	✓	0.6854 <sub>(0.11)</sub>	0.6954 <sub>(0.07)</sub>	0.6929 <sub>(0.03)</sub>	0.6753 <sub>(0.09)</sub>	0.6624 <sub>(0.11)</sub>	0.6621 <sub>(0.06)</sub>
FROZEN	265M	✓	✓	0.6781 <sub>(0.08)</sub>	0.6908 <sub>(0.09)</sub>	0.6842 <sub>(0.01)</sub>	0.6627 <sub>(0.10)</sub>	0.6521 <sub>(0.10)</sub>	0.6530 <sub>(0.02)</sub>
EHR-KnowGen	77M	✓	✓	0.6580 <sub>(0.06)</sub>	0.7085 <sub>(0.05)</sub>	0.6816 <sub>(0.02)</sub>	0.6382 <sub>(0.05)</sub>	0.6724 <sub>(0.06)</sub>	0.6511 <sub>(0.02)</sub>
Clinical CoT									
-w/o TSE	60M		✓	0.6751 <sub>(0.05)</sub>	0.7069 <sub>(0.03)</sub>	0.6905 <sub>(0.03)</sub>	0.6607 <sub>(0.04)</sub>	0.6796 <sub>(0.06)</sub>	0.6612 <sub>(0.02)</sub>
-w/ TSE	85M	✓	✓	0.7011 <sub>(0.04)</sub>	0.6808 <sub>(0.06)</sub>	0.6917 <sub>(0.04)</sub>	0.6971 <sub>(0.05)</sub>	0.6354 <sub>(0.03)</sub>	0.6577 <sub>(0.03)</sub>
LLM Zero-shot									
-LLaMA	7B	✓	✓	0.1357 <sub>(0.11)</sub>	0.0997 <sub>(0.07)</sub>	0.1150 <sub>(0.06)</sub>	0.0435 <sub>(0.09)</sub>	0.1466 <sub>(0.07)</sub>	0.0619 <sub>(0.05)</sub>
-ChatGPT	175B	✓	✓	0.4536 <sub>(0.07)</sub>	0.1458 <sub>(0.05)</sub>	0.2207 <sub>(0.04)</sub>	0.4532 <sub>(0.06)</sub>	0.1831 <sub>(0.06)</sub>	0.2147 <sub>(0.05)</sub>
ClinRaGen	87M	✓	✓	0.7009 <sub>(0.01)</sub>	0.6963 <sub>(0.02)</sub>	<u>0.6989</u> <sub>(0.01)</sub>	0.6868 <sub>(0.03)</sub>	0.6603 <sub>(0.01)</sub>	<u>0.6685</u> <sub>(0.02)</sub>
ClinRaGen*	793M	✓	✓	0.6848 <sub>(0.04)</sub>	0.7429 <sub>(0.02)</sub>	<b>0.7127</b> <sub>(0.02)</sub>	0.6779 <sub>(0.02)</sub>	0.7087 <sub>(0.01)</sub>	<b>0.6893</b> <sub>(0.01)</sub>

Table 1: The performance of comparative methods in the disease diagnosis tasks on MIMIC-III and MIMIC-IV. The best results are highlighted in bold, and the second-best results are marked with an underline.

(with/without TSE) (Kwon et al., 2024), and LLM-based models LLaMA-7B (Touvron et al., 2023) (zero-shot and fine-tuning) and ChatGPT (Open, 2023) (zero-shot). Baseline and implementation details are provided in Appendices A.3 and A.4. For a fair comparison, all baselines (except LLaMA) use Flan-T5-Small as the backbone; our model employs Flan-T5-Small (ClinRaGen) and Flan-T5-Large (ClinRaGen\*) for evaluate effect of varying scales. ChatGPT (GPT-3.5-turbo) serves as our teacher LLM. Results are averaged over five runs with statistical significance determined at  $p < 0.05$  by t-test.

## 4.2 Disease Diagnosis Performance

**Comparison with Baselines:** We evaluate disease diagnosis using micro and macro precision, recall, and F1 scores. Table 1 shows that multimodal models outperform single-modality models, confirming the value of lab test results. Clinical CoT surpasses other baseline models, highlighting rationale distillation’s effectiveness. ClinRaGen (80M) achieves the best performance, with an average F1 score improvement of over 1.1% across all baselines, even

outperforming LLaMA-7B-ft. Furthermore, ClinRaGen\* (793M) improves by over 1.5%, significantly exceeding other baselines. The weak performance of zero-shot LLMs confirms the absence of data leakage. These results demonstrate ClinRaGen’s ability to match or surpass LLMs in clinical tasks through multimodal rationale distillation and the knowledge-augmented attention mechanism.

