# TIME CONDITIONED FORESEEING: TEMPORAL GENERATIVE PRETRAINING FOR EHR FOUNDATION MODELS

**Anonymous authors**Paper under double-blind review

000

001

002

004

006

008 009 010

011 012 013

014

015

016

017

018

019

021

025

026

027

028

029

031

033

035

037

040

041

042

043

044

046

047

048

051

052

#### **ABSTRACT**

Electronic Health Records (EHRs) possess unique characteristics that differ significantly from natural language. However, existing models have overlooked these properties and largely relied on Natural Language Processing (NLP) approaches, resulting in suboptimal performance. To address these limitations, we propose a pretraining method designed to effectively capture the distinctive features of EHRs. First, EHRs contain both clinically critical and less informative numerical ranges. To reflect this, we introduce a Pathology-Focused Binning strategy that emphasizes values with clinical significance. Second, both absolute timestamps and relative time intervals are important in EHRs. To incorporate these temporal aspects, we propose a Dual-Calendar Rotary Positional Embedding (RoPE) that jointly encodes complementary temporal signals. Third, many medical applications require modeling long-term patient interactions. Accordingly, we extend conventional next-token prediction with a Time-Conditioned Foreseeing (TCF) objective, enabling the model to forecast long-range clinical events across multiple temporal horizons. Our approach establishes the first genuine temporal generative EHR model, advancing long-range clinical forecasting. It outperforms existing EHR foundation models on seven diverse downstream tasks and enables realistic and temporally consistent EHR generation. All code and models will be made publicly available in the final version of the manuscript.

# 1 Introduction

Electronic health records (EHRs) are longitudinal records that comprehensively document a patient's medical history. EHRs help clinicians assess patient conditions, coordinate diagnostic and therapeutic interventions, and communicate with other healthcare providers (Häyrinen et al., 2008). One of the key objectives in medical AI is to develop models that can learn from EHRs to perform various clinical tasks. However, building such models is challenging due to the complex temporal dependencies and the predominance of numerical data in EHRs (Nasarudin et al., 2024). Recently, there have been growing efforts to leverage large language model (LLM) training paradigms in building EHR foundation models (Niu et al., 2024). Despite these advances, approaches explicitly designed to model the distinct characteristics of EHRs are still in their early stages of development.

EHRs consist of diverse clinical events—such as examinations, treatments, and diagnoses—that are recorded with associated timestamps. Figure 1 illustrates an example EHR, where events are arranged chronologically, and shows how these events can be transformed into a sentence of tokens. Recent preprocessing approaches for EHRs commonly represent a single clinical event as a Time (T), Feature (F), Value (V) triplet (Tipirneni & Reddy, 2022). Here, the Feature denotes attributes such as diagnosis codes, prescribed medications, or laboratory tests (e.g., Systolic Blood Pressure) and represented as a single token, while the Value corresponds to the result or auxiliary information of the Feature (e.g., 87mmHg). Values are typically numerical but may also be absent, or take the form of heterogeneous modalities such as text, depending on the Feature. Despite the necessity of including all triplet components for a faithful representation of clinical events, as indicated in the "data usage" column of Table 1, even the most recent EHR foundation models often exclude Time or Value information due to modeling complexities (Yang et al., 2024).

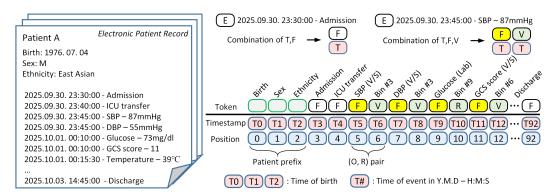


Figure 1: (Left) Extraction of raw patient data from the EHR database in chronological order. (Right) Tokenization of each event (E) with triplet representation, where patient information is placed at the beginning, Features and Values are tokenized, and timestamps remain continuous.

Recent EHR foundation models have improved performance on various downstream tasks through large-scale pre-training. However, most of these models follow standard LLM training paradigms without adapting to the structure and clinical semantics of EHR data (Burkhart et al., 2025), which differ from natural language. For example, converting temporal information into absolute positional embeddings hinders capturing relative intervals and preserving clinically meaningful calendrical information (Likhomanenko et al., 2021). Also, processing numeric Value through *uniform* binning concentrates bins around normal ranges and reduces resolution for pathological states. Moreover, most learning objectives are adopted from language modeling, such as next-token prediction (NTP) or masked language modeling (MLM), without considering EHR-specific characteristics. To address these limitations, we introduce improved binning, temporal embedding, and novel training objectives tailored to EHR data and clinical planning process.

First, we introduce a simple yet effective **Pathology-focused Binning** for Value tokenization. As shown in the "Value Binning" column of Table 1, most EHR models tokenize Value through uniform binning. However, as illustrated in Figure 2A, uniform binning assigns a large amount of bins to physiologic ranges, while allocating only a few bins to clinically important pathologic ranges, thereby limiting the ability to distinguish the severity of abnormalities. Other models rely on false distributional assumptions of Gaussianity, and instead apply standard deviation (std)—based binning (Zhu et al., 2024) or z-normalization (Tipirneni & Reddy, 2022), making them vulnerable to outliers, long-tailed, and dual peaks distributions common in EHR. To address this, we propose a density-based binning that makes no distributional assumptions and focuses on pathological ranges. In this approach, values in high-density physiologic zones are assigned lower weights, whereas values in low-density pathologic zones receive higher weights. This design is suited for all value distributions, and we are the first to apply such binning to EHR models.

Second, we introduce **Dual-Calendar RoPE**, a novel timestamp addressing method for EHRs. Unlike language models, where tokens are assumed to be uniformly spaced, EHRs contain events with highly irregular intervals. Clinically, both relative intervals and calendarical context—e.g., morning/afternoon or weekday/weekend—are important (body temperature is higher in the afternoon, and dialysis complications are common after weekends (Fotheringham et al., 2020)). Also, multiple events may occur at the same time, such as laboratory tests recorded together. As shown in Figure 2B, we partition the dimensions of rotary positional embedding (Su et al., 2024) to jointly encode position and time, assigning calendrical components (e.g., minute, day, month) in increasing units to the time dimension. This enables explicit modeling of distance relations such as "two tests performed at the same time" or "the same test performed at the same hour on different days." The "Time Addressing" column of Table 1 shows that conventional models have not fully addressed crucial temporal information.

Finally, and most importantly, we propose a new learning objective, **Time-Conditioned Foreseeing** (TCF). This objective aligns with the clinical process of treatment planning, and it enables, for the

<sup>&</sup>lt;sup>1</sup>Suppose that SBP, 120mmHg, DBP, 80mmHg are recorded simultaneously. The position dimension provides additional support to prevent the model from confusing results such as SBP, 80mmHg, DBP, 120mmHg.

Table 1: Comparison of recent EHR models by architecture and training objective. **Data usage**: whether the model uses <code>Time</code> and <code>Value</code> information. Value binning: whether values are uniformly binned (STraTS embeds values continuously) and whether bin tokens are shared across <code>Features</code>. **Time addressing**: whether the model considers relative time intervals, calendrical time, and distinguishes concurrent events; same timestamps. **Learning objective**: type of loss, '**Foresee'** - whether model forecasts beyond the next-token, and '**Temporal Generation'** - whether temporal generative modeling is possible. Abbreviations\*: MEP; missing entity prediction, TTE; time to event prediction, TCF; time conditioned foreseeing. Refer to the related works section for details<sup>†</sup>.

	[Data ι	isage]	[Value	Binning]	[	Time addre	ssing]	[Learning Objective]						
Models	Event Timestamp	Numeric Value	Non- uniform	Value sharing	Relative interval	Calendrica time	l Non- concurrency	Туре	Forese	Temporal Generation				
BEHRT (2020) Med-BERT (2021) Foresight (2024) ClinicalMamba (2024) EHR-BERT (2024)	X	X	-	-	X	X	-	NTP&MLM	X	X				
HEART (2024 FM4EHR (2025)	X	О	X	X O	X	X	-	MEP* NTP	X	X				
MOTOR (2024)	О	X	-	-	0	X	X	TTE*	О	X				
STraTS (2022)	О	0	0	О	X	X	-	MSE	X	X				
EHRSHOT (2023)	О	0	X	X	X	X	-	NTP	X	X				
TRADE (2024)	О	0	0	X	X	X	-	MLM	X	X				
EHRMamba (2025)	О	0	X	X	X	X	-	NTP&MLM	X	X				
ETHOS (2024)	О	0	X	О	Δ <sup>†</sup>	X	-	NTP	X	Δ <sup>†</sup>				
OURS	0	0	0	both	0	0	0	TCF*	O	0				

first time, generative temporal modeling of a patient's medical timeline. As shown in the "Learning Objective" column of Table 1, prior models have relied on objectives designed for language models or variants thereof, with the exception of the time-to-event (TTE) objective. Conventional EHR models trained with NTP loss capture only  $P(F_{next} \mid E_{past})$ , without explicitly modeling temporal information. Consequently, they cannot distinguish whether an event occurs minutes later or after many hours, treating both urgent and routine vital sign measurements (short and long time intervals respectively) identically as the 'next token.'

In contrast, TCF explicitly models long-range temporal information, thereby capturing how real-world clinical practice unfolds over time. In NLP, missing a single token disrupts grammar, and consecutive tokens are tightly correlated. By contrast, neighboring EHR events are loosely connected and often exhibit long-range dependencies, such as 8-hour follow-up tests. This reflects clinical practice, where physicians do not always act in real time but instead devise broader clinical plans. TCF embodies this principle: rather than the short-sighted scope of NTP, which predicts only the immediate next event, TCF enables questions such as, "What intervention is needed in the next six hours?" To achieve this, TCF module first generates the next timestamp from the last hidden state. The multiple foreseeing timestamps are then fed back as module inputs, conditioning subsequent token generation. This time-conditioned architecture allows simultaneous learning of  $P(T_{next} \mid E_{past})$  and  $P(F_{foresees} \mid T_{foresees}, E_{past})$ , leading to improved performance.

Our model ranked first across all combinations of the three dataset configurations and seven diverse downstream tasks. Across these tasks, the AUPRC was consistently improved, reaching up to 48% higher than that of the second-best model, highlighting a clinically meaningful improvement given the data imbalance. We also demonstrated that the model generates temporally stable, realistic EHR records and is capable of leveraging the calendrical component in generative modeling.

Our contributions can be summarized as follows:

- Pathology-Focused binning: Introduces density-adjusted binning to the EHR foundation model, focusing on clinically relevant pathologic ranges.
- Dual-Calendar RoPE: Simultaneously represents both calendrical time and positional information, allowing model to capture calendrical periodicity and event concurrency.
- Time Conditioned Foresee Objective: Enables clinically aligned foreseeing training and temporal generative modeling of patient medical timelines.

# 2 RELATED WORKS

EHR foundation models differ from medical specialist LLMs, which rely on patient history texts summarized by clinicians. EHR foundation models learn directly from raw EHR events (Burkhart et al., 2025) and have been applied to various downstream clinical tasks (Table 1).

A common practice in EHR modeling is to represent each EHR event as a triplet of Time, Feature, and Value (Tipirneni & Reddy, 2022; Lee et al., 2023). However, many models exclude temporal and numeric data, as they are difficult to handle in standard language model frameworks. For instance, BEHRT (Li et al., 2020), Med-BERT (Rasmy et al., 2021), and others rely solely on discrete Features, omitting critical information and limiting their utility.

Some models incorporate numeric Values but omit Time. HEART (Huang et al., 2024) discretize Values into uniform bins, mapping Feature–Value pairs to single tokens. This approach inflates the vocabulary size, leading to data sparsity. FM4EHR (Burkhart et al., 2025) addresses this by tokenizing Features and Values separately, allowing tokens to be shared.

In contrast, MOTOR (Steinberg et al., 2024) models Time but not Value, performing survival analysis by treating each feature's occurrence as an endpoint. Its utility is limited by its inability to handle numeric values, low temporal expressiveness based on pre-defined intervals, unrealistic constant hazard assumption, and a quadratic complexity that hinders practical application. Moreover, encoding timestamp as 'days since birth' with RoPE does not account for calendrical time.

STraTS (Tipirneni & Reddy, 2022) tokenizes only the Feature, embedding Value and Time as continuous variables to predict the next value. By modeling only  $P(V_{next} \mid E_{past})$ , it loses important context and cannot support generative modeling.

TRADE (Zhu et al., 2024) and EHRmamba (Fallahpour et al., 2025) used MLM/NTP paradigms, discretizing values and applying absolute positional embeddings to Feature and Value tokens.

ETHOS (Renc et al., 2024) tokenizes time intervals and insert time-interval tokens between events. This coarse discretization limits medical precision, cause cumulative errors, and increases computational cost by lengthening the sequence. Unlike positional embeddings, it requires aggregating all intervening tokens to determine a time duration. More details are provided in Appendix A

To address these limitations, this work designs modeling strategies and learning objectives tailored to the unique characteristics of EHR data.

#### 3 Method

**Pathology-Focused Binning**. First, we estimate the value distribution non-parametrically using a Gaussian Kernel Density Estimator (KDE).  $V_{list}^f$  denotes the list of all Values of Feature f in the training set. We uniformly partition the value range  $[\min(V_{list}^f), \max(V_{list}^f)]$  with  $X = \{x_1, x_2, \ldots, x_P\}$ , where the inverval is  $0.05\sigma$ . At each discrete point  $x_k \in X$ , data density  $\rho(x_k)$  is calcuated with Gaussian convolution kernel from all value  $v_j \in V_{list}^f$ . The density is:

$$\rho(x_k) = \sum_{j=1}^{|V_{list}^f|} K_h(x_k - v_j), \quad s.t. \ K_h(u) = \exp\left(-\frac{u^2}{2(0.1\sigma)^2}\right)$$

This allows us to approximate the local density  $\rho(v)$  for any given value v. Then, we assign a weight w(v) to each value that is inversely proportional to its density, effectively giving greater importance to values in sparser region  $(w(v) \propto \rho(v)^{-N}; N \geq 1)$ . In short, values in sparse regions are assigned larger weights than those in dense regions.

Second, these density-based weights are used to construct the final value bins via weighted percentile binning. In this step, the contribution of each unique value  $v_j$  with a raw count of  $c_j$  is scaled by its weight  $w(v_j)$ , creating a weighted count  $c_j' := c_j \cdot w(v_j)$ .

