

# Early Risk Prediction with Temporally and Contextually Grounded Clinical Language Processing

Rochana Chaturvedi<sup>\*1</sup>, Yue Zhou<sup>2</sup>, Andrew Boyd<sup>2</sup>, Brian T. Layden<sup>2</sup>, Mudassir Rashid<sup>3</sup>,  
Lu Cheng<sup>3</sup>, Ali Cinar<sup>3</sup>, Barbara Di Eugenio<sup>2</sup>

<sup>1</sup>Mathematics and Computer Science Division, Argonne National Laboratory, Lemont, IL, USA

<sup>2</sup>University of Illinois Chicago, Chicago, IL, USA

<sup>3</sup>Illinois Institute of Technology, Chicago, IL, USA

Correspondence: [rchaturvedi@anl.gov](mailto:rchaturvedi@anl.gov)

## Abstract

Clinical notes in Electronic Health Records (EHRs) capture rich temporal information on events, clinician reasoning, and lifestyle factors often missing from structured data. Leveraging them for predictive modeling can be impactful for timely identification of chronic diseases. However, they present core natural language processing (NLP) challenges: long text, irregular event distribution, complex temporal dependencies, privacy constraints, and resource limitations. We present two complementary methods for temporally and contextually grounded risk prediction from longitudinal notes. First, we introduce HiT-GNN, a hierarchical temporal graph neural network that integrates intra-note temporal event structures, inter-visit dynamics, and medical knowledge to model patient trajectories with fine-grained temporal granularity. Second, we propose REVEAL, a lightweight, test-time framework that distills the reasoning of large language models into smaller verifier models. Applied to opportunistic screening for Type 2 Diabetes (T2D) using temporally realistic cohorts curated from private and public hospital corpora, HiT-GNN achieves the highest predictive accuracy—especially for near-term risk—while preserving privacy and limiting reliance on large proprietary models. REVEAL enhances sensitivity to true T2D cases and retains explanatory reasoning. Our ablations confirm the value of temporal structure and knowledge augmentation, and fairness analysis shows HiT-GNN performs more equitably across subgroups.

## 1 Introduction

Modeling disease progression represents one of the most critical challenges in modern healthcare, with particularly high stakes for chronic conditions like type 2 diabetes (T2D), which affects half a billion people worldwide and continues rising (IDF,

2025). The rich information in free-text clinical notes available in EHRs provides an opportunity for NLP predictive modelling. These notes present two natural perspectives for understanding disease progression and suggest distinct yet complementary computational approaches. First, the **temporal progression** of symptoms, diagnoses, treatments, and interventions that are natural for understanding the underlying disease pathology requires structured modeling. Second, capturing the **semantic richness** of clinical notes calls for contextualized language understanding. These two perspectives are particularly relevant for opportunistic screening, where patients with low socioeconomic status often miss routine care and timely screening (Danielson et al., 2023). In such settings, clinicians must act on all historical information available at a given visit to estimate the disease risk. Formally, given patient documents  $d_1, d_2, \dots, d_n$  observed at non-decreasing sequence of patient visits  $1, 2, \dots, n$ , we want to estimate the risk at next visit  $(n + 1)$ .

To address these complementary modeling needs, this work proposes two novel methods that integrate domain-specific modeling with targeted representation learning techniques from a sequence of clinical notes, optimized for low-resource settings. To address structured temporal modeling, we propose **Hierarchical Temporal Graph Neural Network (HiT-GNN)**, a dynamic model that captures patient state evolution with fine-grained temporal granularity and knowledge-enhanced structure. This fine-grained temporal knowledge is crucial for understanding disease pathways. For example, elevated glucose **after** corticosteroid use likely indicates drug-induced hyperglycemia, not T2D, while the same glucose levels **before** any steroid therapy may signal early T2D onset. This distinction motivates our HiT-GNN framework, which extracts event-temporal graphs from individual clinical notes using state-of-the-art methods, and augments them with semantic information from a clin-

<sup>\*</sup>This work was conducted while the author was affiliated with the University of Illinois Chicago.

ical knowledge base. These multi-layered event-temporal graphs are modeled with dynamic graph neural networks (GNNs) to capture patient timeline within and across visits. While GNNs are effective for relational reasoning, their performance depends on rich contextual embeddings; thus, we incorporate clinically pretrained language model (CPLM) embeddings and knowledge-graph embeddings to enhance input representations with rich semantics. This addresses the long-recognized need to incorporate causal and temporal patterns into diagnostic reasoning (Patil et al., 1981), moving beyond modeling patient trajectories as sequences of visits.

Complementing this structured approach, current large language models (LLMs) offer strong localized and context-aware reasoning capabilities that can effectively model the underlying text semantics. However, they often struggle with long-range context integration and structured reasoning—limitations particularly pronounced in the clinical domain, where patient histories span multiple long notes. Additionally, real-world healthcare deployments face privacy and resource constraints, limiting the use of proprietary or large-scale models. Our second method, **Reasoning with Verifier-Aided Labeling (REVEAL)**, is a lightweight, scalable inference-time architecture that preserves LLM-style reasoning within resource-constrained clinical settings. It distills reasoning from a large LLM into a smaller model, scaling performance with interpretability without the full computational cost of training or deploying massive models.

Together, these complementary modeling paradigms—graph-based temporal reasoning and interpretable LLM inference—provide a holistic risk prediction framework grounded in structure and semantics. We evaluate them in the context of opportunistic screening of T2D as a representative use case, using a real private hospital (PH) corpus and a corpus curated from MIMIC-IV (Johnson et al., 2023). Additionally, due to the prevalent demographic biases (Meng et al., 2022; Zhou et al., 2025a; Hall et al., 2015), we conduct a fairness analysis to better understand model behavior across demographic groups.

**Contributions.** (1) We present two complementary representation learning frameworks from longitudinal clinical notes, tailored for low-resource, privacy-sensitive settings: (i) **HIT-GNN**: the first application of temporal relation extraction systems for clinical risk prediction that integrates intra-document temporal relations, inter-visit dynamics,

and medical knowledge, enabling reasoning across both local event structures and longitudinal patient trajectories. (ii) **REVEAL**: an inference-time scaling framework where a smaller LLM validates predictions from a larger frozen LLM, inheriting interpretability and improving accuracy without full retraining.

(2) We demonstrate the translational value of temporally enriched representations in a real clinical application—opportunistic screening for T2D, especially in immediate-risk horizons, where intervention is most impactful.

(3) We curate rigorous datasets from public (MIMIC-IV) and private (PH corpus) sources, exclude post-diagnosis inputs to prevent label leakage and ensure balanced cohorts for fair evaluation.

(4) Extensive ablations of graph architectures and embedding choices confirm value of fine-grained temporal structure and knowledge augmentation.

(5) Our fairness analysis highlights how clinical data can encode bias even without explicitly modeling sensitive attributes, and discusses implications for equitable model design.

These contributions enable a clinically grounded framework that integrates events, timelines, and semantic context to enable low-cost, privacy-conscious decision support.

## 2 Related Literature

**Computational Approaches for Disease Progression Prediction** Despite significant advances in predictive modeling using electronic health records (EHR), many studies continue to rely predominantly on structured EHR data (Lipton et al., 2016; Singhal et al., 2023b; Jiang et al., 2024), which is often noisy (Hersh et al., 2013) and incomplete (Capurro et al., 2014). Among efforts that leverage free-text, recent works either use only the last clinical note (Xu et al., 2023; Nguyen et al., 2024), relying on structured data to capture temporal trends; or use a coarser form of temporal modeling where each note is embedded using a language model (Huang et al., 2020), a bag-of-concepts (Chaturvedi et al., 2023), or topics (Ghassemi et al., 2015), and a BiLSTM models the representations from multiple notes across patient visits. While these studies explore various fusion strategies, a key gap remains: no prior work models temporal relations both within a single clinical note and across multiple notes/visits to forecast long-term disease risk. Our work directly addresses this gap by: (1) lever-

aging the sequence of clinical notes for each patient, and (2) introducing multi-layered modeling to leverage event-temporal information extracted from each note and across the notes from multiple visits to track evolving health trajectories.

**LLM in Healthcare** The integration of Large Language Models (LLMs) into healthcare is a rapidly evolving field (Zhou et al., 2024a; He et al., 2024). Key developments include impressive performance on medical question answering (QA) benchmarks (Singhal et al., 2023a). However, outside of controlled QA settings, LLMs still struggle on tasks requiring complex reasoning for disease diagnosis in real clinical settings (Hager et al., 2024; Wang et al., 2024b), and often exhibit performance disparities across demographic groups (Zhou et al., 2025b). Recent research has demonstrated that scaling test-time compute through verifier-guided search can significantly improve LLMs’ reasoning capabilities and often outperforms simply scaling model size (Snell et al., 2024). However, scaling test-time compute of LLMs in low-resource real-world healthcare tasks is still underexplored. Our work contributes to this literature by exploring verifier-guided reasoning approaches tailored to low-resource clinical prediction task.