**Ablation Study:** We assess the impact of key components in ClinRaGen: (1) *w/o* LAB&KNOW removes lab tests and knowledge input, (2) *w/o* KNOW replaces the knowledge-based vocabulary with a standard one, and (3) *w/o* REASONING excludes rationale distillation while maintaining model structure. Table 2 shows that *w/o* REASONING performs worst, highlighting the importance of rationale distillation. The drop in F1 scores for *w/o* LAB&KNOW confirms the value of multimodal integration, while *w/o* KNOW shows the KA module’s contribution to diagnostic accuracy.

**Model Efficiency:** We evaluate ClinRaGen’s efficiency by comparing model parameters, micro F1 scores (Figure 5), and training times (Table 3).

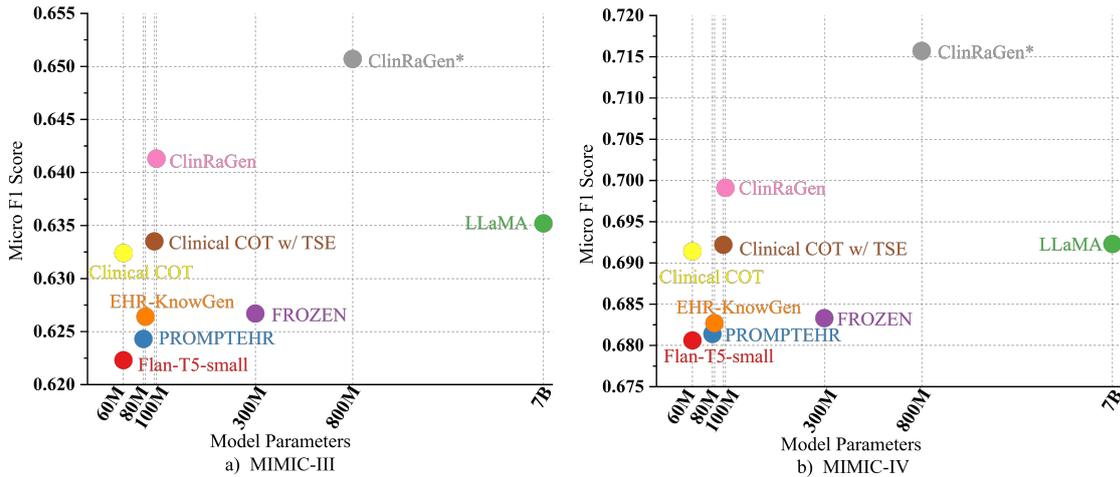


Figure 5: Model Parameter Counts and Micro F1 Scores

Models	Micro F1	Macro F1
<b>MIMIC-III</b>		
ClinRaGen	<b>0.6410</b>	<b>0.6113</b>
w/o LAB&KNOW	0.6323	0.6021
w/o KNOW	0.6349	0.6042
w/o REASONING	0.6255	0.5915
<b>MIMIC-IV</b>		
ClinRaGen	<b>0.6989</b>	<b>0.6685</b>
w/o LAB&KNOW	0.6925	0.6643
w/o KNOW	0.6936	0.6644
w/o REASONING	0.6828	0.6541

Table 2: Ablation studies on disease diagnosis.

ClinRaGen-Small achieves superior diagnostic performance with 80× fewer parameters and less than half the training time of LLaMA. These results highlight the effectiveness of our sequential multimodal distillation paradigm and KA mechanism in enabling efficient and accurate clinical reasoning.

### 4.3 Rationale Generation Performance

**Evaluation Methods:** To assess the quality of generated multimodal rationales and maximize the potential of SLMs, we evaluate ClinRaGen (80M) using five evaluation criteria—*Correctness, Readability, Soundness, Consistency, and Persuasiveness*—based on clinicians and prior research (Lin et al., 2024; Kwon et al., 2024). Scores range from 1 to 5 on a Likert scale (details of criteria defined in Appendix A.5). We conduct both LLM-based and human evaluations. For LLM comparisons, we use Mistral-7B, LLaMA2-7B, and LLaMA3-8B with five-shot prompting. Distilled rationales from ChatGPT serve as ground truth (GT). Comparative LLMs receive time series anomalies and

medical notes, while ClinRaGen directly processes numerical lab test and medical notes. Following Lin et al. (2024); Chiang and Lee (2023), we use GPT-4 to evaluate 1000 randomly selected samples. For human assessment, 15 professional post-graduates rate 100 samples, achieving moderate intra-class (0.637) and inter-class (0.608) agreement, indicating reasonable consistency despite the task’s subjectivity.

Models	Time Cost (Seconds)
Knowledge Retrieval	12,636
LLM-Guided Rationale Generation	604,715
LlaMA – 7B Tuning	259,113
ClinRaGen – 84M Tuning	94,623

Table 3: Time Cost Evaluation.