Bin thresholds are then determined from the cumulative distribution of these weighted counts. As a result, high-weight values from pathologic ranges command a larger share of the percentile space, leading to a finer-grained partitioning in these clinically important areas (Figure 2A). The detailed methodology is described in Appendix B.1.

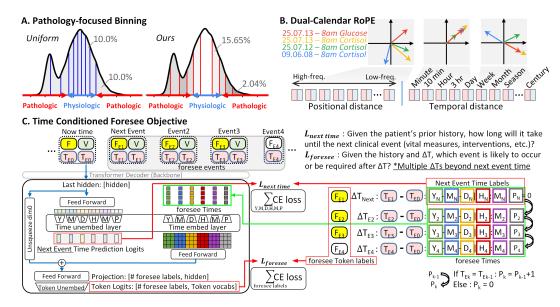


Figure 2: (A) Uniform binning concentrates bins in dense, physiologic ranges. In contrast, our density-based method allocates more bins to medically significant pathologic ranges. (B) Events at the same time are distinguished by their positional distance. Events occurring at the same time on different dates share the same representation for time units below a day but have different representations for units of a day or longer. (C) Illustrates TCF objective of a single timestep (actual model training is fully parallel, like NTP). The TCF objective consists of  $L_{next\ time}$  and  $L_{foresee}$ . The last hidden state is passed through a time head to predict the interval to the next event in a calendrical format ( $L_{next\ time}$ ). Then, the times to multiple future events are re-input and combined with the last hidden state to predict the events at those specific times ( $L_{foresee}$ ).

**Dual-Calendar Rotary Position Embedding**. Second, we propose a novel positional encoding designed for the temporal characteristics of EHR (Figure 2). It jointly models the relative order and calendrical interval by partitioning the dimension of each query and key vector,  $x \in \mathbb{R}^d$ , into a positional component  $x_{pos} \in \mathbb{R}^{d_{pos}}$  and a temporal component  $x_{time} \in \mathbb{R}^{d_{time}}$  ( $d = d_{pos} + d_{time}$ ):

$$x = [x_{pos} \parallel x_{time}]$$

The  $x_{pos}$  component uses a standard RoPE to encode the relative token position, p. With a reduced dimensionality ( $d \to d_{pos}$ ), it employs a truncated frequency spectrum. This strategic choice focuses its role on disambiguating the order of co-occurring events sharing an identical timestamp, while long-range dependencies are handled by the temporal component. The rotation angle is defined as:

$$\theta_{p,i}^{(pos)} = \frac{p}{10000^{2i/d}} \;\;,\;\; i \in \{0,1,...,d_{pos}/2-1\}$$

The core of our method, the  $x_{time}$  component, encodes the second-level timestamp t. This is achieved using a predefined set of semantically meaningful calendrical periods (e.g., minute=60s, hour=3600s,...; see Table 5 for a full list). For each period  $s_j$  in the set, a rotation angle  $\theta_{t,j}^{(time)}$  is calculated as the phase of the event within that period:

$$\theta_{t,j}^{(time)} = \left(\frac{t \pmod {s_j}}{s_j}\right) \cdot 2\pi \ , \ j \in \{0, 1, ..., d_{time}/2 - 1\},$$

The two components are rotated independently using their respective angles and then concatenated to form the final query vector q' (and also for the key). This allows the attention mechanism to simultaneously address both sequential order and calendrical time (More details in Appendix B.2).

$$q' = [\text{RoPE}(q_{pos}, \theta^{(pos)}) \parallel \text{RoPE}(q_{time}, \theta^{(time)})]$$

**Time-Conditioned Foresee Objective (TCF)**. Lastly, we propose a novel learning objective to effectively model the temporal dynamics of EHR data. TCF employs a dual-objective structure (Figure 2C) that simultaneously learns to: (1) predict *when* the next event will occur  $(P(\Delta T_{\text{next}}|E_{\text{past}}))$ ,

and (2) foresee what event will happen at a specified future time ( $P(F_{\text{foresee}} | \Delta T_{\text{foresee}}, E_{\text{past}})$ ), unlike NTP which only models  $P(F_{\text{next}} | E_{\text{past}})$ .

The TCF module is placed after the transformer backbone. It takes the final hidden state  $h_{\text{last}} \in \mathbb{R}^{d_{\text{model}}}$  as input<sup>2</sup> and outputs both a next time prediction loss and a conditioned hidden states for future event prediction.

To generate a calendrical ground-truth label for  $\Delta T_{\rm next}$ , the time delta, expressed in seconds, is transformed into an integer vector of dimension  $N_{\rm scales}$ . Each element of this vector corresponds to a predefined calendrical time unit, ranging from 10-year to 1-minute (e.g.,  $[\alpha, \beta, \gamma, \dots]$  represents a time composed of  $\alpha$  years,  $\beta$  months,  $\gamma$  days, etc.).

To predict  $\Delta T_{\text{next}}$  from  $h_{\text{last}}$ , the last hidden is projected into # $N_{\text{scales}}$  vectors of size  $d_{\text{embed}}$ . Each of these vectors is transformed into time-logit through the unembedding layer.

$$h_{\text{time}} = \text{FFN}_{\text{enc}}(h_{\text{last}}) \in \mathbb{R}^{(N_{\text{scales}} \cdot d_{\text{embed}})} \rightarrow \{h_{\text{time}}^{(i)}\}_{i=1}^{N_{\text{scales}}} \,, \quad \text{time-logits}^{(i)} = h_{\text{time}}^{(i)} \cdot (W_{\text{embed}}^{(i)})^T$$

 $\mathcal{L}_{\text{next\_time}}$  is Cross-Entropy loss between these  $\{\text{time-logit}_{\text{time}}^{(i)}\}_{i=1}^{N_{\text{scales}}}$  and the calendrical  $\Delta T_{\text{next}}$  labels, averaged over  $N_{\text{scales}}$ .

For foreseeing future events, a **Time-Conditioning** process is performed. We aim to predict the Feature of  $N_{\text{foresee}}$  future events. A given future time deltas,  $\Delta T_{\text{foresee}}$ , is first transformed into a vector of integer labels ( $C_{\text{foresee}} \in \mathbb{Z}^{N_{\text{foresee}} \times N_{\text{scales}}}$ ) using the same multi-scale decomposition. These labels are passed through embedding layers to produce a comprehensive time embedding,  $e_{\text{time}} \in \mathbb{R}^{N_{\text{foresee}} \times (N_{\text{scales}} \cdot d_{embed})}$ . Finally, this time embedding is fused with the original hidden state  $h_{\text{last}}$  via a residual connection to produce a time conditioned hidden state,  $h_{\text{conditioned}}$ .

$$h_{\text{conditioned}} = \text{FFN}(\text{LayerNorm}(h_{\text{last}} + \text{FFN}(e_{\text{time}}))) \in \mathbb{R}^{N_{\text{foresee}} \times d_{\text{model}}}$$

This  $h_{\text{conditioned}}$  is projected to token-logit  $\in \mathbb{R}^{N_{\text{foresee}} \times \text{vocab\_size}}$  that predicts the clinical event (Feature) that occur at the corresponding future timestamps.

 $\mathcal{L}_{\text{foresee}}$  is Cross-Entropy loss between the future events and the token-logit, averaged over  $N_{\text{foresee}}$ . Through this dual-objective learning ( $\mathcal{L} = \mathcal{L}_{\text{next.time}} + \mathcal{L}_{\text{foresee}}$ ), our model acquires the ability to accurately and generatively model a patient's entire medical timeline.

So far, we have considered the position where Feature is predicted given the previous events. Modeling Value given the previous events and Feature is carried out in the same manner. Since F and V belong to the same event and thus share the time label, we always have  $\Delta T_{\text{next}}=0$ . Moreover, because V is conditioned on the preceding F, we predict  $V_{\text{now}}$  by modeling

$$P(V_{\text{foresee}} \mid \Delta T_{\text{foresee}}, F_{\text{now}}, E_{\text{past}})$$

while inserting only a zero into  $\Delta T_{\text{foresee}}$ . More detailed explanation and tensor-level parallel processing are provided in Appendix B.3.

# 3.1 Data and Preprocessing

While many EHR models rely on private datasets and often do not release their code or parameters—making reproduction and evaluation difficult—we use a publicly available dataset and provide open-source code throughout all stages. Specifically, we employ the MIMIC-III Clinical Database v1.4 (Johnson et al., 2016a), which contains comprehensive clinical data from over 30,000 patients. We adopt the widely used preprocessing and train/test split pipeline introduced by Harutyunyan et al. (2019). A summary of the dataset is provided in Table 2. Further details are provided in Appendix C. Tasks necessitating clinical judgment, such as defining exclusion criteria and outlier re-

Table 2: Data summary. Parentheses indicate cases where bins are not shared.

MIMIC-III preprocessed	Train / Test
Total Patient #	28,728 / 5,070
Total Hospitalization #	35,730 / 6,295
Total Events #	38,641,175 / 6,744,906
Total Tokens #	77,109,833 / 13,459,430
Avg. length	2,684 / 2,655
Max length	393,337 / 62,759
Unique Tokens #	155 (1,208)
Token # bin	10
Token # ethnicity	10
Token # vital signs	17
Token # laboratory tests	100

moval, were independently reviewed by an internist, an otolaryngologist, and a general physician.

 $<sup>^2</sup>$ In practice, the full last hidden state  $H_{\text{last}} \in \mathbb{R}^{B \times L \times d_{\text{model}}}$  is processed in parallel, similar to the NTP.

#### 3.2 BACKBONE ARCHITECTURE, BASELINE MODELS, AND PRE-TRAINING

We used a Transformer decoder as the backbone for all experiments. For Dual-calendar RoPE, the first 24 dimensions of the 64-dim K and Q vectors encode positional information, and the remaining 40 dimensions encode calendric time. Baseline models were reproduced under identical conditions, including backbone and training data. We mostly followed the original papers' implementations but made necessary modifications where direct application was infeasible (e.g., adapting the MOTOR model to numeric value events). The pre-training input token length was fixed at 2048. Sequences exceeding this length (Appendix Figure 7) were segmented with a 512-token overlap, and we ensured that a single event's F and V were not split at the segmentation point. A detailed description of our model, baselines, and pre-training can be found in Appendix D.

# 

# 4 RESULT

### 4.1 DOWNSTREAM TASK AND FINE-TUNING

Table 3: Results on downstream tasks using EHR datasets with 117, 17, and 6 features. The Test loss column reports the overall test loss for each feature set (lower is better). For 117 features, we report performance on seven downstream tasks ranging from IHM to Vaso. Binary classification tasks are measured by AUROC (ROC) and AUPRC (PRC), while multiclass tasks are evaluated with macro F1 (Ma-f1) and Cohen's Kappa. For tasks with multiple subtasks, both macro and micro AUROC are reported. We trained our model with and without value sharing; in both cases, it outperformed all other baselines. Full downstream task results are provided in Appendix Table 9-11)

	Tasks	Tes	st Loss	(\dagger)	IH	M	Pl	ne	Dec-	death	Dec-	arrest	L	OS	Н	JO	Va	iso
	Metric	117	17	6	ROC	PRC	macro	micro	ROC	PRC	ROC	PRC	Ma-f1	Kappa	macro	micro	ROC	PRC
re	HEART	5.304	5.434	5.835	0.838	0.442	0.717	0.718	0.869	0.205	0.862	0.199	0.150	0.142	0.703	0.701	0.865	0.363
sha	MOTOR	4.945	5.212	5.645	0.872	0.547	0.770	0.773	0.904	0.272	0.889	0.261	0.174	0.163	0.753	0.748	0.891	0.438
e	EHRSHOT	5.841	6.078	6.341	0.801	0.433	0.634	0.633	0.829	0.167	0.802	0.153	0.101	0.115	0.701	0.614	0.867	0.341
/aJr	TRADE	5.260	5.454	6.048	0.828	0.441	0.738	0.738	0.867	0.170	0.857	0.165	0.158	0.151	0.732	0.731	0.869	0.393
	EHRmamba	5.137	5.439	5.926	0.868	0.557	0.690	0.687	0.901	0.277	0.886	0.260	0.150	0.159	0.751	0.753	0.873	0.399
Z	Ours (No share)	4.686	4.907	5.367	0.889	0.607	0.809	0.816	0.928	0.400	0.917	0.388	0.181	0.185	0.776	0.781	0.912	0.498
re	FM4EHR	6.429	6.389	6.397	0.617	0.177	0.530	0.519	0.744	0.075	0.778	0.102	0.023	0.003	0.598	0.653	0.690	0.130
sha	ETHOS	4.971	5.248	5.572	0.859	0.530	0.739	0.746	0.900	0.311	0.890	0.304	0.170	0.165	0.721	0.731	0.890	0.437
lue	STraTS	5.786	5.812	6.071	0.759	0.311	0.656	0.661	0.840	0.141	0.804	0.103	0.123	0.121	0.590	0.598	0.864	0.331
Va	Ours (Share)	4.879	5.043	5.561	0.876	0.559	0.781	0.784	0.910	0.319	0.902	0.310	0.173	0.170	0.749	0.755	0.906	0.470

We evaluated our model on a range of clinical downstream tasks commonly used in EHR model evaluation. These tasks, defined by clinical labels excluded from training, are not direct measures of generative modeling performance but serve as proxies for the quality of patient representations. In addition to the four MIMIC-III benchmark (Harutyunyan et al., 2019) tasks—In-hospital Mortality (IHM), Decompensation-death (Dec-death), Length of Stay (LOS), and Phenotyping (Phe)—we included three additional tasks: Decompensation-arrest (Dec-arrest), Oliguria/Anuria (HUO), and Vasopressor (Vaso) use. Label

counts for all tasks are provided in Appendix Table 8, with

detailed descriptions in Appendix E.1.

Table 4: Ablation study on Pathology-Focused Binning (Binning), Dual-Calendar Rotary Positional Embedding (Embedding), and Time-Conditioned Foreseeing (Objective).

Binning	Embedding	Objective	Test loss
V	V	V	4.686
Uniform	V	V	4.713
Uniform	RoPE	V	4.810
Uniform	RoPE	NTP	5.241

Downstream task-specific prediction heads were attached to the backbone. Since labels must be inferred using only information up to each timestep, a causal mask was applied for all baselines. To evaluate generalization to data with different distributions (e.g., missing lab information), we experimented with three input configurations: all 117 features, 17 vital signs (without lab data), and only 6 vital signs (SBP, DBP, body temperature, heart rate, respiratory rate, SpO2). Please refer to Appendix E.2 for more details.