### 3 Data Curation

Our first dataset, the PH corpus, is curated from private data comprising adults (age  $\geq 18$  years) collected from the University of Illinois Hospital and Health Sciences System (UI Health) between January 2010 and July 2021. We also curate a second corpus from MIMIC-IV (Medical Information Mart for Intensive Care, version 4) to study the generalizability of our methods. MIMIC-IV is a large, publicly available database comprising de-identified clinical notes from ICU encounters at Beth Israel Deaconess Medical Center in Boston from 2008–2019. We exclude diagnosis of other types of diabetes (e.g., Type-1 diabetes, gestational diabetes), and construct the T2D and the non-diabetic (Nod) cohorts using ICD codes—standardized diagnostic codes used to classify medical conditions.

Cohort refinement steps ensure that post-diagnosis data is excluded to prevent label leakage. While we use the date of the first recorded ICD code as a proxy for the onset of type 2 diabetes (T2D), this signal is known to be noisy (Hersh et al., 2013) and often delayed relative to when the diagnosis is first documented in clinical notes. To

mitigate this, we additionally process the T2D cohort’s notes using a large language model (LLM) to identify the earliest explicit mention of a T2D diagnosis in unstructured text. This type of additional filtering is often overlooked in existing works, yet it is crucial for ensuring realistic model evaluation. For example, Zhang et al. (2024) show that a widely used sepsis prediction model frequently issues risk alerts only after clinicians have documented suspicion, undermining its practical utility. Finally, we construct a demographically matched test set for fair evaluation. This is achieved by estimating T2D propensity scores (Rosenbaum and Rubin, 1985) from demographic variables, including Age, Gender, and Race (and Ethnicity in the case of the PH corpus, where this attribute is combined with Race in MIMIC-IV). This is followed by 1:1 greedy nearest-neighbor matching without replacement to select NoD controls.<sup>1</sup> The final PH corpus comprises 3332 patients (712 in the test set), and the final MIMIC-IV subset contains 5802 patients (291 in the test set). Both datasets

pose unique strengths and limitations—PH corpus contains richer longitudinal histories with diverse note types ideal for chronic disease modeling but it is not public. In contrast, while MIMIC-IV is one of the most widely used, large public datasets, it only focuses on ICU notes, omitting much of the hospitalization timeline, and has narrower temporal windows (84.5% patient records comprise single visits (Cui et al., 2025)). Further, MIMIC-IV notes are de-identified, lacking an essential temporal anchor—DATE—considered protected health information (PHI).<sup>2</sup> Figure 1 shows demographic differences: *females* are the majority in PH but a minority in MIMIC-IV; *Black* is the majority in PH, *White* in MIMIC-IV, with *Asian* as the minority. PH corpus records Hispanics under a separate ethnicity attribute. Therefore, evaluating methods on both datasets is important. However, given MIMIC-IV’s limited temporal continuity, its results would likely represent a lower bound on model performance achievable in richer, real-world settings.

### 4 HiT-GNN: Disease Prediction using Hierarchical Temporal Graphs

This section presents our **Hierarchical Temporal Graph Neural Network (HiT-GNN)** for modeling

<sup>1</sup>The data filtering details are provided in Appendix A.

<sup>2</sup>Health Insurance Portability and Accountability Act (HIPAA), a U.S. law, considers dates as PHI.

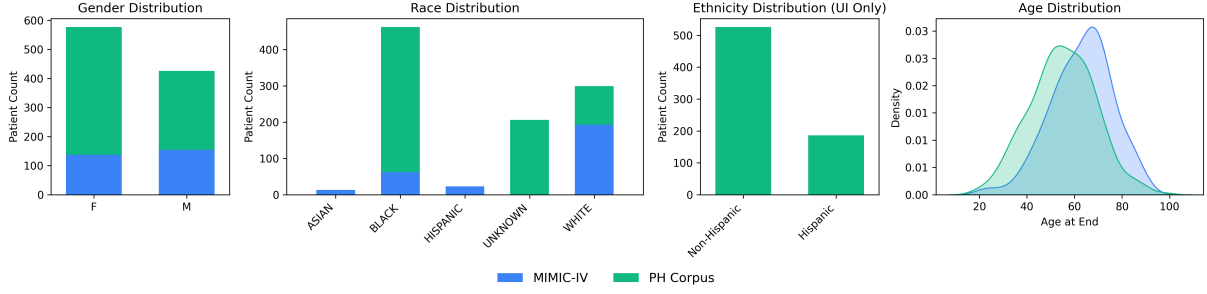


Figure 1: Demographic distribution in PH and MIMIC-IV test sets.

fine-grained temporal information to predict T2D risk. We extract and refine event-temporal graphs from clinical notes (sections 4.1–4.4) and detail our modeling approach in section 4.5.

#### 4.1 Extracting and Aligning Event-Temporal Graphs from Notes

We use state-of-the-art temporal relation extraction models from Chaturvedi et al. (2025) to identify clinical entities (*Problem, Test, Treatment, Clinical Department, Clinical Occurrences, Evidential*), time expressions (*Date, Time, Duration, Frequency*), and their temporal relations (*before, after, overlap*) from each clinical note, forming a temporal graph. We then perform entity normalization to cluster synonymous terms as a single entity and also to link entities to an external knowledge graph.

#### 4.2 Entity Linking/Normalization

**Clinical Entities** If any event is of type *Problem/Treatment/Test/Clinical Department*, we use Metamap (Aronson and Lang, 2010) to map the event text to a unique concept in UMLS Metathesaurus (Bodenreider, 2004).<sup>3</sup> UMLS comprises a vast repository integrating synonymous medical terms to unique concept identifiers (CUI) and also defines various relationships among concepts, including hierarchical (e.g., *is-a*) and associative (e.g., *part-of, related-to*) links. Besides the metathesaurus, UMLS also consists of a Semantic Network. This higher-level abstraction organizes UMLS concepts into broader semantic types, such as Disease or Syndrome, Pharmacologic Substance, Gene, etc. It also defines permissible semantic relationships between semantic types, for example, Pharmacologic Substance *treats* Disease or Syndrome, Finding *associated\_with* Disease or Syndrome. We represent these semantic types as additional nodes

<sup>3</sup>We do not normalize other entity types, as it leads to noise, as per a manual inspection of 150 clinical notes.

in the temporal graph and connect them to their corresponding entity nodes to preserve type hierarchy with *is-a* relations. This results in a heterogeneous graph that combines semantic relations based on medical knowledge with the temporal sequence information determined from the note text.

**Dates** We use Microsoft Recognizers Text (Huang et al., 2017) to normalize dates, and link multiple mentions of same date together.

#### 4.3 Node Representation

To represent a node, we compute its corresponding textual span’s embedding by concatenating the in-context BioMedBERT (Gu et al., 2021) embeddings of the first and last tokens with the span-width embeddings.<sup>4</sup> We represent multiple linked mentions of an entity as a single node in each temporal graph, and initialize it with the mean of the embeddings of all linked mentions. To represent semantic type nodes, we average BioMedBERT embeddings over all tokens in their label (preferred name). We also use knowledge graph embeddings (KG-embeddings) for Concept Unique Identifiers (CUIs) derived from the UMLS knowledge graph for each linked event, as introduced by Maldonado et al. (2019). The final node representations are concatenated KG embeddings and text embeddings (BioMedBERT) over all mentions.

#### 4.4 Reduced Temporal Graph

We apply the Timegraph algorithm (Miller and Schubert, 1990) to obtain a reduced temporal graph. It builds the graph one edge at a time, starting with the most confident edges based on prediction probabilities and dropping inconsistent edges (that lead to a cycle). The final graph contains *before, after, and overlap* edges. We flip the *after* edges to *before* for simplicity. The descriptive statistics of the

<sup>4</sup>We use the fine-tuned model version and width embeddings from Chaturvedi et al. (2025).



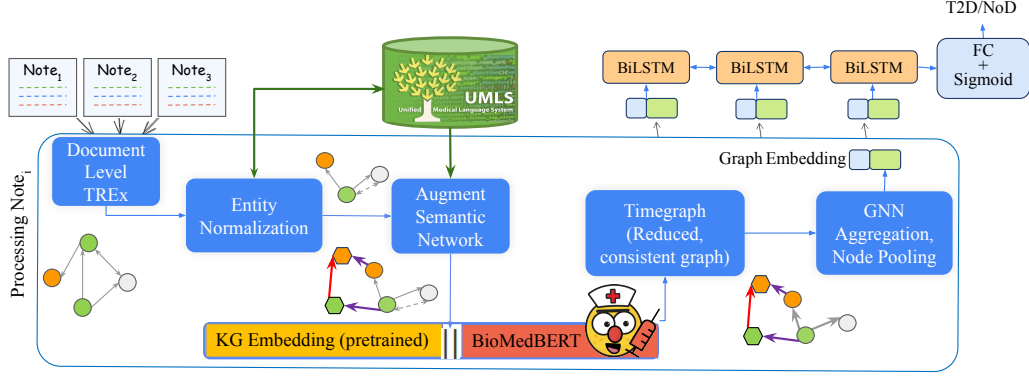


Figure 2: HiT-GNN Architecture: Hierarchical Temporal GNN that models intra- and inter-document temporal dependencies between clinical entities and integrates UMLS knowledge for type 2 diabetes (T2D) risk prediction.

extracted graphs are provided in Table 1.<sup>5</sup>

Parameter	PH-Data	MIMIC-IV
#Patients	3332	5802
Avg. #tokens/doc	1209.2 (637.6)	454.5 (208.3)
Avg. #nodes/doc	114.7 (46.1)	74.8 (27.9)
Avg. #Edges/doc	261.1 (133.4)	158.2 (76.8)
#doc or visits	3.2(1.7)	2.3 (1.4)

Table 1: Data statistics: mean number of tokens per document, nodes, and edges in the extracted temporal graph per document (doc), per patient, and the number of documents per patient. All statistics are averaged over patient IDs, with standard deviations in parentheses.