**Evaluation Results:** Figures 6(a) and (b) show GPT-4 and human evaluations across five criteria, with closely aligned results. LLaMA3 performs best, benefiting from its large scale and pre-training. ClinRaGen ranks second, matching LLaMA3 in readability and correctness while surpassing LLaMA2 and Mistral, which often generate incoherent rationales. Unlike other LLMs relying on anomaly captions, ClinRaGen achieves the second-highest consistency score, demonstrating the KA mechanism’s effectiveness in consist multimodal reasoning. ClinRaGen also outperforms LLaMA2 and Mistral in soundness and persuasiveness, further underscoring our method’s effectiveness. Appendix A.6 further validates rationale quality using BLEU (Papineni et al., 2002) and BERTScore (Zhang et al., 2019).

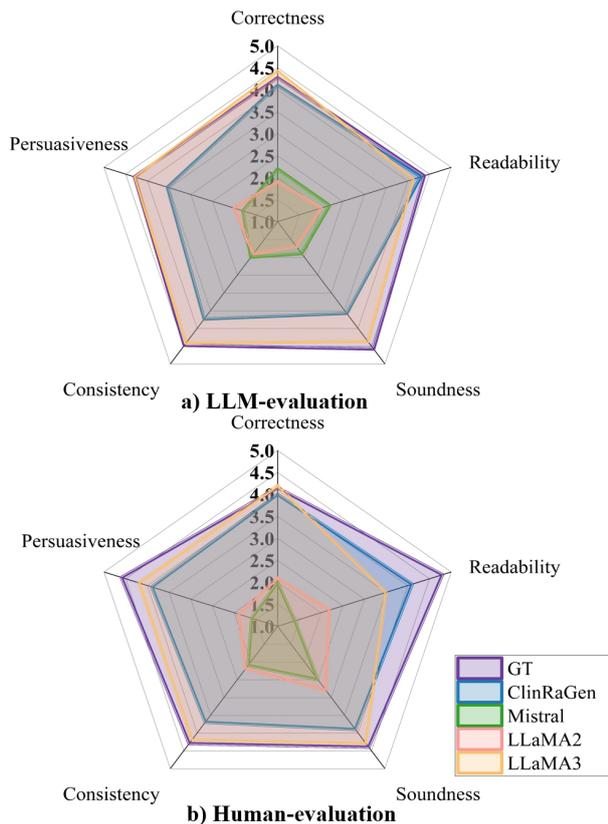


Figure 6: Clinical rationale generation evaluation

**Case Studies:** As illustrated in Figure 7, our model ClinRaGen can produce both medical note-based rationales (e.g., “Based on the medical notes...”) and lab test-based rationales (e.g., “Lab test shows...”), akin to the outputs of teacher LLM. For medical note-based rationale generation, ClinRaGen effectively extracts key medical terms essential for disease diagnosis (highlighted in green). Additionally, for lab test-based rationales, our model accurately identifies abnormal lab test features (highlighted in blue), demonstrating its capability to understand numerical time series lab test data effectively. These results indicate that ClinRaGen competently produces clinically relevant multimodal rationales to support disease diagnosis.

**Further Discussions:** LLMs may introduce bias into distilled clinical rationales. To evaluate ClinRaGen’s correctness, we assess the relevance of key medical terms to diagnosed diseases. In one case, our model identifies *weakness*, *lethargy*, and *basal ganglia hemorrhage* as evidence for *acute cerebrovascular*, while teacher LLM captures only *basal ganglia hemorrhage*, missing relevant symptoms (Unnithan et al., 2023). In another case, while the teacher LLM reports *no disease*, ClinRa-

**Example 1**  
**Disease Diagnosed:** Acute cerebrovascular disease, Complications of surgical procedures or medical care, Disorders of lipid metabolism  
**LLM Rationale:** Based on the medical notes, the patient... The presence of basal ganglia hemorrhage raises suspicion for acute cerebrovascular disease... Lab test shows repeated elevated diastolic and mean blood pressures indicating acute cerebrovascular disease post craniotomy...  
**ClinRaGen Rationale:** Based on the medical notes provided, the patient... The patient exhibited signs of acute cerebrovascular disease, including weakness, lethargy, and basal ganglia hemorrhage. Management included heparin drip, antiplatelet therapy, and blood pressure control... Lab test shows transiently elevated diastolic blood pressure and mean blood pressure post craniotomy for aneurysm clipping...