Table 3 summarizes the results on downstream tasks. To ensure fair comparison, we trained our model with and without value sharing and compared each setting to the corresponding baselines. In

383

384

385

391

392

394

395

396

397 398 399

400

401

402

403

404

405

406

407

408

409

410

411

412

413

414

415

416

417

418

419

420

421

422

423

424

425

426

427

428

429

430

431

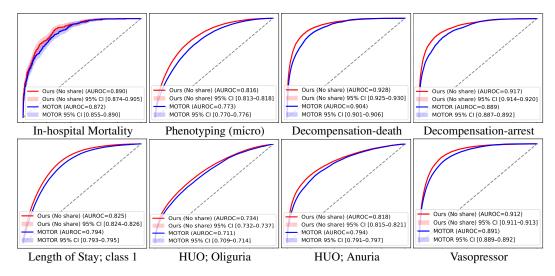


Figure 3: AUROC curves of our model and the second-best baseline, with 95% confidence intervals estimated via bootstrapping. LOS was evaluated as a binary classification for the first class, and Phenotyping was assessed using micro-ROC. For HUO, both oliguria and anuria are presented.

both cases, our model consistently outperformed all baselines across the three input configurations and all downstream tasks. Notably, for the decompensation task, which predicts patient death or arrest up to 24 hours in advance, our model achieved an AUPRC nearly 50% higher than that of the second-best model. Given the severe class imbalance in these tasks (positive:negative ratio of 1:40), this represents a significant improvement in real-world clinical settings where high precision is crucial. Additionally, the ablation study (Table 4) shows that all three of our proposed methods contribute substantially to the performance improvement.

Figure 3 shows the ROC curves with 95% confidence intervals, confirming that our model achieves statistically significant improvements over the second-best model in most tasks. The complete results for all three input configurations can be found in Appendix F.1.

#### 4.2 TEMPORAL GENERATIVE MODELING

Our model is the first to generate fine-grained temporal information and clinical events conditioned on time, demonstrating strong temporal generative modeling of EHR data. To qualitatively assess its effectiveness, we compared it (share ver. for fair comparison) with ETHOS, which, despite limitations in temporal modeling, is one of the few approaches capable of generating temporal information. Since ETHOS outputs time range tokens, timestamps were sampled and rounded to the nearest 5 minutes to match the reso-

```
# [ Ours (Value share) ]
                                     #[ETHOS]
Sex: Female
                                Sex: Female
Ethnicity: Asian
                                Ethnicity: Asian
Age: 81
                                Age: 81
2136-10-02 13:25:04
                                2136-10-02 13:25:04

    ICU transfer

                                  - ICU transfer
2136-10-02 13:49:59
                                2136-10-02 13:49:59
   RR: 12
                                  - RR: 12
2136-10-02 14:00:00
                                2136-10-02 14:00:00
   RR: 18
                                  - RR: 18
  - SBP: 118
                                  - SBP: 118
  - O2 saturation: 99
                                  - O2 saturation: 99
  - MBP: 77
                                  - MBP: 77
  - DBP · 51
                                  - DBP · 51
  - GCS-V : 1
                                  - GCS-V:1
   - GCS · 3
                                   - GCS: 3
                                  - GCS-M · 5
  - GCS-M : 1
  - Temperature: 38.3
                                  - GCS-E: 2
                                  - Temperature: 38.2

    GCS-E : 1

2136-10-02 14:04:00
                                  - HR: 74
  - Potassium (ER): 3.5
                                2136-10-02 14:50:00
  - PO2: 492.0
                                  - RR: 15
  - PEEP: 5.0
                                  - SBP · 96
  - CO2: 26.0
                                  - O2 saturation: 100
  - pH : 7.44
                                  - MBP · 89
   Base excess: 1.0
                                  - DBP: 50
  - Hemoglobin (ER): 10.1
                                  - HR: 78
   PCO2: 37.0
                                2136-10-02 15:35:00
2136-10-02 14:10:00
                                  - RR: 28
  - SBP: 104
                                   O2 saturation: 96
  - RR: 15
                                2136-10-02 16:45:00
  - O2 saturation: 99
                                   - GCS: 9
  - MBP : 69
                                  - GCS-E: 1
   DBP: 45
                                  - DBP: 80
2136-10-02 15:00:00
                                  - GCS-M: 3
   SBP: 125
                                  - O2 saturation: 96
   GCS: 3
                                  - MBP: 86
   O2 saturation: 100
                                  - HR: 76
   MBP: 76
                                  - SBP: 139
  - HR: 80 (truncated rest)
                                  - GCS-V: 1 (truncated rest)
```

Figure 4: Given the initial record (orange), the subsequent medical history is generated (blue). PEEP: Positive end-expiratory pressure, ER: emergency lab.

lution of MIMIC-III (only for ETHOS). For both models, binned measurement values were decoded to actual values by sampling from the empirical distributions of the training data.

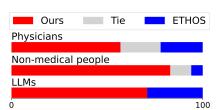


Figure 5: Generated patient EHRs were evaluated by five physicians and five non-medical participants and four LLMs, with 100 comparison responses collected for each category.

Figure 4 presents generated medical history sequences from our model and ETHOS, given the same initial EHR records. From a content perspective, the Glasgow Coma Scale (GCS) should equal the sum of GCS-E/V/M. At 10-02 14:00, our model generated E:1, V:1, M:1; GCS:3, correctly capturing this relationship, whereas ETHOS produced E:2, V:1, M:5; GCS:3, which is inconsistent. Moreover, our model reflected early emergency labs and a variety of tests, followed by routine vital sign checks, while ETHOS generated no labs. From a temporal perspective, our model first performed several tests at short intervals after admission, then naturally returned to an hourly routine. In contrast, ETHOS produced events at irregular intervals and often failed to follow the typical hourly schedule.

We further evaluated 100 generated samples with three evaluator groups: physicians (n=5), non-medical participants (n=5), and commercial LLMs (n=4; ChatGPT (via API, accessed Sep 2025), Gemini 2.5 Flash (via API), 2.5 Pro (via API), Claude 4 Sonnet (via API)). After reviewing up to 10 ground-truth EHR samples, each group assessed subsequent EHR records generated from the same initial records. Figure 5 shows that our model consistently outperformed ETHOS. The LLM input prompts and the generated samples are presented in Appendix F.2.

To verify whether our model effectively integrates calendrical information, we generated vital signs conditioned on time across a 24-hour window (00:00–24:00) based on the same patient history. Figure 6 illustrates that our model generated higher heart rate and temperature values during daytime hours, reflecting realistic circadian variation. In contrast, ETHOS, even for the control variable Height, produced clinically implausible patterns across all cases.

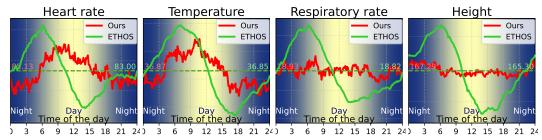


Figure 6: Assessment of the model's ability to capture calendrical temporal patterns. Heart rate, body temperature, and respiratory rate are physiologically higher during the day and lower at night. Using these three features along with height as a control, we let the model sequentially generate predictions across 00:00–24:00, averaged over 1,000 test samples.

#### 5 CONCLUSION

We present a novel approach for modeling the unique characteristics of Electronic Health Record (EHR) data, including irregular time intervals and complex numerical values. This work introduces three key contributions: Pathology-Focused Binning to emphasize clinically significant numerical ranges, Dual-Calendar Rotary Position Embedding (RoPE) to encode relative and absolute calendrical time, and a Time-Conditioned Foreseeing (TCF) training objective. TCF enables temporal generative modeling by predicting future timestamps and forecasting events, reflecting clinical planning. Our model outperforms existing foundation models on seven downstream tasks with up to 48% improvement in AUPRC, while generating realistic and temporally consistent EHRs for long-range clinical forecasting. **Limitations:** There is currently no established metric to evaluate the temporal generative performance of EHR models. Assessing the appropriateness of timing is crucial, making conventional methods used for evaluating LLM generation difficult to apply. Developing quantitative evaluation metrics for EHR generation will be important for advancing EHR foundation models.

# REFERENCES

- Kingma DP Ba J Adam et al. A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 1412(6), 2014.
- Abien Fred Agarap. Deep learning using rectified linear units (relu). arXiv preprint arXiv:1803.08375, 2018.
  - Michael C Burkhart, Bashar Ramadan, Zewei Liao, Kaveri Chhikara, Juan C Rojas, William F Parker, and Brett K Beaulieu-Jones. Foundation models for electronic health records: representation dynamics and transferability. *arXiv preprint arXiv:2504.10422*, 2025.
  - Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pp. 4171–4186, Minneapolis, Minnesota, jun 2019. Association for Computational Linguistics. URL https://aclanthology.org/N19-1423.
  - A Fallahpour, M Alinoori, A Afkanpour, and A Krishnan. Ehrmamba: Towards generalizable and scalable foundation models for electronic health records. In *ML4H*, pp. 291–307, 2025.
  - James Fotheringham, Michael T Smith, Marc Froissart, Florian Kronenberg, Peter Stenvinkel, Jürgen Floege, Kai-Uwe Eckardt, and David C Wheeler. Hospitalization and mortality following non-attendance for hemodialysis according to dialysis day of the week: a european cohort study. *BMC nephrology*, 21(1):218, 2020.
  - Hrayr Harutyunyan, Hrant Khachatrian, David C Kale, Greg Ver Steeg, and Aram Galstyan. Multitask learning and benchmarking with clinical time series data. *Scientific data*, 6(1):96, 2019.
  - Kristiina Häyrinen, Kaija Saranto, and Pirkko Nykänen. Definition, structure, content, use and impacts of electronic health records: a review of the research literature. *International journal of medical informatics*, 77(5):291–304, 2008.
  - Tinglin Huang, Syed Asad Rizvi, Rohan Krishna Thakur, Vimig Socrates, Meili Gupta, David van Dijk, R Andrew Taylor, and Rex Ying. Heart: Learning better representation of ehr data with a heterogeneous relation-aware transformer. *Journal of Biomedical Informatics*, 159:104741, 2024.
  - Alistair Johnson, Tom Pollard, and Roger Mark. MIMIC-III Clinical Database (version 1.4). PhysioNet, 2016a. URL https://doi.org/10.13026/C2XW26. RRID:SCR\_007345.
  - Alistair EW Johnson, Tom J Pollard, Lu Shen, Li-wei H Lehman, Mengling Feng, Mohammad Ghassemi, Benjamin Moody, Peter Szolovits, Leo Anthony Celi, and Roger G Mark. Mimic-iii, a freely accessible critical care database. *Scientific data*, 3(1):1–9, 2016b.
  - Dennis Kasper, Anthony Fauci, Stephen Hauser, Dan Longo, J Jameson, and Joseph Loscalzo. *Harrison's principles of internal medicine*, 19e, volume 1. Mcgraw-hill New York, NY, USA:, 2015.
  - Seyed Mehran Kazemi, Rishab Goel, Sepehr Eghbali, Janahan Ramanan, Jaspreet Sahota, Sanjay Thakur, Stella Wu, Cathal Smyth, Pascal Poupart, and Marcus Brubaker. Time2vec: Learning a vector representation of time. *arXiv preprint arXiv:1907.05321*, 2019.
  - Z Kraljevic, D Bean, A Shek, R Bendayan, H Hemingway, J A Yeung, A Deng, A Baston, J Ross, E Idowu, et al. Foresight-a generative pretrained transformer for modelling of patient timelines using electronic health records: a retrospective modelling study. *Lancet Digit. Health*, 6(4), 2024.
- Kwanhyung Lee, Soojeong Lee, Sangchul Hahn, Heejung Hyun, Edward Choi, Byungeun Ahn, and Joohyung Lee. Learning missing modal electronic health records with unified multi-modal data embedding and modality-aware attention. In *Machine Learning for Healthcare Conference*, pp. 423–442. PMLR, 2023.
  - Y. Li, S. Rao, J. Solares, A. Hassaine, R. Ramakrishnan, D. Canoy, Y. Zhu, K. Rahimi, and G. Salimi-Khorshidi. Behrt: Transformer for electronic health records. *Scientific Reports*, 10 (1):1–10, 2020. doi: 10.1038/s41598-020-62922-y.

- Tatiana Likhomanenko, Qiantong Xu, Gabriel Synnaeve, Ronan Collobert, and Alex Rogozhnikov. Cape: Encoding relative positions with continuous augmented positional embeddings. *Advances in Neural Information Processing Systems*, 34:16079–16092, 2021.
  - Nurul Athirah Nasarudin, Fatma Al Jasmi, Richard O Sinnott, Nazar Zaki, Hany Al Ashwal, Elfadil A Mohamed, and Mohd Saberi Mohamad. A review of deep learning models and online healthcare databases for electronic health records and their use for health prediction. *Artificial Intelligence Review*, 57(9):249, 2024.
  - H Niu, O A Omitaomu, M A Langston, M Olama, O Ozmen, H B Klasky, A Laurio, M Ward, and J Nebeker. Ehr-bert: A bert-based model for effective anomaly detection in electronic health records. *J. Biomed. Inf.*, 150:104605, 2024.
  - Adam Paszke, Sam Gross, Soumith Chintala, Gregory Chanan, Edward Yang, Zachary DeVito, Zeming Lin, Alban Desmaison, Luca Antiga, and Adam Lerer. Automatic differentiation in pytorch. In *NIPS-W*, 2017.
  - Laila Rasmy, Yiqing Xiang, Zhaoyi Xie, Cong Tao, and Daqing Zhi. Med-bert: pretrained contextualized embeddings on large-scale structured electronic health records for disease prediction. *npj Digital Medicine*, 4(1):86, 2021.
  - Pawel Renc, Yugang Jia, Anthony E. Samir, Jaroslaw Was, Quanzheng Li, David W. Bates, and Adam Sitek. Zero shot health trajectory prediction using transformer. *npj Digital Medicine*, 7(1), February 2024. ISSN 2398-6352. doi: 10.1038/s41746-024-01235-0. URL https://www.nature.com/articles/s41746-024-01235-0.
  - E Steinberg, J A Fries, Y Xu, and N Shah. Motor: A time-to-event foundation model for structured medical records. In *ICLR*, 2024.
  - Ethan Steinberg, Kevin Jung, Jason A Fries, Christopher K Corbin, Stephen R Pfohl, and Nigam H Shah. Language models are an effective representation learning technique for electronic health record data. *J. Biomed. Inf.*, 113, 2021.
  - Jianlin Su, Murtadha Ahmed, Yu Lu, Shengfeng Pan, Wen Bo, and Yunfeng Liu. Roformer: Enhanced transformer with rotary position embedding. *Neurocomputing*, 568:127063, 2024.
  - Sindhu Tipirneni and Chandan K Reddy. Self-supervised transformer for sparse and irregularly sampled multivariate clinical time-series. *ACM Transactions on Knowledge Discovery from Data* (*TKDD*), 16(6):1–17, 2022.
  - Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information processing systems*, 30, 2017.
  - Michael Wornow, Rahul Thapa, Ethan Steinberg, Jason Fries, and Nigam Shah. Ehrshot: An ehr benchmark for few-shot evaluation of foundation models. *arXiv preprint arXiv:2307.02028*, 2023. URL https://som-shahlab.github.io/ehrshot-.
  - Zhichao Yang, Avijit Mitra, Sunjae Kwon, and Hong Yu. Clinicalmamba: A generative clinical language model on longitudinal clinical notes. *arXiv preprint arXiv:2403.05795*, 2024.
  - W Zhu, H Tang, H Zhang, H R Rajamohan, S-L Huang, X Ma, A Chaudhari, D Madaan, E Almahmoud, S Chopra, et al. Predicting risk of alzheimer's diseases and related dementias with ai foundation model on electronic health records. *medRxiv preprint medRxiv:2024.04.26.24306180*, 2024.