#### 4.5 Hierarchical Temporal Modeling

We model irregular time-series of temporal graphs across patient visits using HiT-GNN for Type 2 Diabetes prediction. Each visit is a graph  $\mathcal{G}_t = (\mathcal{V}_t, \mathcal{E}_t)$  with temporal (Before/After/Overlap) and semantic (UMLS) edges. HiT-GNN captures intra-graph and inter-graph dependencies via a GNN and a recurrent neural network (RNN).

**Intra-visit Modeling.** Each graph is encoded with a multi-layer GraphSAGE encoder (with mean aggregation, residual connections, and layer normalization). Given initial features  $\mathbf{h}_v^{(0)}$  of node  $v$ ,  $\mathcal{N}(v)$  the neighbors of  $v$ , and learnable parameters  $\mathbf{W}_1^{(k)}, \mathbf{W}_2^{(k)}$ , node updates at layer  $k$  are:

$$\mathbf{h}_v^{(k)} = \text{LayerNorm}\left(\mathbf{W}_1^{(k)}\mathbf{h}_v^{(k-1)} + \mathbf{W}_2^{(k)}\frac{1}{|\mathcal{N}(v)|} \sum_{u \in \mathcal{N}(v)} \mathbf{h}_u^{(k-1)}\right) + \mathbf{h}_v^{(k-1)}$$

<sup>5</sup>The mean number of tokens to get estimated document length is counted based on LLaMA3.1-8b tokenizer, as we also include comparisons with this model.

where  $\mathbf{h}_v^{(k)}$  is the representation of node  $v$  at layer  $k$ . Updated node embeddings are mean-pooled to obtain the graph representation for a visit  $\mathbf{g}_t$ .

**Inter-visit Modeling:** temporal dependencies across visits are modeled as:

$$\mathbf{z}_t = \text{BiLSTM}_\phi(\mathbf{g}_1, \dots, \mathbf{g}_t), \quad \mathbf{z}_T = [\vec{\mathbf{h}}_T; \overleftarrow{\mathbf{h}}_1]$$

Where  $\mathbf{g}_1, \dots, \mathbf{g}_t$  is the sequence of graph encodings over visits. The final forward and backward hidden states from bidirectional LSTM (*BiLSTM*) are concatenated to represent a patient’s trajectory. Finally, a fully connected network with sigmoid ( $\sigma(\cdot)$ ) gives binary outcome ( $\hat{y}$ )—T2D or NoD (No Diabetes):

$$\hat{y} = \sigma(\mathbf{W} \cdot \mathbf{z}_T + \mathbf{b})$$

Figure 2 shows the model architecture, integrating UMLS knowledge graph, document-level temporal information, and cross-visit progression.

#### 5 Reasoning with Verifier-Aided Labeling (REVEAL)

Healthcare applications face significant resource constraints, which limit the deployment of LLMs; yet, maintaining reasoning capabilities is crucial for effective clinical decision-making. Recent advances in test-time compute scaling offer a promising solution: using smaller models enhanced by verification mechanisms, rather than relying solely on larger, computationally expensive ones. To this end, the REVEAL framework combines a reasoner LLM with a smaller verifier LLM, in three steps: (1) a reasoner LLM generates  $N$  reasoning paths with predictions; (2) a fine-tuned verifier evaluates the credibility of these paths and assigns a score from 0–1, and (3) predictions across different reasoning paths are aggregated using verifier scores.

Concretely, at *training time*, the reasoner model generates five stochastic outputs per example, each including a prediction ('true'/'false') and an explanation. A small LLM verifier, fine-tuned with LoRA (Hu et al., 2022), classifies each output as 'correct' or 'incorrect' using the prompt in Figure 8. The normalized probabilities of the output tokens serve as the verifier's confidence score. At *inference time*, the reasoner generates  $N$  ( $N = 10$ ) diverse predictions. Each prediction contains a risk outcome ('true'/'false') and a corresponding explanation. The verifier scores them, and the final outcome is selected by majority vote among the top- $k$  highest-confidence predictions.

## 6 Experiments

This section presents a comprehensive empirical evaluation of HIT-GNN and REVEAL.<sup>6</sup>

### 6.1 Baseline Methods

**Zero-shot Prompting:** The model directly predicts T2D risk from clinical notes using deterministic decoding (temperature=0.0). We extract normalized token probabilities for 'true'/'false' predictions to compute AUC scores (see prompt in Figure 6). We use LLaMA3.1-8b and LLaMA3.2-1, and also compare the performance on MIMIC-IV with prevailing large-scale models, DeepSeek-V3 (671B parameter Mixture-of-Experts model) and GPT-4o.<sup>7</sup>

**Self-Consistency (LLaMA3.1-8B-SC):** We generate  $N$  ( $N = 10$ ) diverse predictions per patient using stochastic sampling (temperature=1.0) and take the majority vote for the final classification.

**Supervised Fine-tuning (SFT):** We fine-tune a 1B LLM (given the low-resource setting) using LoRA for direct binary classification. The model (LLaMA3.2-1B-ft) is fine-tuned to output only 'true' or 'false' labels without explanations due to lack of reasoning data (see prompt in Figure 7).

**CLSTM** We also use concept-based CNN-LSTM baseline (Chaturvedi et al., 2023). Here we model extracted entities (nodes of temporal graphs) with a CNN and max pooling, then process each document representation with an LSTM.

<sup>6</sup>See Appendix B for prompts and experimental details. We will publicly release our code and data curation scripts upon publication.

<sup>7</sup>Not evaluated on PH corpus due to privacy constraints.

### 6.2 Evaluation Metrics

We report precision and recall for both groups, the macro-average  $F_1$ , and ROC-AUC scores.<sup>8</sup> For fairness evaluation, we use the following metrics:

**Demographic Parity Difference (DPD)** : difference in positive prediction rates between a target group ( $Z = z$ ) and others ( $Z \neq z$ ).

$$DPD = P(\hat{Y} = 1 \mid Z = z) - P(\hat{Y} = 1 \mid Z \neq z)$$

Where  $\hat{Y}$  is the predicted label (T2D = 1, NoD = 0), and  $Z$  is a sensitive attribute (e.g., gender, race).

**Equal Opportunity Difference (EOD)** : difference in true positive rates ( $TPR$ ) between the target group and others:  $TPR_{Z=z} - TPR_{Z \neq z}$

Where  $TPR_{Z=z} = P(\hat{Y} = 1 \mid Y = 1, Z = z)$ .

## 7 Results and Discussion

Table 2 presents the main experimental results.<sup>9</sup> The first panel comprises LLM-based models. Regarding zero-shot performance, the LLaMA3.2-1B model fails to identify patients at risk of diabetes, classifying almost all patients into the 'no-risk' (NoD) cohort. The performance improves as the model capacity increases to 8B variant. While self-consistency (SC) has shown promising performance in other tasks and domains, LLaMA3.1-8B-SC performs poorly on both datasets. The LLaMA3.2-1B-ft (fine-tuned) model demonstrates strong performance, achieving a ROC-AUC of 65.06% on the PH corpus subset and 64.41%—the highest among all models—on MIMIC-IV. While this fine-tuned LLM excels as a classifier, it cannot provide reasoning or explanations for its outputs.<sup>10</sup> Our REVEAL model does not significantly outperform the LLaMA3.2-1B-ft model. However, it preserves the reasoning output and improves upon the zero-shot 8B variant in terms of  $F_1$  and T2D group recall, highlighting its sensitivity in capturing true positives for this subgroup. For larger models, including DeepSeek-V3 with impressive clinical decision-making capabilities (Sandmann et al.,

<sup>8</sup>ROC-AUC scores are computed from prediction probabilities. For LLMs, this is the probability of generating 'true' at the predicted token position. For LLM-SC, it is the 'true' token probability, averaged across all runs.

<sup>9</sup>We also perform a pairwise comparison of the top four models using bootstrap sampling to estimate whether the differences are statistically significant, based on non-parametric bootstrap resampling test (see Appendix Table C).

<sup>10</sup>Obtaining gold-standard reasoning data at scale for fine-tuning is costly and labor-intensive, especially from already overburdened clinical experts, limiting the feasibility of fine-tuning with reasoning.