**Example 2**  
**Disease Diagnosed:** Chronic obstructive pulmonary disease and bronchiectasis, Coronary atherosclerosis and other heart disease, Disorders of lipid metabolism, Essential hypertension, Cardiac dysrhythmias  
**LLM Rationale:** Based on the medical notes provided, the patient ... hypertension, or cardiac dysrhythmias in the provided medical notes. Therefore, the rationale for diagnosis is “no disease was diagnosed” due to the focus on the acute neurosurgical intervention and recovery... Lab test shows episodes of significantly elevated fraction inspired oxygen, supportive of hypoxemia following...  
**ClinRaGen Rationale:** Based on the medical notes provided, the patient .... Therefore, the comprehensive rationale for the diagnosed diseases includes acute cerebrovascular disease, disorders of lipid metabolism, and essential hypertension... Lab test shows transiently elevated fraction inspired oxygen and increased mean blood pressure, consistent...

Figure 7: Case studies on disease diagnosis and clinical rationale generation compared with teacher LLM.

Gen correctly identifies conditions like *disorders of lipid metabolism* and *essential hypertension*. These results highlight ClinRaGen’s ability to mitigate LLM biases by capturing time-series variations and integrating structured knowledge.

## 5 Conclusion and Future Work

We present ClinRaGen, a knowledge-augmented framework that enhances SLMs with LLM-derived reasoning and structured medical knowledge for disease diagnosis and multimodal rationale generation. It introduces a knowledge-augmented attention module that jointly unifies time-series and textual EHRs in the same encoding space while injecting domain knowledge for reliable rationale generation and a sequential multimodal distillation paradigm for transferring LLMs’ reasoning capabilities to SLMs. Extensive evaluations on real world datasets, including quantitative and qualitative analyses, show that ClinRaGen enables SLMs to achieve LLM-comparable performance in disease diagnosis and multimodal rationale generation. This work bridges the performance gap between LLMs and SLMs in clinical tasks. Future research will extend ClinRaGen to a broader range of SLM architectures, datasets, and medical applications.

## 494 Limitations

495 While ClinRaGen effectively enhances multimodal  
496 clinical reasoning, certain limitations remain:

- 497 • First, although rationale distillation transfers  
498 reasoning capabilities from LLMs to SLMs,  
499 potential biases in LLM-generated rationales  
500 may persist.
- 501 • Second, the effectiveness of the knowledge-  
502 augmented attention module depends on the  
503 quality and coverage of external knowledge  
504 sources.
- 505 • Lastly, ClinRaGen is evaluated on structured  
506 EHR datasets, and its applicability to unstruc-  
507 tured clinical text or other medical modalities  
508 requires further exploration.

509 Future work will refine knowledge integration, en-  
510 hance bias mitigation strategies, and extend evalua-  
511 tions to diverse clinical settings.

## 512 Ethics Statement

513 **Data Privacy:** The datasets utilized in our re-  
514 search are publicly accessible and feature de-  
515 identified patient data, accessing these datasets still  
516 requires passing the CITI Exam<sup>4</sup> and download  
517 from Physionet<sup>5</sup>. In addition, this study used the  
518 Azure OpenAI service and completed the “opting  
519 out of the review process” agreement.

## 520 References

521 Josh Achiam, Steven Adler, Sandhini Agarwal, Lama  
522 Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo  
523 Almeida, Janko Altenschmidt, Sam Altman, Shyamal  
524 Anadkat, et al. 2023. Gpt-4 technical report. *arXiv*  
525 *preprint arXiv:2303.08774*.

526 Zeming Chen, Alejandro Hernández Cano, Angelika  
527 Romanou, Antoine Bonnet, Kyle Matoba, Francesco  
528 Salvi, Matteo Pagliardini, Simin Fan, Andreas Köpf,  
529 Amirkeivan Mohtashami, et al. 2023. Meditron-70b:  
530 Scaling medical pretraining for large language models.  
531 *arXiv preprint arXiv:2311.16079*.

532 Cheng-Han Chiang and Hung-Yi Lee. 2023. Can large  
533 language models be an alternative to human evaluations?  
534 In *Proceedings of the 61st Annual Meeting of the Asso-*  
535 *ciation for Computational Linguistics (Volume 1: Long*  
536 *Papers)*, pages 15607–15631.

537 Hyung Won Chung, Le Hou, Shayne Longpre, Bar-  
538 ret Zoph, Yi Tay, William Fedus, Yunxuan Li, Xuezhi

Wang, Mostafa Dehghani, Siddhartha Brahma, et al. 2024. Scaling instruction-finetuned language models. *Journal of Machine Learning Research*, 25(70):1–53. 539 540 541

Hrayr Harutyunyan, Hrant Khachatryan, David C. Kale, Greg Ver Steeg, and Aram Galstyan. 2019. Multitask learning and benchmarking with clinical time series data. *Scientific Data*, 6(1):96. 542 543 544 545

Namgyu Ho, Laura Schmid, and Se-Young Yun. 2023. Large language models are reasoning teachers. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 14852–14882. 546 547 548 549 550