**Appendix** 595

594

596

597 598

600

601

602

603

604 605

606 607

608

609

610

611

612

613

614

615

616

617

618

619 620

621

622

623

624

625

626 627

628

629

630

631 632

633

634

635

636

637 638 639

640

641

642

643

644 645

646

647

# ADDITIONAL RELATED WORKS

The pursuit of powerful and versatile foundation models for Electronic Health Records (EHRs) has led to a rapid evolution of modeling techniques. Current models drew heavily from advancements in Natural Language Processing (NLP), treating EHR data as sequences of discrete events. However, the unique characteristics of EHRs—specifically their sparse, irregularly sampled nature, and the continuous numerical values—have necessitated the development of more specialized architectures and learning objectives.

#### EARLY APPROACHES: MODELING CATEGORICAL EVENT SEQUENCES

The initial wave of EHR foundation models adapted the successful Transformer architecture from NLP to the clinical domain. These models primarily focused on learning representations from sequences of medical codes, such as diagnoses, procedures, and medications, while largely omitting numerical and temporal data.

BEHRT (Li et al., 2020) (Transformer for Electronic Health Records) introduced the use of the Bidirectional Encoder Representations from Transformers (BERT (Devlin et al., 2019)) architecture for EHR data. It treats a patient's EHR as a sequence of "sentences," where each sentence is a collection of medical codes from a single visit. BEHRT is pre-trained on a large dataset of patient records using a Masked Language Model (MLM) objective, where the model learns to predict masked medical codes based on their context. An additional task, Next Visit Prediction (NVP), was also used to predict the codes for a subsequent visit. While effective for tasks like disease prediction, BEHRT's exclusion of numerical values and timestamps limits its clinical utility, as it cannot capture disease severity or the precise timing of events.

Med-BERT (Rasmy et al., 2021) followed a similar approach to BEHRT, applying the BERT architecture to structured EHR data. It also represents patient histories as sequences of medical codes and uses an MLM pre-training objective to learn contextualized embeddings. Med-BERT demonstrated strong performance on various downstream tasks, including disease prediction and patient mortality prediction. However, like BEHRT, it does not explicitly model the temporal intervals between visits or the continuous values of lab tests, which are crucial for a comprehensive understanding of a patient's health trajectory.

EHR-BERT (Niu et al., 2024) is another BERT-based model that focuses on detecting anomalies in EHR data. It learns the typical patterns of medical events and flags deviations from these patterns as potential anomalies. While its primary application is in data quality and fraud detection, it shares the same fundamental limitations as other early BERT-based models in its handling of EHR data, as it does not incorporate numerical or temporal information into its core architecture.

Further including code-based models: ClinicalMamba (Yang et al., 2024), and Foresight (Kraljevic et al., 2024), these models established the viability of large-scale pre-training for EHR data and demonstrated the power of the Transformer architecture in capturing the complex relationships between medical events. However, their reliance on a purely categorical representation of patient histories highlighted the need for more sophisticated methods that could incorporate the rich numerical and temporal information present in EHRs.

#### INCORPORATING NUMERIC VALUES

Recognizing the limitations of purely categorical models, subsequent research focused on integrating continuous numerical values, such as lab results and vital signs, into the modeling process. A common approach has been to discretize these values into a set of predefined bins, allowing them to be treated as discrete tokens within the existing language modeling framework.

HEART (Huang et al., 2024) employs this discretization strategy. They convert numeric values into a fixed number of uniform bins (e.g., 10 bins) and create a unique token for each "feature-value" pair. This allows them to capture the magnitude of numerical measurements to some extent. However, this approach has two major drawbacks. First, it leads to a massive increase in the vocabulary size,

as each feature needs its own set of value-specific tokens. This exacerbates data sparsity issues and increases the model's memory footprint. Second, uniform binning can be suboptimal, as it focuses on the non-significant range of clinical values and may not provide sufficient resolution for clinically significant changes.

FM4EHR (Burkhart et al., 2025) (Foundation Models for Electronic Health Records) proposed a more efficient method for handling numeric values. Instead of creating unique tokens for each feature-value pair, FM4EHR separates the tokenization of features and values. This allows different features to share the same set of value tokens, significantly reducing the vocabulary size and mitigating the data sparsity problem. This "value sharing" approach is a key innovation that allows for more scalable and efficient modeling of numerical data. However, FM4EHR still does not explicitly model the temporal aspect of EHR data, relying on the implicit ordering of events in the sequence.

#### ADDRESSING TEMPORAL INFORMATION

The timing of medical events is often as important as the events themselves. Another line of research has focused on developing models that can explicitly receive temporal information of EHR data.

MOTOR (Steinberg et al., 2024) (A Time-to-Event Foundation Model for Structured Medical Records) is a model specifically designed for survival analysis and time-to-event prediction. It processes sequences of medical codes and learns to predict the time to a future event of interest. (Note! it does not take numerical value) MOTOR represents time by discretizing the time horizon into a set of predefined intervals and models the hazard function within each interval. This allows it to capture the temporal dependencies between events and make time-aware predictions. However, MOTOR's primary limitation is that it does not incorporate numerical values, which are often strong predictors of patient outcomes. Additionally, its reliance on predefined time intervals and the assumption of a constant hazard function within each interval can limit its temporal precision.

STraTS (Tipirneni & Reddy, 2022) (Self-Supervised Transformer for Sparse and Irregularly Sampled Multivariate Clinical Time-Series) takes a different approach to modeling time and values. It tokenizes only the categorical features and embeds the time intervals and numerical values as continuous variables. STraTS is trained using a Mean Squared Error (MSE) loss to predict the values of different features at future time points. This allows it to handle irregularly sampled data and make fine-grained predictions. However, by only predicting the next 2-hour value, STraTS loses important contextual information and cannot be used for generative modeling of entire patient trajectories.

TRADE (Zhu et al., 2024) (Predicting Risk of Alzheimer's Diseases and Related Dementias with AI Foundation Model on Electronic Health Records) also incorporates numerical values, and it uses a non-uniform binning. It discretizes values into nine bins based on their standard deviation from the mean. This approach is more clinically plausible than uniform binning, as it can better capture extreme values that are often indicative of disease. However, TRADE does not employ value sharing, which means it still faces the challenge of a large and sparse vocabulary.

*EhrMamba* (Fallahpour et al., 2025) is a recent model that leverages the Mamba architecture, a type of State Space Model (SSM), to efficiently process long EHR sequences. It tokenizes categorical features and uses uniform binning for numerical values. It uses time2vec module Kazemi et al. (2019) to capture temporal dependencies. The use of the Mamba architecture allows *EhrMamba* to scale to much longer patient histories than Transformer-based models, which have a quadratic complexity with respect to sequence length.

ETHOS Renc et al. (2024) (Zero-shot health trajectory prediction using transformer) introduces a novel method for explicitly modeling the time intervals between events. It discretizes the time gaps into 13 logarithmic bins, ranging from minutes to months, and inserts a special "time token" between each event token in the input sequence. This allows the model to explicitly reason about the temporal relationships between events. ETHOS also incorporates numerical values through binning and value sharing. While this explicit time tokening is a significant step forward, it can increase the sequence length and computational cost. Moreover, the discretization of time still imposes a limit on the model's temporal precision.

# B ADDITIONAL METHOD

#### B.1 PATHOLOGY-FOCUSED BINNING

This numerical value tokenization method is designed to create a granular representation for clinically important pathologic values. This non-parametric approach assigns greater resolution to sparse, low density value ranges without making distributional assumptions. The process consists of two main stages: (1) value weight assignment via Kernel Density Estimation, and (2) weighted percentile binning using these assigned weights.

# B.1.1 VALUE WEIGHT ASSIGNMENT VIA KERNEL DENSITY ESTIMATION

The core principle is to assign low weights to values in high-density (physiologic) regions and high weights to values in low-density (pathologic) regions. This is achieved by estimating the data density for each medical feature and assigning a weight inversely proportional to this density.

For a feature with a set of values V and standard deviation  $\sigma$ , we define a set of discrete representative points  $X = \{x_1, x_2, \dots, x_M\}$ . These points span the feature's range  $[\min(V), \max(V)]$  and are spaced at uniform intervals of  $0.05\sigma$ .

At each discrete point  $x_k \in X$ , we estimate the data density  $\rho(x_k)$  by applying a Gaussian convolution kernel. This is a form of Kernel Density Estimation (KDE), where the density at  $x_k$  is the sum of influences from all unique data values  $v_j \in V$ . The density is:

$$\rho(x_k) = \sum_{j=1}^{|V|} K_h(x_k - v_j)$$

Here,  $K_h(u)$  is an unnormalized Gaussian kernel defined as:

$$K_h(u) = \exp\left(-\frac{u^2}{2h^2}\right)$$

The bandwidth h, which controls the smoothness of the density estimate, is set to  $h=0.1\sigma$  to capture local variations.

From this density, we calculate a raw weight  $w_{\text{raw}}(x_k) = 1/(\rho(x_k) + \epsilon)$  for each discrete point. These weights are then normalized and clipped to produce the final weight:

$$w_{\text{final}}(x_k) = \min\left(\frac{w_{\text{raw}}(x_k)}{\min_{k'} w_{\text{raw}}(x_{k'})}, w_{\text{max}}\right)$$

where  $w_{\text{max}}$  is a predefined ceiling (e.g., 10). Finally, each original unique value  $v_j$  is assigned the weight of its nearest discrete point,  $w(v_j) = w_{\text{final}}(x_{k^*})$ , where  $k^* = \arg\min_k |v_j - x_k|$ .

#### B.1.2 WEIGHT CALCULATION AND NORMALIZATION

With weights assigned, we proceed to the binning stage. The goal is to partition the feature's values into B bins such that regions with higher weights are given more bins.

We start by calculating a **weighted count**  $c'_i$  for each unique value  $v_j$ :

$$c_i' = c_i \cdot w(v_i)$$

This new count reflects the value's clinical importance as determined by its rarity. Next, we compute the total weighted count for the feature,  $C'_{\text{total}} = \sum_{j=1}^{|V|} c'_j$ .

The bin thresholds are then determined from the cumulative distribution of these weighted counts. For a set of sorted unique values  $v_1 < v_2 < \cdots < v_{|V|}$ , the cumulative weighted count up to value  $v_k$  is  $S_k = \sum_{j=1}^k c_j'$ . The threshold for the p-th bin (where  $p \in \{1, 2, \dots, B-1\}$ ) is set to the first value  $v_k$  whose cumulative weighted share meets or exceeds the p/B percentile:

$$T_p = \min\{v_k \mid \frac{S_k}{C'_{\text{total}}} \ge \frac{p}{B}\}$$

757

758

759

This procedure ensures that value ranges containing high-weight (pathologic) data contribute more significantly to the cumulative sum. As a result, a smaller span of these values is needed to cross a percentile boundary, leading to a denser allocation of bin thresholds in these clinically important regions.