Model	PH-Data						Mimic-IV					
	AUC	F <sub>1</sub>	T2D		NoD		AUC	F <sub>1</sub>	T2D		NoD	
			P	R	P	R			P	R	P	R
LLaMA3.2-1B (0-shot)	46.80	35.46	39.29	3.09	49.56	<b>95.22</b>	40.00	33.10	0	0	49.48	1.00
LLaMA3.1-8B (0-shot)	57.73	56.98	57.49	53.93	56.61	60.11	53.71	51.03	61.54	27.21	52.65	82.64
DeepSeek-V3 (0-shot)	-	-	-	-	-	-	55.02	47.85	69.23	18.37	52.38	91.67
GPT-4o (0-shot)	-	-	-	-	-	-	57.75	51.81	<b>75.56</b>	23.13	54.07	92.36
LLaMA3.1-8B-SC	55.62	55.49	56.29	50.28	55.08	60.96	52.12	43.43	66.67	12.24	51.14	<b>93.75</b>
LLaMA3.2-1B-ft	65.06	61.62	59.53	<b>78.93</b>	<b>68.75</b>	46.35	<b>64.41</b>	59.17	63.89	46.94	57.38	72.92
REVEAL	57.98	65.24	66.77	60.96	64.08	69.66	53.13	51.60	55.21	36.05	51.79	70.14
CLSTM	68.24	65.62	63.19	76.69	70.36	55.34	63.08	59.62	63.96	48.30	57.78	72.22
HiT-GNN	<b>72.24</b>	<b>67.28</b>	<b>71.48</b>	58.43	64.85	76.69	62.97	<b>62.88</b>	63.27	<b>63.27</b>	<b>62.50</b>	62.50
#patients	712		356		356		291		147		144	

Table 2: Performance of HiT-GNN and REVEAL against different baselines.

2025), the gains are less pronounced: it achieves a similar AUC to the much smaller LLaMA3.1-8B, but with lower  $F_1$  and T2D recall. GPT-4o attains a higher AUC but exhibits worse T2D recall than both LLaMA3.1-8B and REVEAL, and an overall weaker performance than LLaMA3.2-1B-ft.

The second panel summarizes results with structured approaches. The CLSTM model provides a strong baseline, significantly outperforming LLMs in AUC and T2D recall on the PH corpus and performing comparably on MIMIC-IV. HiT-GNN performs even better, attaining the highest AUC on the PH corpus and the highest T2D recall on MIMIC-IV, showing superior effectiveness in identifying at-risk patient. These improvements over other top-performing models (REVEAL, LLaMA3.2-1B-ft, and CLSTM) are statistically significant. In contrast, no other model shows significant improvements over HiT-GNN across key metrics (AUC,  $F_1$ , and T2D recall) on either corpus. For example, while LLaMA3.2-1B-ft attains the highest AUC on MIMIC-IV, its advantage over CLSTM and HiT-GNN is not statistically significant. HiT-GNN also outperforms both the larger models (DeepSeek-V3 and GPT-4o) across key metrics. We next present additional analyses and ablation studies.

### 7.1 Prediction Horizons

We compare performance across prediction horizons by dividing the time from last pre-diagnosis visit to diagnosis into 3-month windows, combining each T2D subgroup with all NoD controls to compute AUC.<sup>11</sup> The final window aggregates values above the 95th percentile to avoid sparsity. Results (Figure 3) for the three best-performing models (HiT-GNN, REVEAL, and LLaMA3.2-1B-ft) show HiT-GNN performs more robustly across win-

dows and outperforms LLMs in near-term prediction ( $\leq 1.5$  years). This performance aligns with EHR realities where longitudinal completeness is rare and supports timely preventive care.

### 7.2 HiT-GNN Ablations

**Node Embeddings:** Figure 4a shows HiT-GNN AUC with different node embeddings.<sup>12</sup> Initializing the nodes with in-context BioMedBERT embeddings (‘Text-only’) outperforms CUI embeddings (‘KG-only’), combining both further improves performance across datasets and metrics.

**Subgraph Variations:** Figure 4b shows HiT-GNN AUC across input subgraphs.<sup>13</sup> The first (‘Temporal’) captures the extracted temporal graphs, the second (‘KG’) adds semantic nodes and edges from UMLS and removes the temporal (*before*, *overlap*) edges, and the last combines both. While KG-graph is more important for the PH corpus; for MIMIC-IV, both are significant.

### 7.3 Fairness Analysis

Table 4 presents results on fairness metrics. In the PH corpus, all models exhibit negative demographic parity (DPD) for *Hispanics*; this is more pronounced with HiT-GNN, which also yields a slight negative equal opportunity difference (EOD) ( $-0.08$ ). Across gender, HiT-GNN shows a low DPD bias against *Females* and no opportunity bias, while LLM models show moderate-to-high biases against *Males*. Across racial groups, all models are negatively biased against the minorities (*Unknown* (NI) in the PH corpus and *Asians* in MIMIC-IV) and slightly favor the majority *Black* subgroup in the PH corpus. LLaMA3.2-1B-ft strongly favors this group in MIMIC-IV also. LLMs are somewhat

<sup>11</sup>Appendix D Figure 9 shows trends for T2D group recall.

<sup>12</sup>Appendix D, Figure 10 shows T2D recall.

<sup>13</sup>Appendix D, Figure 11 shows T2D group recall trends.

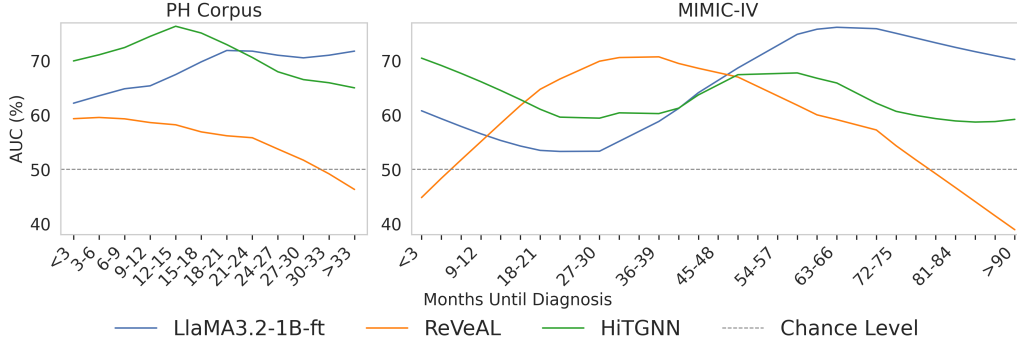


Figure 3: AUC as a function of the prediction horizon, evaluated over consecutive 3-month windows.

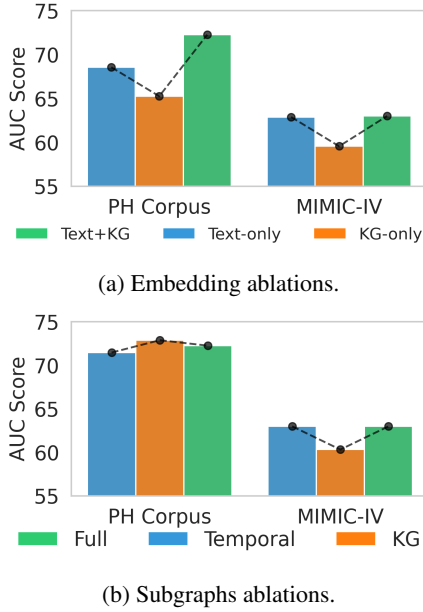


Figure 4: HiT-GNN performance ablations.

negative towards the *White* minority in the PH corpus, while all models are slightly negative against this group in MIMIC-IV, where it has a majority. HiT-GNN shows a high positive bias for *Hispanics* in MIMIC-IV, and LLMs show low-to-moderate bias. Overall, HiT-GNN is relatively fairer.

#### 7.4 Computational Efficiency

Table 3 summarizes the resource usage across models (excluding GPT-4o and DeepSeek-V3, which are accessed via API: GPT-4o via the OpenAI API, and DeepSeek-V3 via the TogetherAI API). A single A100 GPU with 80 Gb RAM is used for all experiments (accelerated inference is used for LLaMA experiments).<sup>14</sup> Training the HiT-GNN

<sup>14</sup>Because HiTGNN relies on MetaMap’s single-threaded API for concept linking, preprocessing is not parallelized. This one-time step—temporal graph extraction, UMLS linking, TimeGraph, and embeddings lookup—takes 42 seconds per PH note and 15.6 seconds per MIMIC-IV note.

model takes approximately 1.5 minutes, and inference per patient takes .007 seconds. For CLSTM, training the model takes 12 minutes, and inference on one patient takes .02 seconds. For inference with the LLaMA model using accelerated inference on a single A100 GPU, the average time required with explanations is approximately 6 seconds. Fine-tuning the model takes approximately 18 hours, and inference with this version takes 0.2 seconds. REVEAL verifier training takes 30 hours (excluding training data preparation, which takes approximately 30 seconds per patient for obtaining explanations from LLaMA3.1-8B along 5 reasoning paths). Inference along 10 reasoning paths from this model takes 62 seconds (obtaining and verifying reasoning along 10 paths). The full LLaMA3.2-1B model requires 2.4 GB of memory, while LLaMA3.1-8B requires 15 GB. The LORA adaptors require an additional 20 MB. In contrast, CLSTM takes 54 MB and HiT-GNN is even lighter at 44 Mb, demonstrating impressive scalability in terms of both speed and memory usage.<sup>15</sup>

## 8 Conclusion

This paper presents two complementary approaches for clinical risk prediction from longitudinal notes: HiT-GNN, a temporally grounded dynamic GNN augmented with clinical knowledge graphs, and REVEAL, a test-time LLM scaling framework for interpretable predictions. HiT-GNN achieves superior predictive performance and robustness, especially on the immediate-risk hori-

<sup>15</sup>Both the CLSTM and HiT-GNN models utilize the state-of-the-art, open-source end-to-end temporal relation extraction models SPANTREX and GRAPHTREX from Chaturvedi et al. (2025). Our reported results utilize SPANTREX, which requires 0.44 GB of disk space. We also obtain similar performance using GRAPHTREX, which outperforms SPANTREX on the temporal relation extraction benchmarks but not on our task, and is larger at 4.35 GB.