Cheng-Yu Hsieh, Chun-Liang Li, CHIH-KUAN YEH, Hootan Nakhost, Yasuhisa Fujii, Alex Jason Ratner, Ranjay Krishna, Chen-Yu Lee, and Tomas Pfister. 2023. Distilling step-by-step! outperforming larger language models with less training data and smaller model sizes. In *The 61st Annual Meeting Of The Association For Computational Linguistics*. 551 552 553 554 555 556 557

Pengcheng Jiang, Cao Xiao, Minhao Jiang, Parminder Bhatia, Taha Kass-Hout, Jimeng Sun, and Jiawei Han. 2025. Reasoning-enhanced healthcare predictions with knowledge graph community retrieval. *ICLR*. 558 559 560 561

Alistair EW Johnson, Lucas Bulgarelli, Lu Shen, Alvin Gayles, Ayad Shammout, Steven Horng, Tom J Pollard, Sicheng Hao, Benjamin Moody, Brian Gow, et al. 2023. MIMIC-IV, a freely accessible electronic health record dataset. *Scientific data*, 10(1):1. 562 563 564 565 566

Alistair EW Johnson, Tom J Pollard, Lu Shen, H Lehman Li-Wei, Mengling Feng, Mohammad Ghassemi, Benjamin Moody, Peter Szolovits, Leo Anthony Celi, and Roger G Mark. 2016. MIMIC-III, a freely accessible critical care database. *Scientific data*, 3(1):1–9. 567 568 569 570 571

Minki Kang, Seanie Lee, Jinheon Baek, Kenji Kawaguchi, and Sung Ju Hwang. 2024. Knowledge-augmented reasoning distillation for small language models in knowledge-intensive tasks. *Advances in Neural Information Processing Systems*, 36. 572 573 574 575 576

Taeyoon Kwon, Kai Tzu-iunn Ong, Dongjin Kang, Seungjun Moon, Jeong Ryong Lee, Dosik Hwang, Beomseok Sohn, Yongsik Sim, Dongha Lee, and Jinyoung Yeo. 2024. Large language models are clinical reasoners: Reasoning-aware diagnosis framework with prompt-generated rationales. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 38, pages 18417–18425. 577 578 579 580 581 582 583 584

Siddhartha Laghuvarapu, Zhen Lin, and Jimeng Sun. 2024. Codrug: Conformal drug property prediction with density estimation under covariate shift. *Advances in Neural Information Processing Systems*, 36. 585 586 587 588

Hongzhan Lin, Ziyang Luo, Wei Gao, Jing Ma, Bo Wang, and Ruichao Yang. 2024. Towards explainable harmful meme detection through multimodal debate between large language models. In *Proceedings of the ACM on Web Conference 2024*, pages 2359–2370. 589 590 591 592 593

Shuai Niu, Jing Ma, Liang Bai, Zhihua Wang, Li Guo, and Xian Yang. 2024. Ehr-knowgen: Knowledge- 594 595