# Algorithm 1 Density-Based Value Weight Assignment

```
777
         Require:
778
              item_counters: A map from item ID \rightarrow {value: count}.
779
              w_{max}: Maximum weight threshold (e.g., 10).
780
         Ensure:
781
              value_weights: A map from item ID \rightarrow {value: weight}.
782
           1: Initialize value_weights \leftarrow \emptyset
783
           2: for all item_id, counter in item_counters do
784
                  V \leftarrow \text{sorted unique values from counter}
           3:
785
           4:
                  C \leftarrow corresponding counts for each value in V
           5:
                  AllValues \leftarrow list of all values repeated by their counts
786
           6:
                  \sigma \leftarrow \text{StandardDeviation}(AllValues)
787
           7:
                  if \sigma = 0 then
788
                       item_weights \leftarrow \{v: 1.0 \text{ for } v \in V\}
           8:
789
           9:
                       value_weights[item_id] ← item_weights
790
         10:
                      continue
791
                                                                            > Set bandwidth for the Gaussian kernel
                  h \leftarrow 0.1 \times \sigma
         11:
792
                  interval \leftarrow 0.05 \times \sigma
         12:
793
                  SplitPoints \leftarrow Generate points from min(V) to max(V) with interval
         13:
794
                                                                                795
         14:
                  densities \leftarrow \emptyset
                  for all split point x in SplitPoints do
796
         15:
                      \rho(x) \leftarrow \frac{1}{|V|} \sum_{j=1}^{|V|} C_j \cdot \exp\left(-\frac{(x-V_j)^2}{2h^2}\right)
797
         16:
798
                       densities.append(\rho(x))
         17:
799
                                                                 ▶ Step 2: Calculate, Normalize, and Clip Weights
800
                  w_{raw} \leftarrow 1.0/(\text{densities} + 10^{-10})
         18:
801
                  w_{norm} \leftarrow w_{raw} / \min(w_{raw})
         19:
802
         20:
                  w_{final} \leftarrow \text{clip}(w_{norm}, 1.0, w_{max})
803
                                                                         ▶ Step 3: Assign weights to original values
804
         21:
                  item_weights \leftarrow \emptyset
805
         22:
                  for all value v in V do
806
         23:
                       closest\_idx \leftarrow arg min_k |SplitPoints_k - v|
807
         24:
                       item_weights[v] \leftarrow w_{final}[closest\_idx]
808
         25:
                  value_weights[item_id] ← item_weights
809
         26: return value_weights
```

#### 810 Algorithm 2 Weighted Percentile Binning 811 Require: 812 item\_counters: A map from item ID $\rightarrow$ {value: count}. 813 value\_weights: The output from Algorithm 1. 814 B: The desired number of bins (e.g., 100). 815 816 bin\_thresholds: A map from item ID $\rightarrow$ a list of B-1 thresholds. 817 1: **Initialize** bin\_thresholds $\leftarrow \emptyset$ 2: for all item\_id, counter in item\_counters do 818 $V \leftarrow \text{sorted unique values from counter}$ 819 3: 4: 820 ▶ Apply weights to counts for percentile calculation 821 5: if value\_weights is provided then 822 $C' \leftarrow [\text{counter}[v_j] \times \text{value\_weights}[\text{item\_id}][v_j] \text{ for } v_j \in V]$ 6: 823 7: □ Uniform binning case 824 8: $C' \leftarrow [\texttt{counter}[v_i] \text{ for } v_i \in V]$ 825 thresholds $\leftarrow$ a list of size B-19: 10: if N > B then 11: 828 12: 829 **for** p = 1 to B - 1 **do** 13: 830 $target\_count \leftarrow C'_{total} \times p/B$ 14: $idx \leftarrow \mathsf{FindFirstIndexWhere}(C'_{cumulative} > target\_count)$ 831 15: 832 16: thresholds $[p-1] \leftarrow V[idx+1]$ 833 17: else ▶ Apply specific assignment for sparse values (centering or striding) 834 18: thresholds ← Generate thresholds based on sparse assignment logic 835 19: bin\_thresholds[item\_id] ← thresholds 836 20: return bin\_thresholds

#### B.2 DUAL-CALENDAR ROTARY POSITION EMBEDDING

837 838

839 840

841

842

843

844

845

846

847

848

849

850

851 852

853 854

855

856

858

859

860

861 862

863

To address the unique temporal characteristics of Electronic Health Record (EHR) data—namely, the highly irregular event intervals and the clinical significance of calendrical time—we propose **Dual-Calendar Rotary Position Embedding (RoPE)**. This method extends the conventional RoPE by partitioning the embedding dimension within each attention head to jointly encode both the **relative sequence order** of tokens and their **absolute calendrical time**.

For a given query or key vector  $x \in \mathbb{R}^{d_k}$  in an attention head, we partition it into two subspaces: a positional component  $x_{pos} \in \mathbb{R}^{d_{pos}}$  and a temporal component  $x_{time} \in \mathbb{R}^{d_{time}}$ , where  $d_k = d_{pos} + d_{time}$ .

$$x = [x_{pos} \parallel x_{time}]$$

Each component is then rotated using a specialized RoPE variant before being concatenated back together.

#### **B.2.1** Positional Dimension Encoding

The  $x_{pos}$  component, corresponding to the first  $d_{pos}$  dimensions, employs the standard RoPE formulation. With a reduced dimensionality of  $d_{pos}$ , this component does not rescale its rotational frequencies to cover a wide positional range. Instead, it effectively **truncates the frequency spectrum**, retaining the high-frequency rotations corresponding to the initial dimensions of a standard RoPE. This strategic choice is predicated on the observation that its primary role is now to disambiguate the order of co-occurring events that share an identical timestamp. For a token at position m, the rotation angle  $\theta_{p,i}^{(pos)}$  is defined as:

$$\theta_{p,i}^{(pos)} = \frac{p}{\text{base}^{2i/d}} \;\;,\;\; i \in \{1,2,...,d_{pos}/2\}$$

The task of modeling long-range temporal dependencies is thus naturally offloaded to the Calendar-Time dimension, which is explicitly designed for this purpose.

#### B.2.2 CALENDAR-TIME DIMENSION ENCODING

The core novelty of our approach lies in the encoding applied to the  $x_{time}$  component. This component is designed to encode an event's absolute timestamp, t, by capturing its periodicity across multiple, clinically relevant time scales. This is achieved using a predefined set of **semantically meaningful calendrical periods**,  $S = \{s_1, s_2, \ldots, s_{d_{time}/2}\}$ , as detailed in Table 5. These periods capture periodicities ranging from short-term diurnal patterns to long-term annual and multi-year trends.

For each period  $s_j$  from this set, we calculate a unique rotation angle  $\theta_{t,j}^{(time)}$  that represents the phase of the event within that specific period. The formula for the rotation angle is:

$$\theta_{t,j}^{(time)} = \left(\frac{t \pmod {s_j}}{s_j}\right) \cdot 2\pi$$

This mechanism produces a **multi-scale temporal representation**. Events occurring at the same time of day but on different dates will share the exact same rotation for the 'day' period, allowing the model to easily learn periodical patterns.

# **B.2.3** Integration and Application

Finally, the two rotated components are concatenated to form the final query and key vectors. The full transformation for a query vector  $q = [q_{pos} \parallel q_{time}]$  to its rotated form q' is:

$$q' = [\text{RoPE}(q_{pos}, \theta^{(pos)}) \parallel \text{RoPE}(q_{time}, \theta^{(time)})]$$

An identical transformation is applied to the key vector k. By equipping the self-attention mechanism with this dual-encoding strategy, our model can simultaneously reason about the sequential flow of information and the absolute, cyclical context of clinical events.

Table 5: Predefined Calendrical Periods for Temporal Encoding

Category	Period Name	<b>Duration (seconds)</b>
Short-term	5 minutes	300
	10 minutes	600
	30 minutes	1,800
	1 hour	3,600
	3 hours	10,800
	12 hours	43,200
Mid-term	1 day	86,400
	2 days	172,800
	1 week	604,800
	2 weeks	1,209,600
	1 month	2,629,746
	1 season (3 months)	7,889,238
	6 months	15,778,476
Long-term	1 year	31,556,952
	2 years	63,113,904
	4 years	126,227,808
	10 years	315,569,520
	30 years	946,708,560
	100 years	3,155,695,200
	300 years	9,467,085,600

#### B.3 TIME CONDITIONED FORESEE OBJECTIVE

To elucidate the mechanics of our proposed Time-Conditioned Foresee (TCF) module, we provide a step-by-step explanation. This description follows the flow of information from the initial input—the

final hidden state of a backbone model—to the module's dual outputs. The detailed computational flow is presented in Algorithm 3.

We define the primary dimensions:

- B denotes the batch size.
- L denotes the sequence length.
- $d_{\text{model}}$  represents the hidden dimension of the backbone model's output.
- $d_{\text{embed}}$  is the dimensionality of our internal time embeddings, set to 32.
- $N_{\text{scales}}$  is the number of time scales used for decomposition, 11 in our implementation.
- $N_{\text{foresee}}$  is the number of future timestamps provided for the foresee objective, set to 10.
- $C_i$  is the number of discrete categories for the *i*-th time scale.

The process begins with the final hidden state from the backbone decoder for each token in the sequence.

**Input:** The last hidden state tensor,  $H_{\text{last}}$ .

- $H_{\text{last}} \in \mathbb{R}^{B \times L \times d_{\text{model}}}$
- B.3.1 A. Next Event Time Prediction Path ( $\mathcal{L}_{\text{Next\_time}}$ )

This pathway is responsible for predicting the time until the next event.

- 1. Initial Projection The input hidden state  $H_{\text{last}}$  is passed through a two-layer Feed-Forward Network (FFN), denoted as FFN<sub>enc</sub>, to create a representation for time prediction.
  - Input:  $H_{\text{last}} \in \mathbb{R}^{B \times L \times d_{\text{model}}}$
  - Output: An intermediate time-focused tensor,  $H_{\text{time}}$ .
    - $H_{\text{time}} = \text{FFN}_{\text{enc}}(H_{\text{last}}) \in \mathbb{R}^{B \times L \times (N_{\text{scales}} \cdot d_{\text{embed}})}$
- 2. Logit Generation The tensor  $H_{\text{time}}$  is conceptually partitioned into  $N_{\text{scales}}$  segments. Each segment is used to compute the logits for its corresponding time scale by multiplying it with the respective time embedding weight matrix.
  - Input:  $H_{\text{time}}$ , treated as  $N_{\text{scales}}$  tensors  $\{H_{\text{time}}^{(i)}\}_{i=1}^{N_{\text{scales}}}$ , where each  $H_{\text{time}}^{(i)} \in \mathbb{R}^{B \times L \times d_{\text{embed}}}$ .
  - **Operation:** For each scale i, we compute logits:  $\text{Logits}^{(i)} = H_{\text{time}}^{(i)} \cdot (W_{\text{embed}}^{(i)})^T$ , where  $W_{\text{embed}}^{(i)} \in \mathbb{R}^{C_i \times d_{\text{embed}}}$ .
  - Output: A set of  $N_{\text{scales}}$  logit tensors.
    - Logits<sup>(i)</sup>  $\in \mathbb{R}^{B \times L \times C_i}$
- **3. Ground-Truth Label Decomposition** To compute the loss, these logits are compared against ground-truth labels. Instead of regressing a continuous time value, we transform the ground-truth time delta (in seconds) into a set of categorical integer labels. This is achieved through a deterministic process analogous to a mixed-radix conversion, using the time scales defined in Table 6.

For instance, assume a ground-truth time delta  $\Delta T_{\rm next}$  of **34,586,130 seconds**. The conversion to a vector of  $N_{\rm scales}$  integer labels proceeds sequentially from the largest time scale to the smallest, using integer division to find the label and the modulo operator to find the remainder for the next step.

```
1. year10: 34,586,130 // 315,360,000 = 0. Remainder: 34,586,130. \rightarrow Label: 0
```

- 2. **year1:** 34,586,130 // 31,536,000 = 1. Remainder:  $3,050,130. \rightarrow$  **Label: 1**
- 3. month3: 3,050,130 //7,948,800 = 0. Remainder:  $3,050,130. \rightarrow$  Label: 0
- 4. month1: 3,050,130 // 2,678,400 = 1. Remainder:  $371,730. \rightarrow$  Label: 1

Table 6: Time scales for multi-scale decomposition, along with their duration in seconds and the number of categories for classification.

Time Scale	<b>Duration in Seconds</b>	Num. of Categories
year10	315,360,000	10
year1	31,536,000	10
month3	7,948,800	4
month1	2,678,400	3
week1	604,800	5
day1	86,400	7
hour6	21,600	4
hour1	3,600	6
minute10	600	6
minute1	60	10

- 5. **week1:** 371,730 //604,800 = 0. Remainder:  $371,730. \rightarrow$  **Label: 0**
- 6. day1: 371,730 // 86,400 = 4. Remainder:  $27,330. \rightarrow$  Label: 4
- 7. hour6: 27,330 // 21,600 = 1. Remainder:  $5,730. \rightarrow$  Label: 1
- 8. hour1:  $5{,}730 // 3{,}600 = 1$ . Remainder:  $2{,}130. \rightarrow Label: 1$
- 9. minute10: 2,130 //600 = 3. Remainder:  $330. \rightarrow Label: 3$
- 10. minute1: 330 //60 = 5. Remainder: 30.  $\rightarrow$  Label: 5
- 11. **position:** A mechanism to account for events occurring simultaneously at the same time. For the next event label, it is set to 0; for subsequent foresee labels, it is set to +1 if the time is the same as the previous one, and 0 otherwise. → **Label: 0**

Ultimately, the continuous value of 34,586,130 seconds is converted into the following vector of ten integer labels, which constitutes the ground-truth  $Y_{\text{next time}}^{(i)}$  for this example:

By training the model to predict these categorical labels for each time scale, we transform a difficult regression task into a series of more stable and effective classification tasks. The final scalar loss,  $\mathcal{L}_{\text{next.time}}$ , is the average Cross-Entropy loss calculated between the generated logits and these decomposed ground-truth labels.

#### B.3.2 B. TIME-CONDITIONING PATH FOR FORESEE OBJECTIVE

This pathway conditions the hidden state on a set of specified future timestamps to predict upcoming events.