Model	Training Time (minutes)	Inference Time (per patient) (seconds)	Memory Usage	Input
HiT-GNN	1.5	0.007	44 MB	5 notes
CLSTM	12	0.02	54 MB	5 notes
LLaMA3.1-1B (with explanations)	pre-trained	6	2.4 GB	2 notes
LLaMA3.1-8B (with explanations)	pre-trained	6	15 GB	2 notes
LLaMA3.1-1B-ft (no explanations)	18 hours	0.2	2.42 GB	2 notes
REVEAL (with explanations)	30 hours	62	17.42 GB	2 notes

Table 3: Comparison of computational efficiency across models for PH-MIMIC-IV corpora. The temporal graph extraction and UMLS linking are one-time preprocessing costs incurred for HiT-GNN and CLSTM models. It takes approximately 48 seconds per note for the PH corpus and 15.6 seconds per note for MIMIC-IV. The table excludes the LLM pretraining time and one-time data preparation costs for fine-tuning REVEAL model; the time shown is using 5 reasoning paths for training, and 10 for inference. LLMs perform best with the last two notes, while HiT-GNN and CLSTM perform best with the last five.

		PH corpus				MIMIC-IV			
		G	GNN	RVL	ft	G	GNN	RVL	ft
Ethnicity	DPD	H	-0.12	-0.05	-0.02				
		N	0.12	0.05	0.02				
	EOD	H	-0.08	0	0.01				
		N	0.08	0	-0.01				
Gender	DPD	F	-0.04	0.03	0.03	F	-0.02	0.04	0.02
		M	0.04	-0.03	-0.03	M	0.02	-0.04	-0.02
	EOD	F	0	0.06	0.06	F	0	-0.10	0.19
		M	0	-0.06	-0.06	M	0	0.10	-0.19
Race	DPD	B	0.09	0.10	0.04	B	0.03	0.03	0.12
		W	0	-0.09	-0.03	W	-0.05	-0.03	-0.04
		NI	-0.11	-0.06	-0.03	H	0.21	0.07	0.02
					A	-0.21	-0.10	-0.31	
	EOD	B	0.08	0.08	0.06	B	-0.05	-0.06	0.16
		W	0	-0.09	-0.03	W	-0.03	-0.01	-0.05
		NI	-0.11	-0.02	-0.03	H	0.25	0.15	-0.05
					A	-0.14	-0.03	-0.32	

Table 4: DPD/EOD for subgroups (G): H (Hispanic), N (non-Hispanic), F (Female), M (Male), B (Black), W (White), NI (Unknown), A (Asian). Models: GNN (HiT-GNN), RVL(REVEAL), ft (LLaMA3.2-1B-ft).

zon, while remaining computationally efficient, with ablations confirming the importance of both fine-grained temporal structure and external knowledge enrichment. REVEAL provides a promising balance between performance and explainability through verified rationales. Notably, even large models like GPT-4o show limited zero-shot performance on this complex clinical task. Demographic fairness analysis indicates how clinical data can encode bias without explicit sensitive attribute modeling. Our temporally realistic cohorts, created by filtering post-diagnosis notes, avoid data leakage often overlooked in prior work.

By advancing beyond static health snapshots to dynamic patient representations, this work captures the complexity of real-world clinical trajectories through event-centric modeling that offers low-

resource, privacy-preserving alternatives adaptable to cross-institutional settings. Finally, our framework could also be potentially extended to other chronic diseases, as well as to non-clinical domains that involve longitudinal reasoning over rich textual data, such as financial forecasting or conflict prediction using news reports.

## 9 Ethics Statement

Use of the PH corpus was approved by the institutional review board (IRB) at the University of Illinois, and MIMIC-IV was reviewed by the IRB at Beth Israel Deaconess Medical Center. Both datasets were obtained after completing the required training and accessed through secure protocols in accordance with institutional and data use policies (see <https://physionet.org/content/mimiciv/1.0/> for MIMIC-IV).

While our models achieve strong results, they are not intended to be used as standalone diagnostic tools and must be deployed with clinicians at the center. The necessary next steps are ongoing validation with clinical collaborators, user studies, and explainability audits. Future systems should incorporate active learning loops that capture clinician feedback and investigate bias mitigation to ensure equitable early detection across all subgroups.

## 10 Limitations

While the PH corpus offers a more holistic and longitudinal view of a patient’s clinical journey, it cannot be publicly shared due to privacy constraints. To promote reproducibility, we will make our data curation scripts and code publicly available. Future work could address several key areas. First, mitigating demographic and institutional biases remains an important direction to ensure equi-

table model performance. Second, aligning graph-derived insights with LLM-generated rationales may enhance interpretability and foster clinician trust. Additionally, incorporating multimodal data by integrating text-extracted temporal graphs with structured EHR and medical imaging could further enrich the representations. Cross-institutional adaptation to handle diverse entity types and formats is another promising direction.

## References

2025. *IDF Diabetes Atlas*, 11 edition. International Diabetes Federation, Brussels, Belgium.
- Alan R Aronson and François-Michel Lang. 2010. An overview of metamap: historical perspective and recent advances. *Journal of the American Medical Informatics Association*, 17(3):229–236.
- Olivier Bodenreider. 2004. The unified medical language system (umls): integrating biomedical terminology. *Nucleic acids research*, 32(suppl\_1):D267–D270.
- Daniel Capurro, Erik van Eaton, Robert Black, and Peter Tarczy-Hornoch. 2014. Availability of structured and unstructured clinical data for comparative effectiveness research and quality improvement: a multisite assessment. *EGEMS*, 2(1).
- Rochana Chaturvedi, Peyman Baghersahi, Sourav Medya, and Barbara Di Eugenio. 2025. [Temporal relation extraction in clinical texts: A span-based graph transformer approach](#). In *Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 25765–25788, Vienna, Austria. Association for Computational Linguistics.
- Rochana Chaturvedi, Mudassir Rashid, Brian T Layden, Andrew Boyd, Ali Cinar, and Barbara Di Eugenio. 2023. Sequential representation of sparse heterogeneous data for diabetes risk prediction. In *2023 IEEE International Conference on Bioinformatics and Biomedicine (BIBM)*, pages 831–834. IEEE.
- Hejie Cui, Alyssa Unell, Bowen Chen, Jason Alan Fries, Emily Alsentzer, Sanmi Koyejo, and Nigam Shah. 2025. Timer: Temporal instruction modeling and evaluation for longitudinal clinical records.
- Kirstie K Danielson, Brett Rydzon, Milena Nicosia, Anjana Maheswaren, Yuval Eisenberg, Janet Lin, and Brian T Layden. 2023. Prevalence of undiagnosed diabetes identified by a novel electronic medical record diabetes screening program in an urban emergency department in the us. *JAMA Network Open*, 6(1):e2253275–e2253275.
- Marzyeh Ghassemi, Marco A. F. Pimentel, Tristan Naumann, Thomas Brennan, David A. Clifton, Peter Szolovits, and Mengling Feng. 2015. A multivariate timeseries modeling approach to severity of illness assessment and forecasting in icu with sparse, heterogeneous clinical data. In *Proceedings of the Twenty-Ninth AAAI Conference on Artificial Intelligence*, AAAI’15, page 446–453. AAAI Press.
- Yu Gu, Robert Tinn, Hao Cheng, Michael Lucas, Naoto Usuyama, Xiaodong Liu, Tristan Naumann, Jianfeng Gao, and Hoifung Poon. 2021. Domain-specific language model pretraining for biomedical natural language processing. *ACM Transactions on Computing for Healthcare (HEALTH)*, 3(1):1–23.
- Paul Hager, Friederike Jungmann, Robbie Holland, Kunal Bhagat, Inga Hubrecht, Manuel Knauer, Jakob Vielhauer, Marcus Makowski, Rickmer Braren, Georgios Kaissis, and 1 others. 2024. Evaluation and mitigation of the limitations of large language models in clinical decision-making. *Nature medicine*, 30(9):2613–2622.
- William J Hall, Mimi V Chapman, Kent M Lee, Yessenia M Merino, Tainayah W Thomas, B Keith Payne, Eugenia Eng, Steven H Day, and Tamera Coyne-Beasley. 2015. Implicit racial/ethnic bias among health care professionals and its influence on health care outcomes: a systematic review. *American journal of public health*, 105(12):e60–e76.
- Kai He, Rui Mao, Qika Lin, Yucheng Ruan, Xiang Lan, Mengling Feng, and Erik Cambria. 2024. [A survey of large language models for healthcare: from data, technology, and applications to accountability and ethics](#). *Preprint*.
- William R Hersh, Mark G Weiner, Peter J Embi, Judith R Logan, Philip RO Payne, Elmer V Bernstam, Harold P Lehmann, George Hripcsak, Timothy H Hartzog, James J Cimino, and 1 others. 2013. Caveats for the use of operational electronic health record data in comparative effectiveness research. *Medical care*, 51(8 0 3):S30–S37.
- Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, Weizhu Chen, and 1 others. 2022. Lora: Low-rank adaptation of large language models. *ICLR*, 1(2):3.
- Kexin Huang, Abhishek Singh, Sitong Chen, Edward Moseley, Chih-Ying Deng, Naomi George, and Charolotta Lindvall. 2020. Clinical xlnet: Modeling sequential clinical notes and predicting prolonged mechanical ventilation. In *Proceedings of the 3rd Clinical Natural Language Processing Workshop*, pages 94–100.
- Wenhao Huang, Zijia Lin, Chris McConnell, and Börje F. Karlsson. 2017. [Recognizers-Text: Recognition and resolution of numbers, units, and date/time entities expressed across multiple languages](#).
- Pengcheng Jiang, Cao Xiao, Adam Richard Cross, and Jimeng Sun. 2024. Graphcare: Enhancing healthcare predictions with personalized knowledge graphs. In *The Twelfth International Conference on Learning Representations (ICLR)*.