<sup>4</sup><https://about.citiprogram.org/>

<sup>5</sup><https://physionet.org/>

596	enhanced multimodal learning for disease diagnosis generation. <i>Information Fusion</i> , 102:102069.	651
597		652
598	AI Open. 2023. Chatgpt (mar 14 version)[large language model].	653
599		654
600	Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. Bleu: a method for automatic evaluation of machine translation. In <i>Proceedings of the 40th annual meeting of the Association for Computational Linguistics</i> , pages 311–318.	655
601		656
602		657
603		658
604		659
605	Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. 2019. Language models are unsupervised multitask learners. <i>OpenAI blog</i> , 1(8):9.	660
606		661
607		662
608		663
609	Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J Liu. 2020. Exploring the limits of transfer learning with a unified text-to-text transformer. <i>The Journal of Machine Learning Research</i> , 21(1):5485–5551.	664
610		665
611		666
612		667
613		668
614		669
615	Wenqi Shi, Ran Xu, Yuchen Zhuang, Yue Yu, Haotian Sun, Hang Wu, Carl Yang, and May D Wang. 2024. Medadapter: Efficient test-time adaptation of large language models towards medical reasoning. <i>arXiv preprint arXiv:2405.03000</i> .	670
616		671
617		672
618		673
619		674
620	Karan Singhal, Shekoofeh Azizi, Tao Tu, S Sara Mahdavi, Jason Wei, Hyung Won Chung, Nathan Scales, Ajay Tanwani, Heather Cole-Lewis, Stephen Pfohl, et al. 2023. Large language models encode clinical knowledge. <i>Nature</i> , 620(7972):172–180.	675
621		676
622		677
623		678
624		679
625	Jiwoong Sohn, Yein Park, Chanwoong Yoon, Sihyeon Park, Hyeon Hwang, Mujeen Sung, Hyunjae Kim, and Jaewoo Kang. 2024. Rationale-guided retrieval augmented generation for medical question answering. <i>arXiv preprint arXiv:2411.00300</i> .	680
626		681
627		682
628		683
629		684
630	Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. 2023. Llama: Open and efficient foundation language models. <i>arXiv preprint arXiv:2302.13971</i> .	685
631		686
632		687
633		688
634		689
635	Maria Tsimpoukelli, Jacob L Menick, Serkan Cabi, SM Eslami, Oriol Vinyals, and Felix Hill. 2021. Multimodal few-shot learning with frozen language models. <i>Advances in Neural Information Processing Systems</i> , 34:200–212.	690
636		691
637		692
638		693
639		694
640	A K A Unnithan, J M Das, and P Mehta. 2023. Hemorrhagic stroke. In <i>StatPearls</i> . StatPearls Publishing.	695
641		696
642	HP Vinutha, B Poornima, and BM Sagar. 2018. Detection of outliers using interquartile range technique from intrusion dataset. In <i>Information and Decision Sciences: Proceedings of the 6th International Conference on FICTA</i> , pages 511–518. Springer.	697
643		698
644		699
645		
646		
647	Zifeng Wang and Jimeng Sun. 2022. Promptehr: Conditional electronic healthcare records generation with prompt learning. In <i>Proceedings of the 2021 Conference on Empirical Methods in Natural Language Process-</i>	
648		
649		
650		
	<i>ing</i> , pages 2873 – 2885. Association for Computational Linguistics.	
	Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, et al. 2022. Chain-of-thought prompting elicits reasoning in large language models. <i>Advances in Neural Information Processing Systems</i> , 35:24824–24837.	
	Xi Yang, Aokun Chen, Nima PourNejatian, Hoo Chang Shin, Kaleb E Smith, Christopher Parisien, Colin Compas, Cheryl Martin, Mona G Flores, Ying Zhang, et al. 2022. Gatortron: a large clinical language model to unlock patient information from unstructured electronic health records. <i>arXiv preprint arXiv:2203.03540</i> .	
	Susan Zhang, Stephen Roller, Naman Goyal, Mikel Artetxe, Moya Chen, Shuohui Chen, Christopher Dewan, Mona Diab, Xian Li, Xi Victoria Lin, et al. 2022. Opt: Open pre-trained transformer language models. <i>arXiv preprint arXiv:2205.01068</i> .	
	Tianyi Zhang, Varsha Kishore, Felix Wu, Kilian Q Weinberger, and Yoav Artzi. 2019. Bertscore: Evaluating text generation with bert. In <i>International Conference on Learning Representations</i> .	
	Xinlu Zhang, Chenxin Tian, Xianjun Yang, Lichang Chen, Zekun Li, and Linda Ruth Petzold. 2023. Alpacare: Instruction-tuned large language models for medical application. <i>arXiv preprint arXiv:2310.14558</i> .	
	<b>A Appendix</b>	
	<b>A.1 Lab Test Anomaly Caption</b>	
	To caption lab test results into textual descriptions, we initially employ the Inter Quartile Range (IQR) anomaly detection method (Vinutha et al., 2018) to identify anomalous lab test features. Subsequently, we craft multiple text templates to caption these anomalies. These templates are delineated in Table 4.	
	<b>A.2 Prompts for Multimodal Rationale Generation Via ChatGPT</b>	
	The overall procedure for ChatGPT generating clinical rationale is illustrated in Figure 3. The specific medical note-based rationale prompt and lab test-based rationale prompt are detailed as follows.	
	<b>Medical note-based rationale prompt</b> for ChatGPT,	
	“Below is an instruction that describes examples of generating the rationale of disease diagnosis; please refer to the examples style to generate the Output from the Input:	
	### Instruction:	
	There are some examples please to refer:	

Condition: If the lab test value is not an abnormal value:
Prompt: {Lab features} is normal all the time.
Condition: If the lab test value is an abnormal value higher than the standard:
Prompt: {Lab features} is higher than normal {number of times} times.
Condition: If the lab test value is an abnormal value lower than the standard:
Prompt: {Lab features} is lower than normal {number of times} times.
Condition: If the lab test value is abnormal, it includes both higher and lower than the standard value:
Prompt: {Lab features} is higher than normal {number of times} times and lower than normal {number of times} times.

Table 4: Lab test anomaly caption template.