- 1. Input Foresee Timestamps The module receives future time deltas from the  $N_{\text{foresee}}$  future events relative to the current timestamp at each position.
  - Input: A tensor of future time deltas,  $\Delta T_{\text{foresee}} \in \mathbb{Z}^{B \times L \times N_{\text{foresee}}}$ .
- **2. Time Embedding** Each time delta in  $\Delta T_{\text{foresee}}$  is decomposed into  $N_{\text{scales}}$  integer labels (as demonstrated above). These labels are used to look up corresponding vectors from the embedding tables,  $\{W_{\text{embed}}^{(i)}\}_{i=1}^{N_{\text{scales}}}$ , which are then concatenated. \*Note, we share the weights for embedding and unembedding timestamps.
  - Input: Decomposed time labels,  $C_{\text{foresee}} \in \mathbb{Z}^{B \times L \times N_{\text{foresee}} \times N_{\text{scales}}}$ .
  - Output: A dense time embedding tensor,  $E_{\text{time}}$ .
    - $E_{\text{time}} \in \mathbb{R}^{B \times L \times N_{\text{foresee}} \times (N_{\text{scales}} \cdot d_{\text{embed}})}$

- **3. Projection of Time Embedding** The concatenated embedding  $E_{\text{time}}$  is projected to the model's hidden dimension via FFN<sub>dec</sub>. • Input:  $E_{\text{time}} \in \mathbb{R}^{B \times L \times N_{\text{foresee}} \times (N_{\text{scales}} \cdot d_{\text{embed}})}$ • Output: A processed time conditioning tensor,  $H_{\text{time\_cond}}$ . -  $H_{\text{time\_cond}} = \text{FFN}_{\text{dec}}(E_{\text{time}}) \in \mathbb{R}^{B \times L \times N_{\text{foresee}} \times d_{\text{model}}}$ **4. Time-Conditioning via Fusion** The final step fuses the original hidden state  $H_{\text{last}}$  with the pro-cessed time conditioning tensor  $H_{\text{time\_cond}}$ . To align their dimensions for the element-wise addition,  $H_{\text{last}}$  is first expanded by inserting a new dimension. This prepares it for broadcasting across the  $N_{\text{foresee}}$  dimension, allowing each of the  $N_{\text{foresee}}$  time embeddings to condition the single original hidden state. • Inputs: - Original hidden state:  $H_{\text{last}} \in \mathbb{R}^{B \times L \times d_{\text{model}}}$ – Time conditioning tensor:  $H_{\text{time\_cond}} \in \mathbb{R}^{B \times L \times N_{\text{foresee}} \times d_{\text{model}}}$ 
  - Operation: The fusion is performed via a residual connection. First,  $H_{\text{last}}$  is unsqueezed, and then added to  $H_{\text{time\_cond}}$ .
    - $H'_{\text{last}} = \text{Unsqueeze}(H_{\text{last}}, \text{dim} = 2) \in \mathbb{R}^{B \times L \times 1 \times d_{\text{model}}}$
    - $H_{\text{fused}} = H'_{\text{last}} + H_{\text{time\_cond}}$  // Broadcasting occurs along the  $N_{\text{foresee}}$  dimension.
  - Output: The final time-conditioned hidden state,  $H_{\rm conditioned}$ , after Layer Normalization and a final FFN block.
    - $H_{\text{conditioned}} \in \mathbb{R}^{B \times L \times N_{\text{foresee}} \times d_{\text{model}}}$

SUMMARY OF MODULE OUTPUTS

The TCF module produces two primary outputs:

- 1. **Next Time Loss** ( $\mathcal{L}_{next.time}$ ): A scalar value for backpropagation.
- 2. Conditioned Hidden State ( $H_{\text{conditioned}}$ ): A tensor of shape  $\mathbb{R}^{B \times L \times N_{\text{foresee}} \times d_{\text{model}}}$ , which serves as the input to the final prediction head for calculating the foresee loss,  $\mathcal{L}_{\text{foresee}}$ . The ground-truth label corresponding to each conditioned hidden vector is the Feature token of the actual clinical event that occurred at the given foresee timestamp.

```
1080
               Algorithm 3 Time-Conditioned Foresee (TCF) Module
1081
1082
                 1: H_{\text{last}} \in \mathbb{R}^{B \times L \times d_{\text{model}}}: Last hidden states from the backbone model.
                 2: T_{\text{current}} \in \mathbb{R}^{B \times L}: Absolute timestamps for each hidden state in H_{\text{last}}.
1084
                 3: T_{\text{next}} \in \mathbb{R}^{B \times L}: Absolute timestamps of the next event for each position.
                 4: T_{\text{foresee}} \in \mathbb{R}^{B \times L \times N_{\text{foresee}}}: A set of absolute future timestamps for conditioning.
                 5: W: All trainable weights, including FFNs and embedding tables \{W_{\text{embed}}^{(i)}\}_{i=1}^{N_{\text{scales}}}.
1087
                 6: Periods: A dictionary mapping each time scale to its duration in seconds.
1088
               Ensure:
1089
                 7: \mathcal{L}_{next\_time} \in \mathbb{R}: The loss for next event time prediction.
                 8: H_{\text{conditioned}} \in \mathbb{R}^{B \times L \times N_{\text{foresee}} \times d_{\text{model}}}: Hidden states conditioned on T_{\text{foresee}}.
1090
                 9: function TCF_MODULE(H_{last}, T_{current}, T_{next}, T_{foresee}, W, Periods)
                                                                                                                       > Part A: Next Event Time Prediction
1093
                             \Delta T_{\text{next}} \leftarrow T_{\text{next}} - T_{\text{current}}
               10:
1094
                             Y_{\text{next}} \leftarrow \text{DECOMPOSETIME}(\Delta T_{\text{next}}, Periods)
               11:
1095
                             H_{\text{time}} \leftarrow \text{FFN}_{\text{enc}}(H_{\text{last}})
                                                                                                                                      \triangleright Shape: (B, L, N_{\text{scales}} \cdot d_{\text{embed}})
               12:
                             H_{\text{time}} \leftarrow \text{RESHAPE}(H_{\text{time}}, (B, L, N_{\text{scales}}, d_{\text{embed}}))
               13:
               14:
                             \mathcal{L}_{\text{total}} \leftarrow 0
                             for i = 1 \rightarrow N_{scales} do
               15:
                                   H_{\text{time}}^{(i)} \leftarrow H_{\text{time}}[:,:,i,:]
\text{Logits}^{(i)} \leftarrow H_{\text{time}}^{(i)} \cdot (W_{\text{embed}}^{(i)})^T
1099
               16:
                                                                                                                                                     \triangleright Shape: (B, L, d_{embed})
1100
               17:
1101
                                   Y_{\text{next}}^{(i)} \leftarrow Y_{\text{next}}[:,:,i]
               18:
1102
                                   \mathcal{L}_{\text{total}} \leftarrow \mathcal{L}_{\text{total}} + \text{CROSSENTROPYLOSS}(\text{Logits}^{(i)}, Y_{\text{next}}^{(i)})
1103
               19:
1104
               20:
                             \mathcal{L}_{\text{next\_time}} \leftarrow \mathcal{L}_{\text{total}}/N_{\text{scales}}
                                                                                               ▶ Part B: Time-Conditioning for Foresee Objective
1105
                             T'_{\text{current}} \leftarrow \text{UNSQUEEZE}(T_{\text{current}}, \text{dim} = 2)
                                                                                                                                                             \triangleright Shape: (B, L, 1)
               21:
1106
                             \Delta T_{\text{foresee}} \leftarrow T_{\text{foresee}} - T'_{\text{current}}
               22:
1107
                             C_{\text{foresee}} \leftarrow \text{DECOMPOSETIME}(\Delta T_{\text{foresee}}, Periods)
                                                                                                                                      \triangleright Shape: (B, L, N_{\text{foresee}}, N_{\text{scales}})
               23:
1108
               24:
                             E_{\text{time\_list}} \leftarrow []
1109
                             for i = 1 \rightarrow N_{\text{scales}} do
               25:
1110
                                    \begin{array}{l} C = 1 & \text{PN}_{\text{scales}} \text{ } \textbf{uo} \\ C_{\text{foresee}}^{(i)} \leftarrow C_{\text{foresee}}[:,:,:,i] \\ E_{\text{time}}^{(i)} \leftarrow \text{Lookup}(W_{\text{embed}}^{(i)},C_{\text{foresee}}^{(i)}) \end{array} 
               26:
1111
               27:
                                                                                                                                      \triangleright Shape: (B, L, N_{\text{foresee}}, d_{\text{embed}})
1112
                                   Append E_{\text{time}}^{(i)} to E_{\text{time\_list}}
1113
               28:
1114
                                                                                                                       \triangleright Shape: (B, L, N_{\text{foresee}}, N_{\text{scales}} \cdot d_{\text{embed}})
               29:
                             E_{\text{time}} \leftarrow \text{Concatenate}(E_{\text{time\_list}}, \text{dim} = -1)
1115
                                                                                                                                       \triangleright Shape: (B, L, N_{\text{foresee}}, d_{\text{model}})
               30:
                             H_{\text{time\_cond}} \leftarrow \text{FFN}_{\text{dec}}(E_{\text{time}})
1116
                             H'_{\text{last}} \leftarrow \text{UNSQUEEZE}(H_{\text{last}}, \text{dim} = 2)
                                                                                                                                                 \triangleright Shape: (B, L, 1, d_{\text{model}})
               31:
                             H_{\text{fused}} \leftarrow \text{LAYERNORM}(H_{\text{last}}' + H_{\text{time\_cond}})
1117
               32:
               33:
                             H_{\text{conditioned}} \leftarrow \text{LAYERNORM}(H_{\text{fused}} + \text{FFN}_{\text{final}}(H_{\text{fused}}))
1118
               34:
                             return \mathcal{L}_{\text{next\_time}}, H_{\text{conditioned}}
1119
1120
               35: function DECOMPOSETIME(\Delta T, Periods)
                                                                                                                     ▶ Helper function for time decomposition
1121
               36:
                             R \leftarrow \Delta T
1122
                             Labels \leftarrow []
               37:
1123
                             {f for}\ {\it scale}\ {\it in}\ {\it REVERSED}(Periods.{\it keys}())\ {f do}
               38:
1124
               39:
                                    L_{\text{scale}} \leftarrow R // Periods[\text{scale}]
                                                                                                                                                               ▶ Integer division
1125
                                    R \leftarrow R \% \ Periods[scale]
               40:
                                                                                                                                                            Modulo operation
1126
               41:
                                    L_{\text{scale}} \leftarrow \text{CLAMP}(L_{\text{scale}}, \min = 0, \max = C_{\text{scale}} - 1)
1127
               42:
                                    Prepend L_{\text{scale}} to Labels
1128
               43:
                             return STACK(Labels)
1129
```

# C DATA AND PREPROCESSING

#### C.1 DATASET: MIMIC-III

MIMIC-III (Johnson et al., 2016a;b) is an openly accessible resource that contains de-identified clinical data from more than 40,000 individuals admitted to the intensive care units of Beth Israel Deaconess Medical Center between 2001 and 2012. The dataset encompasses a wide range of information, including patient demographics, hourly vital sign recordings, laboratory measurements, administered treatments and procedures, prescribed medications, clinical notes, radiology reports, and outcomes such as in-hospital and post-discharge mortality.

The MIMIC-III database was de-identified in compliance with Health Insurance Portability and Accountability Act (HIPAA) standards through data cleansing and systematic date shifting. To preserve clinical intervals, patient-specific dates were consistently shifted into the future by a random offset, placing admissions within the years 2100–2200 while retaining the original time of day, weekday, and approximate seasonality. For patients older than 89, dates of birth were modified such that their recorded ages exceed 300 years, thereby masking their true age in accordance with HIPAA requirements. This modification provides a suitable framework for our Dual-Calendar RoPE, which is designed to address calendrical time, to operate effectively.

#### C.2 FEATURE SELECTION

Electronic Health Records (EHRs) are rich with events that have a numerical Value, a characteristic that distinguishes them from natural language. Consequently, our experiments focused on events that possess a numerical Value. We utilized 17 vital sign features from the CHARTEVENTS.csv file, following the selection in Harutyunyan et al. (2019), and the top 100 most frequently measured laboratory tests from LABEVENTS.csv as the events for our study. These features are detailed in Table 7. In addition, patient events such as "hospital admission", "ICU transfer (ICU in)", "ICU discharge (ICU out)", and "hospital discharge" were utilized.

#### C.3 FURTHER PREPROCESSING

In addition to the preprocessing of Harutyunyan et al. (2019), we made the following modifications: (1) removed outlier values in the laboratory data based on independent evaluations by three physicians and standardized the measurement units; (2) excluded hospitalization episodes with fewer than 10 events; and (3) added an anchor token with a timestamp of January 1st, 00:00 of the same year before each admission token to serve as a calendrical time reference. As a result, the lengths of patients' medical histories follow the distribution shown in Figure 7.

#### D BACKBONE, BASELINES, PRE-TRAINING DETAIL

# D.1 BACKBONE ARCHITECTURE

The backbone model used in this study follows a standard **decoder-only transformer** architecture. To minimize performance variations caused by differences in backbone models and to quantitatively assess the effectiveness of our proposed training methodology, we used the same backbone across all experiments. However, for models trained with Transformer encoders using the masked language modeling approach (HEART, TRADE), we removed the causal mask during pre-training so that they could be used as encoder models. The backbone details are as follows:

- Vocabulary Size: 166 (1219 if not share bin)
- Embedding and Hidden Dimension ( $d_{model}$ ): 512
- Number of Decoder Layers (N): 6
- Number of Attention Heads: 8
- Dimension per Head: 64
- Dimension of K, Q: 64 (Ours: first 24: positional RoPE / last 40: calendrical time RoPE)

1188 Table 7: Selected features from MIMIC-III used in this study. 1189 1190 **Features** Category 1191 Ethnicity - # 10 1192 White, White - Russian, White - other European, Asian, Asian - Chinese, (from ADMISSIONS.csv) Hispanic or Latino, Hispanic/Latino - Dominican, Black/Cape Verdean, 1193 Black/African American, Others or Unknown 1194 1195 Vital Signs - # 17 1196 Capillary refill rate (CRR), Systolic blood pressure (SBP), Mean blood (from CHARTEVENTS.csv) 1197 pressure (MBP), Diastolic blood pressure (DBP), Fraction of inspired oxygen (FiO2), Heart rate (HR), Respiratory rate (RR), Glasgow coma 1198 scale eye response (GCS-E), Glasgow coma scale motor response (GCS-1199 M), Glasgow coma scale verbal response (GCS-V), Glasgow coma scale (GCS), Serum glucose, O2 saturation, Blood pH, Body temperature, 1201 Height, Weight. 1202 1203 **Laboratory Tests** - # 100 Hematocrit, Potassium, Sodium, Creatinine, Chloride, Blood urea nitrogen, 1204 (from LABEVENTS.csv) Bicarbonate, Platelets, Anion gap, White blood cell count, Hemoglobin 1205 chemistry, Mean corpuscular hemoglobin concentration, Red blood cell count, Mean corpuscular hemoglobin, Mean corpuscular volume, Red Cell 1207 Distribution Width, Magnesium, Calcium Total, Phosphate, Base excess, 1208 CO2 (ETCO2, PCO2, etc.), Partial pressure of oxygen, Partial pressure of 1209 carbon dioxide, Partial thromboplastin time, Prothrombin time INR, Pro-1210 thrombin time, Calcium Free, Bilirubin Total, Alanine aminotransferase, 1211 Asparate aminotransferase, Alkaline phosphate, Potassium blood gas, Lac-1212 tate, Lymphocytes, Neutrophils, Monocytes, Eosinophils, Basophils, Albumin, Creatine Kinase, Oxygen blood gas, Urine Specific Gravity, Creatine 1213 Kinase-MB, Lactate dehydrogenase, Urine Protein, Urine Urobilinogen, 1214 Urine Ketone, Urine Color, Urine Appearance, Urine Blood, Urine Biliru-1215 bin, Urine Nitrite, Urine Leukocyte, Hematocrit blood gas, Hemoglobin 1216 blood gas, Troponin-T, Positive end-expiratory pressure, Urine Yeast, 1217 Urine White blood cell count, Urine Red blood cell count, Urine Epithelial 1218 cells, Band Neutrophils, Urine Bacteria, Sodium blood gas, Lipase, Amy-1219 lase, Estimated GFR, Hypochromia, Anisocytosis, Macrocytosis, Lympho-1220 cytes Atypical, Metamyelocytes, Myelocytes, Microcytes, Poikilocytosis, Vancomycin (blood), Chloride blood gas, Polychromasia, Functional Fibrinogen, Bilirubin Direct, Bilirubin Indirect, Platelet Smear, Urine Creatinine, Thyroid-stimulating hormone, Urine Sodium, Triglycerides, Granulocyte count, CK-MB Index, Phenytoin (blood), Alveolar-arterial gradi-1224 ent, Cholesterol Total, Urine osmolality, Osmolality, Uric acid, Choles-1225 terol HDL, Iron, Cholesterol ratio Total/HDL, Ferritin, Transferrin, Iron 1226 binding capacity, HbA1C, Nucleated red cells, Cholesterol, Ovalocyte, 1227 Urine Hyaline casts, Urine mucous, Cortisol, Urine urea nitrogen, Hap-1228 toglobin, Protein (Total), Vitamin B12, Benzodiazepine Screen, Barbi-1229 turate Screen, Tricyclic Antidepressant Screen, Troponin-I, Urine potas-1230 sium, Tacrolimus level, Schistocytes, Reticulocyte count, Ethanol, Urine 1231 Chloride, Acetaminophen, Urine Cocaine, Urine Benzodiazepine screen, 1232 Urine Amphetamine screen, Urine Opiate screen, Urine Barbiturate screen, 1233 Urine Methadone, Bicarbonate blood gas, Salicylate, Urine Total protein, Teardrop cells, Cyclosporin, Folate, Burr cells, Sedimentation rate, Digoxin, Thyroxine, Globulin, Urine protein/creatine ratio, NT-proBNP, Urine Amorphous cristal, C-reactive protein, Large platelets, Urine Granu-1236 lar casts, Gentamicin, Target cell, Transitional epithelial cells, Fibrin degra-1237 dation, CSF Lymphs.