- Alistair EW Johnson, Lucas Bulgarelli, Lu Shen, Alvin Gayles, Ayad Shammout, Steven Horng, Tom J Pollard, Sicheng Hao, Benjamin Moody, Brian Gow, and 1 others. 2023. MIMIC-IV, a freely accessible electronic health record dataset. *Scientific data*, 10(1):1.
- Zachary Chase Lipton, David C. Kale, Charles Elkan, and Randall C. Wetzel. 2016. [Learning to diagnose with LSTM recurrent neural networks](#). In *4th International Conference on Learning Representations, ICLR 2016, San Juan, Puerto Rico, May 2-4, 2016, Conference Track Proceedings*.
- Ramon Maldonado, Meliha Yetisgen, and Sanda M Harabagiu. 2019. Adversarial learning of knowledge embeddings for the unified medical language system. *AMIA Summits on Translational Science Proceedings*, 2019:543.
- Chuiheng Meng, Loc Trinh, Nan Xu, James Enouen, and Yan Liu. 2022. Interpretability and fairness evaluation of deep learning models on MIMIC-IV dataset. *Scientific Reports*, 12(1):7166.
- Stephanie A Miller and Lenhart K Schubert. 1990. Time revisited 1. *Computational Intelligence*, 6(2):108–118.
- Tuan Dung Nguyen, Thanh Trung Huynh, Minh Hieu Phan, Quoc Viet Hung Nguyen, and Phi Le Nguyen. 2024. Carer-clinical reasoning-enhanced representation for temporal health risk prediction. In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pages 10392–10407.
- Ramesh S Patil, Peter Szolovits, and William B Schwartz. 1981. Causal understanding of patient illness in medical diagnosis. In *Computer-Assisted Medical Decision Making*, pages 272–292. Springer.
- Paul R Rosenbaum and Donald B Rubin. 1985. Constructing a control group using multivariate matched sampling methods that incorporate the propensity score. *The American Statistician*, 39(1):33–38.
- Sarah Sandmann, Stefan Hegselmann, Michael Fajarski, Lucas Bickmann, Benjamin Wild, Roland Eils, and Julian Varghese. 2025. Benchmark evaluation of deepseek large language models in clinical decision-making. *Nature Medicine*, pages 1–1.
- Karan Singhal, Tao Tu, Juraj Gottweis, Rory Sayres, Ellery Wulczyn, Le Hou, Kevin Clark, Stephen Pfohl, Heather Cole-Lewis, Darlene Neal, and 1 others. 2023a. [Towards expert-level medical question answering with large language models](#). *Preprint*.
- Pankhuri Singhal, Lindsay Guare, Colleen Morse, Anastasia Lucas, Marta Byrska-Bishop, Marie A Gueraty, Dokyoon Kim, Marylyn D Ritchie, and Anurag Verma. 2023b. Detect: Feature extraction method for disease trajectory modeling in electronic health records. *AMIA Summits on Translational Science Proceedings*, 2023:487.
- Charlie Snell, Jaehoon Lee, Kelvin Xu, and Aviral Kumar. 2024. [Scaling llm test-time compute optimally can be more effective than scaling model parameters](#). *Preprint*, arXiv:2408.03314.
- Guanchu Wang, Junhao Ran, Ruixiang Tang, Chia-Yuan Chang, Chia-Yuan Chang, Yu-Neng Chuang, Zirui Liu, Vladimir Braverman, Zhandong Liu, and Xia Hu. 2024b. [Assessing and enhancing large language models in rare disease question-answering](#). *Preprint*.
- Yongxin Xu, Kai Yang, Chaohe Zhang, Peinie Zou, Zhiyuan Wang, Hongxin Ding, Junfeng Zhao, Yasha Wang, and Bing Xie. 2023. Vecocare: Visit sequences-clinical notes joint learning for diagnosis prediction in healthcare data. In *IJCAI*, volume 23, pages 4921–4929.
- Shao Zhang, Jianing Yu, Xuhai Xu, Changchang Yin, Yuxuan Lu, Bingsheng Yao, Melanie Tory, Lace M Padilla, Jeffrey Caterino, Ping Zhang, and 1 others. 2024. Rethinking human-ai collaboration in complex medical decision making: A case study in sepsis diagnosis. In *Proceedings of the CHI Conference on Human Factors in Computing Systems*, pages 1–18.
- Hongjian Zhou, Fenglin Liu, Boyang Gu, Xinyu Zou, Jinfa Huang, Jing Wu, Yiru Li, Sam S Chen, Peilin Zhou, Junling Liu, and 1 others. 2024a. [A survey of large language models in medicine: Progress, application, and challenge](#). *Preprint*.
- Yue Zhou, Barbara Di Eugenio, and Lu Cheng. 2025a. [Unveiling performance challenges of large language models in low-resource healthcare: A demographic fairness perspective](#). In *Proceedings of the 31st International Conference on Computational Linguistics*, pages 7266–7278, Abu Dhabi, UAE. Association for Computational Linguistics.
- Yue Zhou, Barbara Di Eugenio, and Lu Cheng. 2025b. Unveiling performance challenges of large language models in low-resource healthcare: A demographic fairness perspective. In *Proceedings of the 31st International Conference on Computational Linguistics*, pages 7266–7278, Abu Dhabi, UAE. Association for Computational Linguistics.

## A Data Curation Details

### A.1 PH Data

We show the data-preparation steps in Table 5. The initial data is 164,910 patients, 15,140 of whom were diagnosed with type 2 diabetes, and 3,708,168 clinical notes. We exclude all post-diagnosis notes (these include notes after 3 days before the day of diagnosis). The 3-day pre-diagnosis window is considered to account for possible lag in transforming information from notes to ICD codes in structured tables. This filter ensures we only consider pre-diagnosis information, leaving us with around 3 million notes. We then exclude the rare

demographic groups. These include race groups Native Hawaiian/Other Pacific Islander/American Indian/Alaska Natives/Asian, and unknown Gender/Ethnicity. We filter the data by note type and author, as it includes all providers and note types, some of which are uninformative. This subset of notes is determined with the help of a senior expert endocrinologist and is provided below:

**Important Note Types** Inpatient Progress Note, Progress Notes, Family Medicine Note, History and Physical Note, Emergency Department Note, Discharge Summary, Nutrition Note, Endocrinology Note, Specialty Pharmacy Services, General Eye Note, Diet Instructions, ED Notes, Outpatient Pharmacy Note, Discharge Note, Dialysis Rounding, Diabetes Education.

**Important Author Types** Unspecified, Anesthesiologist, Care Coordinator, Dentist, Dietician Student, Fellow, Licensed Clinical Social Worker, Licensed Practical Nurse, Medical Assistant, Medical Student, Mental Health Counselor, Midwife, Nurse Practitioner, Nursing Student, Occupational Therapist, Occupational Therapy Assistant, Occupational Therapy Student, Optometrist, Oral Surgeon, Pharmacist, Physical Therapist, Physical Therapy Assistant, Physical Therapy Student, Physician, Physician Assistant, Registered Dietitian.

We then remove notes shorter than 100 words, following previous work in health-risk prediction (Ghassemi et al., 2015). Upon manual inspection of some of the notes in this filtered subset from the pre-diagnosis date, we find that a diagnosis of T2D is already present but not entered in the structured records, or it is entered much later. As discussed previously, such ICD-coding errors are common (Hersh et al., 2013). This motivated us to apply an additional filter using a popular large language model (LLM), LLaMA3.1-8B. LLMs are exceptionally good at extracting explicit information already present in the text. We apply two different prompts (these are variations of the prompt in Figure 5), asking the model to identify if an explicit diagnosis is already present in the note of the T2D cohort (due to scalability challenges, we do not apply the LLM filter to the NoD cohort). Manual inspection of 50 samples reveals higher accuracy if the model answers yes to both prompts. We readjust the earliest diagnosis date accordingly and drop all notes on or after that date. This further reduces our T2D subset to 1717 patients. We split this data into a training-testing split in an 80:20

ratio, stratified by the diagnosis group.