700	<i>Example 1, Example 2, Example ...</i>	<i>Please review the patient’s medical notes, laboratory test anomaly results, and existing rationales in the medical record. Construct a concise, one-sentence rationale, limited to max 50 words, that accurately describes a diagnosed condition based on descriptions of laboratory test abnormalities (Start with "Lab test shows..."). Pay close attention to potential inaccuracies in the lab descriptions.</i>	729
701	### Input:		730
702	### Medical note: [M]		731
703	### Diagnosed diseases: [D]		732
704	<i>Please review the patient’s medical records. Adhere to the provided format to craft a succinct 100-word rationale for diagnosing these conditions (Start with "Based on the medical notes..."). If the diagnosis indicates "no disease was diagnosed," the rationale must state "no disease was diagnosed." Otherwise, provide a comprehensive rationale for the diagnosis.</i>		733
705			734
706		### Response:	735
707		### Output:	736
708		### Lab test-based rationale: [R <sup>t</sup> ]	737
709			738
710			739
711		<b>A.3 Baseline Details</b>	740
712	### Response:	• <b>Flan-T5:</b> Flan-T5 is introduced in the scaling instruction-fine-tuning method for language models (Chung et al., 2024). It is trained on comprehensive datasets designed for tasks like summarization, question answering, and reasoning, enhancing its chain-of-thought capabilities.	741
713	### Output:		742
714	### Medical note-based Rationale: [R <sup>n</sup> ]		743
715	<b>Lab test-based rationale prompt</b> for ChatGPT, we denote the lab test anomalies as $T^*$ :		744
716			745
717	<i>“Below is an instruction that describes examples of generating the rationale of disease diagnosis, please refer to the examples style to generate the Output from the Input:</i>		746
718			747
719		• <b>PROMPTEHR:</b> PROMPTEHR (Wang and Sun, 2022) innovates generative modelling for EHRs through conditional prompt learning; in this experiment, we focus on applying it, particularly on disease diagnosis.	748
720	### Instruction:		749
721	<i>There are some examples please to refer:</i>		750
722	<i>Example 1, Example 2, Example ...</i>		751
723	### Input:		752
724	### Medical note: [M]		753
725	### Descriptions of lab test abnormalities: [T <sup>*</sup> ]		754
726	### Diagnosed diseases: [D]		755
727	### Medical note-based rationale: [R <sup>n</sup> ]		756
728		• <b>LLaMA:</b> The LLaMA-7B model (Touvron et al., 2023), a prominent large language model, employs Reinforcement Learning with Human Feedback (RLHF) and instructional tuning, showcasing its adaptability across diverse NLP tasks. This study applied both zero-shot and fine tuning for disease diagnosis.	757
		• <b>FROZEN:</b> The FROZEN framework (Tsim-poukelli et al., 2021) stands out in multimodal vision-language modeling for few-shot learning. Here, it’s tailored to disease diagnosis,	758
			759
			760
			761
			762
			763

analyzing both lab test results and medical notes.

- **EHR-KnowGen**: As a leading model in EHR multimodal learning, EHR-KnowGen (Niu et al., 2024) specializes in generating disease diagnoses. This study excludes external knowledge to maintain a balanced evaluation.
- **Clinical CoT**: Clinical CoT (Kwon et al., 2024) integrates clinical reasoning into a diagnostic framework for EHRs using prompt-based learning methods distilled from GPT. To ensure a fair comparison, we incorporate the same time series encoder (TSE) as used in our model for multimodal processing.
- **ChatGPT**: ChatGPT (Open, 2023) is a state-of-the-art LLM optimized for conversational applications, such as dialogue, summarization, and text completion.

#### A.4 Implementation Details

For our experiments, we utilized version 2.0.1 of the PyTorch framework, running on a CUDA 11.7 setup. The training processes were conducted using the DeepSpeed<sup>6</sup> framework. We opted for the AdamW optimizer, starting with a learning rate of  $1e^{-5}$  and incorporating a weight decay of 0.05. We implemented a warm-up phase that spanned 10% of the training period. The experimental setup included two NVIDIA A100 GPUs, each with 80 GB of memory. To process time series data consistently, we padded all lab test results to a standard length of 1,000 time steps, dividing the data into 125 patches, where each patch included 8-time steps.

#### A.5 Rationale Evaluation Metrics

We defined the rationale evaluation metrics for the LLM and human evaluation as follows: 1). *Correctness*: how medically accurate the rationale supports the diagnosis results. 2). *Readability*: the extent to which a clinical rationale adheres to proper grammar and structural rules. 3). *Soundness*: the logical coherence and insight provided by the clinical rationale. 4). *Consistency*: the degree of alignment between the clinical rationale derived from medical notes and lab test results. 5). *Persuasiveness*: the effectiveness of the clinical rationale in convincing the reader of its validity.