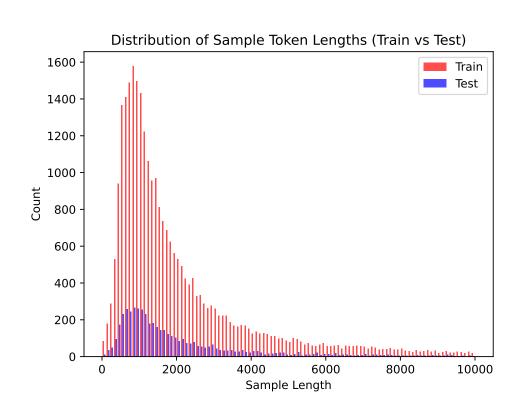


Figure 7: Length distribution of tokenized histories per patient. Most are within 2048 tokens. Lengths exceeding 10k tokens are not shown.

- Feed-Forward Network (FFN) Inner Dimension (dff): 2048
- **Total Parameters (Backbone)**: 19,001,344 (19,060,736 if not share bin)
- Activation: ReLU (Agarap, 2018)

Each decoder layer is composed of two main sub-layers: a **multi-head self-attention** block and a **feed-forward network**. For training stability, the model adopts a **Pre-Layer Normalization (Pre-LN)** structure, where Layer Normalization is applied to the input of each sub-layer. A **residual connection** is then employed around each of the two sub-layers.

#### D.2 BASELINE MODELS

We compared our model with baselines that utilized either Time or Value in their training: HEART, FM4EHR, MOTOR, STraTS, EHRSHOT, TRADE, EHRMamba, and ETHOS. The backbone for each model was standardized as described above, while other methodologies (binning, tokenization, positional embedding, learning objective, etc.) followed their original papers. The reproduction details are as follows.

#### **HEART** (Total Parameters: 20,204,180)

- Backbone: Uses a Transformer encoder due to its MLM-based loss.
- Bin Sharing: No, uses Value-Feature paired tokenization.
- **Binning**: 10-uniform binning.
- **Positional Embedding**: Absolute positional embedding—a learned positional embedding of the visit index (0 for the patient's first visit, 1 for the second, and so on).

each position was altered.

FM4EHR (Total Parameters: 19,001,536)

• Backbone: Transformer decoder.

tokens across different features.

• Learning Objective: Uses NTP loss.

• **Binning**: 10-uniform binning.

explicit time information.

1296

1297

1298

1299

1300

1301 1302

1303 1304

1305

1306

1307 1308

1309

1310

1311

1312

1313	MOTOR (Total Parameters: 20,341,952)
1315	Backbone: Transformer decoder.
1316	• Methodology Adaptation: This model was originally designed for feature-only events
1317	(e.g., diagnosis codes). To adapt it for continuous values, we set the "occurrence of an
1318	abnormal measurement result" as the endpoint for its time-to-event loss. An abnormal value
1319	was defined by either being outside the medical normal range (per Harrison's Principles of
1320	Internal Medicine (Kasper et al., 2015), evaluated by three clinicians) or being a statistical
1321	outlier (e.g., top/bottom 5%).
1322	• Bin Sharing: No, uses Value-Feature paired tokenization.
1323	• Binning: 10-uniform binning was applied to the Value.
1324	• <b>Positional Embedding</b> : The model converts each event's timestamp into 'days since birth' and applies this value in its rotary positional embedding.
1326	• Learning Objective: Uses the Time-to-event loss from the original paper.
1327 1328	zemining oxjective costs and rame to even ross from the original paper.
1329	STraTS (Total Parameters: 19,374,801)
1330	
1331	<ul> <li>Architecture: As this is an older paper, its structure is not suitable for parallel training. We therefore modified the architecture while retaining the core ideas. We used the last hidden</li> </ul>
1332	state of the decoder at each time step as the event embedding, instead of the original Fusion
1333	Self-Attention mechanism.
1334	• Value and Time Embedding: This model does not use binning. Instead, it embeds Fea-
1335	tures via a look-up table and continuously embeds Value and Time (in hours) via a 2-layer
1336	fully connected layer. The resulting embeddings are summed to form the final event em-
1337	bedding.
1338	<ul> <li>Positional Embedding: No additional positional embedding is used beyond the Time information included in the event embedding.</li> </ul>
1340	• Demographics: Unlike typical models, STraTS encodes demographic information (gender,
1341 1342	race, age) with a separate MLP and concatenates it to the last hidden state.
1343	• Learning Objective: The final embedding is used to predict the value of events occurring
1344	within two hours of each event, trained with a mean squared error loss.
1345	TVID (VIOTE (T. 1.1.)
1346	EHRSHOT (Total Parameters: 20,278,892)
1347	Backbone: Transformer decoder.
1348	• Bin Sharing: No, uses Value-Feature paired tokenization.
1349	• Binning: 10-uniform binning.
	25

• Learning Objective: For each visit, event tokens (Value-Feature paired) are masked with

a probability of  $p_{\text{mask}} = 0.15$ . The model is trained via a multi-class classification loss to

predict the masked token. This is trained separately for different event types (V/S, Lab).

Additionally, for unmasked events, values are altered with a probability of  $p_{\text{anomaly}} = 0.05$ ,

and the model is trained with a binary classification loss to identify whether the value at

• Bin Sharing: Yes, Feature and Value are tokenized separately. Numeric values share bin

• Positional Embedding: Uses rotary positional embedding based only on position, without

1403

1350	• Positional Embedding: Following CLMBR-T Steinberg et al. (2021), it uses Rotary posi-
1351	tional embedding based on position order. Time information is not used beyond ordering
1352	for this component.
1353	• Time Information: A 5-dimensional time vector is concatenated to the to-
1354	ken embedding vector. This vector consists of the z-normalized values of
1355	[age, $\log(\text{age})$ , days since admission, $\log(\text{days since admission})$ , first admission indicator].
1356	The final concatenated vector length is 512.
1357	
1358	• Learning Objective: Next token prediction modeling.
1359 1360	TRADE (Total Parameters: 21,617,664)
1361	• Backbone: Uses a Transformer encoder due to its MLM-based loss.
1362	• Bin Sharing: No, uses Value-Feature paired tokenization.
1363	
1364	• <b>Binning</b> : 9-standard deviation-based binning. For each feature's value distribution, thresh-
1365	olds are set by adding $\{-10, -3, -1, -0.5, 0.5, 1, 3, 10\}$ standard deviations to the mean,
1366	creating 9 bins. This method can be sensitive to outliers, so additional clipping was
1367	performed on 17 vital sign data points based on physician guidelines. The clipping
1368	ranges are as follows: CRR:[0,1], SBP:[0,400], MBP:[0,300], DBP:[0,300], FiO2:[0,1],
1369	HR:[0,200], RR:[0,100], GCS-E:[1,4], GCS-M:[1,6], GCS-V:[1,5], Glucose:[0,1200], O2
1370	saturation:[0,100], Body temperature:[20, 45], Height:[0,1000], Weight:[0,1000].
1371	• Positional Embedding: Absolute positional embedding—performs learned positional em-
1372	bedding using three types of sequential and temporal information: 1) The index of the
1373	current hospital admission, 2) The number of days passed since admission, and 3) The
1374	current age. These are all integers, passed through an embedding layer, and then summed.
1375	• Learning Objective: Uses a standard MLM methodology, masking each token with $p =$
1376	0.2 and using the last hidden state of the masked position to predict the pre-mask label via
1377	a classification loss.
1378	
1379	EHRmamba (Total Parameters: 21,770,752)
1380	
1381	<ul> <li>Backbone: Although the original paper uses Mamba, we applied the same Transformer decoder backbone for a fair comparison.</li> </ul>
1382 1383	Bin Sharing: No, uses Value-Feature paired tokenization.
1384	• <b>Binning</b> : 10-uniform binning (the original paper used 5-uniform binning, but we matched
1385	the bin count for a fair comparison).
1386	•
1387	• <b>Positional Embedding</b> : Uses four types of absolute positional embeddings, which are
1388	summed: (1) Learned PE based on the hospital visit number, (2) Learned PE based on token type (3) Time embedding based on are using the Time West model (Kazemi et al.
1389	token type, (3) Time embedding based on age using the Time2Vec model (Kazemi et al., 2019), and (4) Position-based sin/cos positional embedding from Vaswani et al. (2017).
1390	
1391	• Learning Objective: Uses a next token prediction loss.
1392	
1393	ETHOS (Total Parameters: 20,050,944)
1394	
1395	• Backbone: Transformer decoder.
1396	• Bin Sharing: Yes, Feature and Value are tokenized separately. Numeric values share bin
1397	tokens across different features.
1398	• Binning: 10-uniform binning.
1399	•
1400	• <b>Time Information</b> : Following the original paper, the time interval between each event is
1400	converted into one of 13 discrete tokens, which are inserted into the sequence between
	event tokens.
1402	

• Positional Embedding: Uses a learned positional embedding based on position.

• Learning Objective: Uses NTP loss.

#### D.3 PRE-TRAINING DETAILS

Pre-training was conducted using the training dataset. The hyperparameters were fixed as follows: the batch size was 64, and the number of training epochs was 50. We used the Adam (Adam et al., 2014) optimizer with a learning rates of  $\{5\times10^{-4},1\times10^{-4},5\times10^{-5},1\times10^{-5}\}$ . A 50-step warmup was employed, followed by a cosine annealing schedule that reduced the learning rate to 1/100 of its initial value. Gradient clipping was applied with a threshold of 1.0. The training was performed using either 4 NVIDIA A40 GPUs or 2 NVIDIA RTX PRO 6000 Blackwell GPUs, with Distributed Data Parallel training at a per-GPU batch size of 16 or 32, respectively. The implementation was based on Python version 3.12 and PyTorch (Paszke et al., 2017) version 2.8.0.

# E DOWNSTREAM TASKS AND FINE-TUNING

#### E.1 DOWNSTREAM TASKS

Existing EHR foundation models generally lack generative capabilities, and thus evaluating performance on diverse clinical downstream tasks has been a common practice (Fallahpour et al., 2025; Huang et al., 2024; Burkhart et al., 2025; Renc et al., 2024). While our model possesses strong generative properties, we follow the line of prior work and perform downstream task evaluations to assess the quality of patient representations at each timestamp. We first adopted the four widely used downstream tasks from the MIMIC-III benchmark (Harutyunyan et al., 2019): In-hospital Mortality (IHM), Decompensation-death (Dec-D), Length of Stay (LOS), and Phenotyping (Phe). For these tasks, labels were obtained following the original preprocessing pipelines. To further examine whether the model can capture patient states beyond simple deterioration, we added three additional tasks: Decompensation-arrest (Dec-A), Oliguria/Anuria (HUO), and Vasopressor (Vaso) use. Dec-A is similar to Dec-D, but includes arrest events from CHARTEVENTS.csv in addition to death as decompensation events; the task aims to predict deterioration 24 hours in advance. HUO labels indicate whether the patient currently exhibits oliguria or anuria. Specifically, oliguria is defined as urine output below 0.5 mL/kg/hr for at least 6 hours, and anuria as below 0.1 mL/kg/hr for at least 6 hours. Labels were set for patients with available weight and hourly urine output data, while patients with incomplete information were excluded from training. Vaso labels indicate whether a patient is currently receiving vasopressors. Positive labels were assigned if administration records for Vasopressin, Dobutamine, Epinephrine, Norepinephrine, or Dopamine were present in INPUTEVENTS\_CV.csv or INPUTEVENTS\_MV.csv; negative labels were assigned if other medications were administered but no vasopressors were recorded.

All downstream tasks were measured at the hospitalization level following the original papers (Harutyunyan et al., 2019) (i.e., if a patient had multiple hospital admissions, each admission was treated as a separate EHR sequence). Table 8 provides detailed statistics, including the number of hospitalizations available for each task and the label distribution for each class.

#### E.2 FINE-TUNING DETAILS

The downstream tasks were trained by attaching a task-specific prediction head to the last hidden state of the backbone, applied uniformly to our model and all baselines. Since labels must be inferred using only information available up to each timestep, a causal mask was employed. We froze the pretrained model and trained the prediction heads simultaneously. To evaluate generalization performance under varying input distributions (e.g., absence of laboratory data), we experimented with three settings: (i) the full set of 117 variables (17 vital signs + 100 laboratory measurements), (ii) only 17 vital signs, and (iii) a reduced set of 6 vital signs (SBP, DBP, body temperature, heart rate, respiratory rate, SpO<sub>2</sub>). For training, the original training set was split into train/validation subsets with an 85:15 ratio; the train subset was used for optimization, while the validation subset was used for early stopping and hyperparameter search. As in pre-training, the learning rate lr was selected from  $5 \times 10^{-4}$ ,  $1 \times 10^{-4}$ ,  $5 \times 10^{-5}$ ,  $1 \times 10^{-5}$ . The batch size was fixed 100. Training was performed for 5 epochs with a 50-step warm-up followed by cosine annealing that decayed the learning rate to 1/100 of its initial value.