Filtering Step	#NoD	#T2D	#Notes
Initial	164,910	15,140	3,708,168
Pre-diagnosis notes	164,910	4,138	3,146,140
Non-rare demographics	149,673	3,878	2,909,417
Important Note and Author-types	108,635	2,369	1,451,626
Note length	106,222	2,302	1,050,839
LLM for diagnosis date adjustment	106,222	1,717	1,039,043

Table 5: Data filtering steps with the number of non-Diabetic (NoD) and type 2 diabetes (T2D) patients, and the number of notes after each step for the PH corpus.

#### Prompt

You are a helpful assistant who can read and identify required information from given clinical notes, either mentioned explicitly or in ICD code format.

**Required information:** a diagnosis of **Diabetes** for the given patient.

Please reply True if an explicit diagnosis is mentioned in the note.

Please answer False if:

- there is no diagnosis,
- the note negates the diagnosis (e.g., “no diagnosis of diabetes” or “diabetes: negative”),
- if the diagnosis is associated with someone other than the patient (e.g., family).

Do not output anything other than True/False.

**NOTE:** [NOTE\_TEXT]

**Output:**

Figure 5: Prompt for identifying mention of type 2 diabetes in the initial set of pre-diagnosis notes identified based on ICD codes.

We then apply propensity-weighted matching to construct a treatment (T2D) and control group (NoD) matched test sample (Rosenbaum and Rubin, 1985). We use the demographic attributes to match the covariates. For this, we first train a logistic regression model to estimate propensity scores (likelihood of being in the T2D group, given the demographic attributes) (the model achieves an AUC score of 73.40%). We use a greedy nearest-neighbor 1:1 matching without replacement to con-



struct the final control group (NoD). For each treated sample (T2D), we find an unmatched control (NoD) with the closest propensity score, ensuring each control is used only once. Table 6 presents the covariate balance before and after the matching.

Covariate	Before	After
AGE	0.683	0.004
Binarized GENDER (M=1)	0.028	0
Binary RACE NI	-0.057	0
Binarized RACE W	-0.270	0.005
Binarized ETHNICITY (Y=1)	0.002	0

Table 6: Covariate balance regarding standardized mean differences (SMD) before and after propensity score matching for PH corpus test set.

## A.2 MIMIC-IV T2D Data

We select the treatment group as the subset of patients with ICD code E11 (Type 2 diabetes) among the diagnoses. We exclude patients with ICD codes from related diseases (E08-E13, O24) to exclude other types of diabetes. We coarsify race as Hispanic if Race is Hispanic/Latino, BLACK/AFRICAN AMERICAN if race is either Black/African, and drop patients with race group listed as ‘other’. Thereafter, we only retain pre-diagnosis notes. MIMIC-IV notes contain discharge summaries from which we retain “Chief Complaint”, “History of Present Illness”, and “Discharge Instructions” sections and drop all patients with more than 20 admissions following earlier works. We then drop duplicated notes for a patient, and notes that are too short (<.05 quantile) and too long (>0.95 quantile). Finally, we apply an LLM filter on the T2D subgroup to correct the diagnosis date. We use both LLaMA3.1-8B and GPT-4o and manual cleaning to construct a robust subset. Finally, we apply propensity-weighted matched sampling to get a balanced control group. The data filtering steps are detailed in Table 7. The AUC score of the logistic regression model to estimate the propensity score is 63.64%. The covariate balance before and after matching is presented in Table 8. The data is split into training and test in a 90:10 ratio.

Filtering Step	#Notes	#NoD	#T2D
Initial Data	4,756,326	180,640	
No other types of diabetes	4,578,799	178,524	
race not ‘other’	4,033,399	156,513	
Assign cohort	4,033,399	142,412	14,101
pre-diagnosis notes	3,220,793	142,412	5,952
Drop if more than 20 admissions	235,727	59,566	5,847
No Duplicates	184,421	54,803	5,272
Notes length in .05-.95 quantile	175,201	54,170	5,207
No duplicates across sections	174,995	54,170	5,207
LLM + matched controls Filter	14,957	2,909	2,909

Table 7: Data filtering steps with the number of non-Diabetic (NoD) and type 2 diabetes (T2D) patients, and number of notes after each step for MIMIC-IV subset.

Covariate	Before	After
AGE	0.29	0.001
Binarized GENDER (M)	0.12	-0.001
Binary RACE (B)	0.316	0
Binarized RACE (H)	0.110	0.001
Binarized RACE (W)	-0.334	0

Table 8: Covariate balance in terms of standardized mean differences (SMD) before and after propensity score matching on MIMIC-IV test set.

## B Implementation Details

### B.1 Prompts

Figures 6–8 present the prompts used for various LLM-based experiments.

### B.2 Experimental Settings

**GNN** For GNN experiments, we use graphs from the last five notes (we experimented with 1,2,5,10). We train the models with 3-fold stratified cross-validation and save the best model over epochs on each fold. For inference on the held-out test set, we take the average of prediction probabilities from all three models and make the final prediction based on that.<sup>16</sup> SAGEConv operator from PyTorch Geometric with default settings. We also experiment with other GNN variants such as graph attention network (GAT), relational graph convolutional operator (R-GCN), and heterogeneous GNN using heterogeneous graphs. None of these variants performs comparably. We experiment with 1 – 5 GNN layers and find the best performance with  $k = 2$  layers. For aggregating the nodes for graph representation, we also experiment with max pooling and pooling only the nodes corresponding to document creation time (DCT/anchor nodes); neither

<sup>16</sup>We also tried a usual training-validation split in an 80:20 ratio and get almost similar results.

System Prompt: Reasoner

You are a helpful medical assistant. Based on the patient's information and the history of medical notes from one or more visits, assess if this patient has a likelihood of being diagnosed with Type 2 diabetes in the near future.

Make the assessment based on **concrete** medical evidence and how the patient's condition has progressed. Do not solely rely on age as a risk factor.

**# Prediction Format:**  
**## Risk of Type 2 Diabetes:** \*\*[True/False]  
**## Explanation:** [Provide a brief explanation of the reasoning for the prediction.]

User Prompt

Case History:

Patient is {age} y.o. {race} {ethnicity} {gender}.

Note Date: {date1}  
{note-type1} note  
{note-text1}

Note Date: {date2}  
{note-type2} note  
{note-text2}

*Note: {ethnicity} and {note-type} are excluded for the MIMIC-IV. {ethnicity} is also excluded for the PH corpus if the value is 'NI'.*

Figure 6: Prompts for inference from a reasoning model for type 2 diabetes prediction.

System Prompt: Fine-tuned Reasoner

You are a helpful medical assistant. Based on the patient's information and the history of medical notes from one or more visits, assess if this patient has a likelihood of being diagnosed with Type 2 diabetes in the near future.

You must respond with either:

- true
- false

No explanation is needed. Only output one of the two labels above.

Figure 7: System prompt used for fine-tuning and inference from a reasoning model for type 2 diabetes Risk Prediction. The User prompt is the same as in Figure 6.

System Prompt: Fine-tuned Verifier

You are a medical reasoning assistant. Given a sequence of clinical notes from a patient and an analysis predicting the risk of an imminent type 2 diabetes (T2D) diagnosis, judge if the analysis is correct or incorrect.

You must respond with either:

- correct
- incorrect

No explanation is needed. Only output one of the two labels above.

User Prompt

Case History:

Patient is {age} y.o. {race} {ethnicity} {gender}.

Note Date: {date1}  
{note-type1} note  
{note-text1}

Note Date: {date2}  
{note-type2} note  
{note-text2}

Analysis  
{Reasoner Output}

Figure 8: Prompts used for fine-tuning the verifier model in the REVEAL framework that tests if a type 2 diabetes risk prediction model has made a correct prediction based on given input and model reasoning.

works well. In place of an RNN to model multi-document graphs, we also perform simple attention-based aggregation of cross-visit graph representations and transformer-based aggregation, but both underperform compared to BiLSTM.

**LLM** Due to computational constraints, we only fine-tune a smaller LLaMA3.2-1B model instead of the larger LLaMA3.1-8B variant, which exhausts memory even with a batch size of 1 on an A100 80GB GPU. For REVEAL, we use LLaMA3.1-8B as the zero-shot reasoner and LLaMA3.2-1B as the fine-tuned verifier. We experimented with 1, 2, 3, 5, 8, and 10 notes and found that including the last two notes yields the best predictive performance. Since the LLM looks at detailed full-text notes, longer context with additional notes causes it to underperform. On the other hand, since GNN looks at a structured summary of the notes in the form of temporal graphs, it can incorporate additional information effectively and gives the best performance with five notes. We also evaluated the effect of demographic verbalization on LLM. We observed a slight improvement in accuracy when demographic details were included, while demographic integration reduces the performance of our GNN-based models.

All fine-tuning experiments are conducted on LLaMA3.2-1B using a batch size of 1 due to the length of input sequences and memory limitations. We train the model for 30 epochs using the AdamW optimizer with a learning rate of  $6e^{-06}$  and a weight decay of 0.01.<sup>17</sup> A linear learning rate schedule without warm-up is applied. Gradient accumulation is used with a factor of 2 to stabilize updates, given the small batch size. The loss function is standard cross-entropy with mean reduction. In our setup, LoRA is applied to the query and value projection layers (q\_proj and v\_proj) within the transformer blocks. We use a rank of 8, a LoRA scaling factor (alpha) of 16, and a dropout rate of 0.1 during training. This configuration results in only a small fraction of the model parameters being updated—approximately 0.07% (approximately 850K out of 1.24B)—making the training process efficient (31 minutes per epoch for PH corpus and 32 minutes for MIMIC-IV) and tractable on a single A100 80GB GPU while still allowing the model to adapt to the binary classification task.