Evaluation scores based on Likert scale:

<sup>6</sup><https://github.com/microsoft/DeepSpeed>

1. Strongly disagree 810
2. Disagree 811
3. Neither agree nor disagree 812
4. Agree 813
5. Strongly agree 814

#### A.6 Rationale Evaluation with BLUE and BERTScore.

In addition to the criteria defined for evaluating rationale performance, Table 5 presents the performance of our model, ClinRaGen, alongside various baselines, using both BLEU (Papineni et al., 2002) and BERTScore (Zhang et al., 2019) on the MIMIC-III and MIMIC-IV datasets. The results show that ClinRaGen outperformed all other models in both metrics across the datasets. The latest open-source LLM, LLaMA3, ranked second, while Mistral exhibited the poorest performance. These results are consistent with those from LLM and human evaluations.

Models	MIMIC-III		MIMIC-IV	
	BLEU	BERTScore	BLEU	BERTScore
Mistral	0.0163	0.7348	0.0532	0.7625
LLaMA2	0.1441	0.8714	0.2357	0.8808
LLaMA3	0.1641	0.8804	0.2568	0.8919
ClinRaGen	<b>0.2689</b>	<b>0.8972</b>	<b>0.2963</b>	<b>0.9044</b>

Table 5: Rationale evaluation with BLEU and BERTScore.

#### A.7 Discussion on SLMs Selection

In this section, we discuss our choice of Flan-T5 as the base SLM for our research, focusing on the following aspects: 1). **Required CoT ability**: Flan-T5 (Chung et al., 2024) has been extensively instruction-tuned on numerous datasets and hundreds of tasks, endowing it with strong zero-shot, few-shot, and Chain-of-Thought (CoT) abilities that outperform the original T5 (Raffel et al., 2020). In contrast, other SLMs, such as OPT (Zhang et al., 2022) and GPT-2 (Radford et al., 2019), lack these robust CoT capabilities, which is crucial as a initialization ability for further clinical reasoning distillation. 2). **Maximizing SLM potential for practical usage**: Although other instruction-finetuned SLMs (e.g., Flan-PaLM) exist, they have substantially larger parameter counts (ranging from 8B to 540B), which is not practical in real world clinical applications and not our target SLMs to investigate.

848 We selected Flan-T5-Small (80M) and Flan-T5-  
849 Large (780M) as our base models to maximize the  
850 potential of SLMs for accurate disease diagnosis  
851 and LLM-comparable multimodal reasoning, while  
852 maintaining cost-effectiveness in practical applica-  
853 tions. Although we currently use Flan-T5, future  
854 work will explore a broader range of SLM architec-  
855 tures to further enhance accuracy.

## 856 **A.8 Discussion on Teacher LLMs Selection**

857 In this section, we discuss our choice of ChatGPT  
858 (GPT-3.5-turbo) as the teacher LLM for our re-  
859 search, focusing on the following aspects: 1) **High**  
860 **quality clinical rationales** : Although ChatGPT  
861 is known for its strong language modeling capa-  
862 bilities, its generated rationales may still contain  
863 noise and bias—issues that are critical in precision  
864 medicine. To address this, we incorporate external  
865 medical domain knowledge and introduce a novel  
866 knowledge-augmented attention mechanism during  
867 multimodal clinical rationale generation. Our ex-  
868 tensive experiments (Section 4.3) show that ClinRa-  
869 Gen effectively mitigates incomplete or incorrect  
870 diagnoses and rationales distilled from the teacher  
871 LLM, thereby reducing the impact of bias in LLM-  
872 generated outputs for SLMs distillation. 2) **Test**  
873 **set leakage** : The PhysioNet Credentialed Data  
874 Use Agreement prohibits the use of MIMIC-series  
875 data in public LLMs’ training and applications<sup>7</sup>,  
876 ensuring that test set leakage is not an issue with  
877 ChatGPT. Furthermore, the poor performance of  
878 ChatGPT under zero-shot prompting (as shown  
879 in Table 1) indicates that MIMIC-III and MIMIC-  
880 IV data were not used in its training. 3) **More**  
881 **powerful LLM for evaluation**: While we did not  
882 choose the most powerful LLM as our teacher, our  
883 current teacher LLM sufficiently enhances SLM  
884 capabilities in disease diagnosis and clinical ratio-  
885 nale generation. Our evaluations (Sections 4.2 and  
886 4.3), supported by quantitative metrics, compar-  
887 isons with a superior LLM, and human assessments,  
888 confirm the effectiveness of using ChatGPT as the  
889 teacher LLM through our multimodal rationale dis-  
890 tillation paradigm and knowledge-augmented at-  
891 tention mechanism to improve SLMs’ accuracy in  
892 disease diagnosis and modality-consistent rationale  
893 generation. In future work, more powerful LLMs  
894 can be seamlessly integrated with our method to  
895 further enhance evaluation accuracy and robust-  
896 ness.

---

<sup>7</sup><https://physionet.org/news/post/gpt-responsible-use>