In the fine-tuning experiments, the pretrained

<sup>17</sup>We experimented with learning rates of  $5e^{-05}$  and  $1e^{-06}$  that had lower performance.

LLaMA3.2-1B model initially performs near chance level, with a balanced accuracy of 54.50% on the validation set and 0.8% of outputs falling outside the valid label set. After fine-tuning, the best model achieves an accuracy of 70.50% on the validation set with no invalid predictions. On MIMIC-IV, the initial validation set accuracy is much lower at 44.00%, with 18.1% of predictions being invalid. This improves to 62.50% and zero invalid predictions after fine-tuning.

For the REVEAL model, to ensure proper supervision, we restrict training to examples where the 8B reasoner produced both ‘true’ and ‘false’ predictions across the five runs. This guarantees that each case includes both valid and invalid reasoning paths, enabling the verifier to learn discriminative patterns. We have 5550 training examples for the PH corpus and 9180 for MIMIC-IV. A 5% random held-out validation set stratified by the ‘correct/incorrect’ label is reserved from each. After each epoch, the model is evaluated on the validation set, and the best-performing model checkpoint is saved. The verifier 1B model starts with a base accuracy of 43.2% and 1% invalid predictions on the PH corpus, and the accuracy improves to 84.9% and no invalid predictions in 30 epochs. For MIMIC-IV, the accuracy improves from 58.8% initially to 67.5%.

## C Statistical Significance

To assess whether the observed differences in performance metrics presented in Table 2 are statistically significant, we perform non-parametric bootstrap resampling and compare each pair of the top models—CLSTM, REVEAL, HIT-GNN, and LLaMA3.2-1B-ft. Specifically, for each metric (e.g., macro- $F_1$ , AUC, precision, recall), we repeatedly sample subsets of the test set with replacement and compute the performance difference between model pairs over 10,000 bootstrap replicates. This generates an empirical distribution of differences for each metric and model pair. From this distribution, we computed the mean difference, the sample standard deviation, and a 95% confidence interval (CI) using the 2.5<sup>th</sup> and 97.5<sup>th</sup> percentiles. Statistical significance is determined based on the proportion of bootstrapped differences falling below or above zero, yielding a two-tailed p-value. Significance levels are reported using the convention:  $p < 0.05$  (\*),  $p < 0.01$  (\*\*), and  $p < 0.001$  (\*\*\*). This approach makes no parametric assumptions and ro-

Model Pair	PH Corpus						MIMIC-IV					
	AUC	F <sub>1</sub>	T2D		NoD		AUC	F <sub>1</sub>	T2D		NoD	
			P	R	P	R			P	R	P	R
ReVeAL vs. LLaMA3.2-1B-ft	-0.07* (0.03)	0.04 (0.02)	0.07*** (0.02)	-0.18*** (0.03)	-0.05* (0.03)	0.23*** (0.03)	-0.11* (0.05)	-0.08* (0.04)	-0.09 (0.06)	-0.11* (0.06)	-0.06* (0.03)	-0.03 (0.05)
CLSTM vs LLaMA3.2-1B-ft	0.03 (0.02)	0.04 (0.02)	0.04** (0.02)	-0.02 (0.03)	0.02 (0.03)	0.09** (0.03)	-0.01 (0.03)	0 (0.01)	0 (0.02)	0.01 (0.01)	0 (0.01)	-0.01 (0.02)
CLSTM vs ReVeAL	0.10*** (0.03)	0 (0.02)	-0.04 (0.02)	0.16*** (0.03)	0.06** (0.03)	-0.14*** (0.03)	0.10 (0.05)	0.08 (0.04)	0.09 (0.06)	0.12* (0.06)	0.06 (0.03)	0.02 (0.05)
ReVeAL vs. HiT-GNN	-0.14*** (0.03)	-0.02 (0.02)	0.02 (0.02)	-0.16*** (0.03)	-0.07** (0.03)	0.11*** (0.03)	-0.10* (0.05)	-0.11** (0.04)	-0.08 (0.05)	-0.27*** (0.05)	-0.11*** (0.03)	0.08 (0.05)
LLaMA3.2-1B-ft vs. HiT-GNN	-0.07*** (0.02)	-0.06** (0.02)	-0.05** (0.02)	0.02 (0.03)	-0.03 (0.03)	-0.12*** (0.03)	0.01 (0.03)	-0.04 (0.03)	0.01 (0.04)	-0.16*** (0.04)	-0.05* (0.03)	0.10** (0.04)
CLSTM vs HiT-GNN	-0.04* (0.02)	-0.02 (0.02)	-0.02 (0.02)	0.0 (0.02)	-0.01 (0.02)	-0.03 (0.03)	0 (0.03)	-0.03 (0.03)	0.01 (0.03)	-0.15*** (0.04)	-0.05* (0.03)	0.10* (0.04)

Table 9: Pairwise model performance differences for top three models overall (first model minus second), based on bootstrap resampling test. Statistical significance: \*  $p < 0.05$ , \*\*  $p < 0.01$ , \*\*\*  $p < 0.001$ .

bustly quantifies uncertainty in model comparisons.  
The results are summarized in Table 9.

## D Ablation Results across T2D Group Recall

Figure 9 presents T2D group recall trends across the three best models (LLaMA3.2-1B-ft, REVEAL and HiT-GNN) with varying prediction windows. Again, HiT-GNN performs either comparably (for PH data) or better (for MIMIC-IV) in the immediate prediction horizons.

Figures 10 and 11 show consistent improvements in T2D group recall from both embeddings and both subgraphs, respectively, across the two corpora.



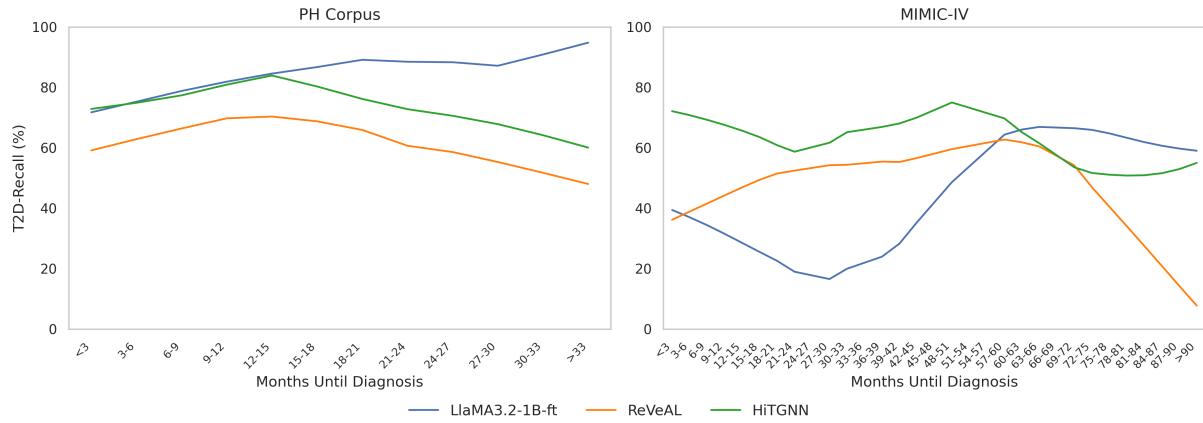


Figure 9: T2D recall as a function of the prediction horizon, evaluated over consecutive 3-month windows.

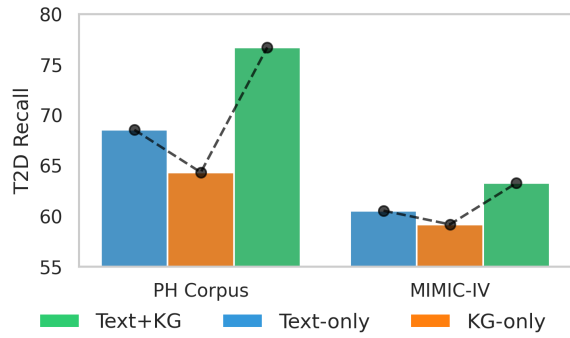


Figure 10: T2D recall performance variation of HiTGNN with different node embedding approaches.

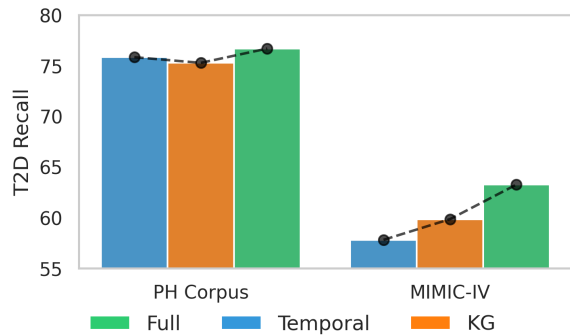


Figure 11: T2D recall performance variation by restricting to specific subgraphs. The KG model includes all nodes extracted from text, all KG nodes and associated edges, but no temporal edges. Temporal model includes only the extracted temporal graphs (excluding KG-nodes and edges). The third contains both.