Each task-specific prediction head consisted of a two-layer MLP. To account for label imbalance, task losses were weighted according to label frequencies in the training dataset. For efficiency,

Table 8: Prediction time and label counts for each downstream task. For binary classification, the positive/negative counts are shown; for multiclass classification, the counts for each class are shown. For tasks with multiple sub-tasks, each sub-task's label counts are shown.

Task	Prediction point	# Hospitalization	# Sub task & # Class	Train labels	Test labels
IHM	48h after ICU adm.	17,903 / 3,236	1/2	2424 / 15479	374 / 2,862
Dec-D	Hourly	35,365 / 6,237	1/2	61,018 / 2,847,424	9,684 / 513,552
Dec-A	Hourly	35,365 / 6,237	1/2	65,298 / 2,843,144	10,239 / 512,997
	•			[cls1: 790,196, cls2: 503,423,	[cls1: 139,682, cls2: 90,478
				cls3: 316,774, cls4: 215,075,	cls3: 56,289, cls4: 38,795
LOS	Hourly	35,523 / 6,265	1 / 10	cls5: 158,987, cls6: 124,146,	cls5: 28,542, cls6: 22,225
	•			cls7: 100,890, cls8: 84,241,	cls7: 18,077, cls8: 15,145
				cls9: 312,111, cls10: 319,619]	cls9: 55,997, cls10: 60,710
				t1: 7,644 / 28,029	t1: 374 / 2862
				t2: 2660 / 33013	t2: 1331 / 4945
				t3: 3657 / 32016	t3: 415 / 5861
				t4: 11434 / 24239	t4: 675 / 5601
				t5: 4821 / 30852	t5: 2028 / 4248
				t6: 4675 / 30998	t6: 831 / 5445
				t7: 7349 / 28324	t7: 789 / 5487
				t8: 2565 / 33108	t8: 1337 / 4939
				t9: 9550 / 26123	t9: 442 / 5834
				t10: 11497 / 24176	t10: 1683 / 4593
				t11: 3444 / 32229	t11: 2074 / 4202
				t12: 6869 / 28804	t12: 593 / 5683
Phe	End of stay	35,563 / 6,273	25 / 2	t13: 10362 / 25311	t13: 1205 / 5071
	•			t14: 14922 / 20751	t14: 1813 / 4463
				t15: 9617 / 26056	t15: 2653 / 3623
				t16: 2573 / 33100	t16: 1667 / 4609
				t17: 4775 / 30898	t17: 495 / 5781
				t18: 3170 / 32503	t18: 819 / 5457
				t19: 1816 / 33857	t19: 556 / 5720
				t20: 1435 / 34238	t20: 355 / 5921
				t21: 3080 / 32593	t21: 272 / 6004
				t22: 4970 / 30703	t22: 570 / 5706
				t23: 6468 / 29205	t23: 852 / 5424
				t24: 5118 / 30555	t24: 1111/5165
				t25: 2779 / 32894	t25: 874 / 5402
шо	TT1	11 001 / 2 107	2/2	t1: 148737 / 891153	t1: 56102 / 162429
HUO	Hourly	11,891 / 2,187	2/2	t2: 148737 / 1053830	t2: 26682 / 191849
Vaso	Hourly	35,438 / 6,249	1/2	202,765 / 2,747,736	36,306 / 494,215

downstream tasks were conducted in a multi-task setting where all seven tasks were jointly optimized; the final objective was defined as the average of the task-specific losses. Since probing does not allow training of new tokens, we instead introduced a <Birth> token at each timestep, serving the same role as the <SOS> token, and used its representation for task prediction. For baseline models, we preserved their original binning and embedding procedures, while unifying the learning objective to the downstream tasks. All downstream tasks were conducted on a single NVIDIA A40 GPU.

### F RESULTS

#### F.1 DOWNSTREAM TASK RESULTS

Tables 9, 10, and 11 report the results of all downstream tasks. Table 9 presents the loss, number of features used, and selected pretraining/downstream-task learning rate for each experiment, while the remaining results are shown in Tables 10 and 11.

#### F.2 EHR GENERATION AND EVALUATION

We generated EHRs with a low temperature of 0.2. Below are examples generated from the same initial EHR, shown in order for our model (Figure 8) and the ETHOS model (Figure 9).

Figure 10 shows the prompt we used for LLM-based evaluation.

Table 9: Meta Information

Method	Value sharing	Downstream # features	test-loss	valid-loss	train-loss	Pre-training LR	Downstream LR
HEART	X	117	5.3043	5.3681	5.3714	0.0005	0.0005
HEART	X	17	5.4340	5.4604	5.4713	0.0005	0.0005
HEART	X	6	5.8346	5.8703	5.8702	0.0005	0.0005
MOTOR	X	117	4.9454	4.9020	4.9118	0.0001	0.0005
MOTOR	X	17	5.2117	5.1570	5.1803	0.0001	0.0005
MOTOR	X	6	5.6451	5.6455	5.6519	0.0001	0.0005
EHRSHOT	X	117	5.8406	5.9138	5.9241	0.0005	0.0005
EHRSHOT	X	17	6.0777	6.1018	6.1022	0.0005	0.0005
EHRSHOT	X	6	6.3406	6.3509	6.3325	0.0005	0.0005
TRADE	X	117	5.2599	5.2480	5.2807	0.0005	0.0005
TRADE	X	17	5.4538	5.4529	5.4839	0.0005	0.0005
TRADE	X	6	6.0481	6.0643	6.0702	0.0005	0.0005
EHRmamba	X	117	5.1366	5.1603	5.2125	0.0005	0.0005
EHRmamba	X	17	5.4389	5.4228	5.4678	0.0005	0.0005
EHRmamba	X	6	5.9261	5.9404	5.9573	0.0005	0.0005
Ours (No value share)	X	117	4.6861	4.6386	4.6273	0.0005	0.0005
Ours (No value share)	X	17	4.9070	4.8621	4.8664	0.0005	0.0005
Ours (No value share)	X	6	5.3675	5.3538	5.4045	0.0005	0.0005
FM4EHR	О	117	6.4288	6.4611	6.4344	0.0005	0.0001
FM4EHR	O	17	6.3888	6.4391	6.3969	0.0005	0.0005
FM4EHR	O	6	6.3972	6.4394	6.4069	0.0005	0.0005
ETHOS	O	117	4.9710	5.0374	4.9378	0.0005	0.0001
ETHOS	O	17	5.2479	5.2044	5.1124	0.0005	0.0001
ETHOS	O	6	5.5724	5.5714	5.5213	0.0005	0.0001
STraTS	O	117	5.7857	5.8492	5.8505	0.0001	0.0005
STraTS	O	17	5.8123	5.8335	5.8446	0.0001	0.0005
STraTS	O	6	6.0711	6.0760	6.0834	0.0001	0.0005
Ours (Value share)	О	117	4.8786	4.8954	5.0276	0.0005	0.0005
Ours (Value share)	O	117	5.0430	5.0862	5.1887	0.0005	0.0001
Ours (Value share)	O	6	5.5613	5.6554	5.7523	0.0005	0.0005

1	56	6
	56	
	56	
	56	
1	57	
1	57	
1	57	
1	57	
1	57	
1	57	
1		
	57 57	
1	57 57	
1	57 57	
1	57	
	58	
	58	
	58	
1	58	3
1	58	4
1	58	5
1	58	6
1	58	7
1	58	8
1	58	9
	59	
	59	
	59	
	59	
	59	
	59	
	59 50	
	59	
	59	
	59	
	60	
1	60	1
1	60	2
1	60	3
1	60	4
	60	
	60	
	60	
	60	
	60	
	61	
	61	
	61	
	61	
1	61	4
1	61	5
1	61	6
1	61	7
	61	
1	61	0

PP.16	0.2543	0.1899	0.4827	0.2838	0.2329	0.2162	0.1566	0.1767	0.4197	0.1983	0.1757	0.2658	9061.0	0.1624	0.5912	0.3246	0.2589	0.1572	0.1305	0.1305	0.3125	0.2395	0.1779	0.2206	0.2239	0.2087	0.4477	0.3210	0.2386
PR.16	0.7242	6313	8503	.7348	51697	1899'	5575	21097	8214	96394	2665	7.462	.6245	5794		8197.0		5518	2000	2000	2992	2969'0	1219	0.6735	62991	96990	3.8345	0.7680	66290
PP.15 F	0.1210	_	_	_	_	Ĭ	_	_	_	_	_	_	Ĭ	.0907	-	) 1618 (	_	ll .				0.1173 (				) 6960'	. 1927 (	_	Ĭ
PR.15 F	0.6260	_	_	_	_	_	_	_	_	_	_	_	_	-	ľ	0.6922 (	_	-	_	_	_	0.6596	_		_	0.5547 (	0.7677 C	_	Ĭ
PP.14 P.	0.4322 0	Ĭ	Ĭ	_	_	_	_	_	_	_	_	_	_	Ĭ	ľ	0.4612 0	_					0.4109 0				3324 0	0.4967 0	-	-
PR.14 P.	0 96690	-	-	_	_	_	_	_	_	-	-	_	-	-	ľ	0.7133 0.	_	-	Ī	7	_	0.6732 0.	_		_	_	0.7534 0.	_	_
P-P.13 P.	0.5210 0,0.4917 0,0	_	Ĭ	_	_	_	_	Ĭ	_	_	Ĭ	_	_	.4181 0.		.5842 0.		-	-	_	_	0.5597 0.0	_	_	_	-	0.5782 0.	_	_
P-R.13 P	0.5826 0.5	_	_	_	_	_	_	_	_	_	_	_	_	_	ľ	0.6883 0.5	_	∥ັ	Ĭ	_	_	0.6587 0.3	_	_	_	_	0 8799.0	-	-
PP.12 P	0.4657 0.4397 0.5													3363 07	ľ	0.5027 0.0	_					0.4855 0.4				3715 0.0	_	0.4728 0.4	Ĭ
PR.12 P	0.6655 0.0	_	_	_	_	_	_	_	_	-	_	_	_	_	-	0.7271 0.3	_	-	-	_	_	0.7011 0.4	_	_	_	0.6474 0.3	.0. 8727.0	_	Ĭ
P.P.11 P.4	0.3330 0.6	-	_	_	_	_	_	_	_	_	_	_	_	-	-	0.4059 0.7	_	-	Ī	7	_	0.3013 0.7	-	Ī	_	_	0.3019 0.7	_	Ĭ
_	0.7287 0.3													.5617 0.2	-	0.7600 0.4	_					0.6726 0.3					_	3.6360 0.2	Ĭ
10 PR.11		_	_	Ĭ	Ĭ	Ĭ	_	Ĭ	Ĭ	Ĭ	Ĭ	Ĭ	Ĭ	0.1266 0.5	ľ	0.4509 0.7	_	-	_	_	_	0.3158 0.6	_		_	0.1459 0.5	0.4538 0.6	_	Ĭ
.10 PP.10	53 0.4154	_	_	_	_	_	_	_	_	_	-	_	_	_	-	_	_	-	Ī	7	_	_	_		-	_	_	_	_
9 P-R.10	20 0.8353 93 0.8415														-	30 0.8538	_					78 0.7556					40 0.8403	_	Ĭ
.9 P.P9	17 0.5820 09 0.5593														-	57 0.6630	_	∥ -	-	_	_	21 0.6178	_	-	_	-	79 0.7040	-	_
:8 P.R.9	08 0.7417	_	Ī	_	_	_	_	_	Ī	_	Ī	_	_	96 0.7142	-	01 0.8057	_	ll .				19 0.7721					81 0.8279	_	Ĭ
.8 PP.8	11 0.4408 10 0.4164	_	_	_	_	_	_	_	_	_	_	_	_	_	Ι-	89 0.5301	-	ll .				32 0.4919					1809'0 16	_	-
7 PR.8	77 0.7211 45 0.6910														ľ	55 0.7689	_	∥ັ	_	Ĭ	Ĭ	21 0.7432	Ĭ	Ĭ	Ĭ	_	85 0.8091	Ĭ	Ĭ
.7 PP.7	25 0.1677 86 0.1545														ľ	52 0.2255	_	ll .				83 0.1721					72 0.1885	_	Ĭ
:6 P-R.7	53 0.7025	_	_	-	-	-	-	-	-	-	-	-	-	_	ľ	40 0.7652	_					93 0.7383					53 0.7572	_	_
.6 PP.6	74 0.2753 67 0.2811	-	_	_	_	_	_	_	_	_	_	_	_	-	-	59 0.3440	_	-	-	_	_	35 0.3293	_	_	_	-	00 0.3453	-	_
5 P.R.6	50 0.6174 05 0.6167	_	_	-	-	-	-	-	-	-	-	-	-	_	-	12 0.6859	_	-	_	-	-	44 0.6535	-	-	-	_	14 0.6800	-	_
5 P-P.5	85 0.2150 79 0.2005														-	05 0.2812	_					67 0.2044					52 0.2714		
4 P.R.5	90 0.6785 96 0.6679	_	_	-	-	-	-	-	-	-	-	-	-	82 0.5843	-	44 0.7305	-	-	_	-	-	797.00 0.6767	-	-	-	_	59 0.7152	-	-
.4 P-P.4	70 0.2590 70 0.2286	_	_	-	-	-	-	-	-	-	-	-	-		-	28 0.3644	_	-	Ī	7	_	33 0.2830	_	_	_	-	12 0.4959	_	_
3 P.R.4	95 0.7270 45 0.6870													52 0.5616		11 0.7828		ll .				86 0.7233					24 0.8512		
R.3 PP.3	13 0.4195 87 0.4145	_	_	_	_	_	_	_	_	_	_	_	_	85 0.3952	Ι-	28 0.5311	-	∥ -	-	_	_	14 0.4886	_	_	_		62 0.5124	_	_
-i-	30 0.6313													_	Ι.	51 0.7228		ll .								08 0.6863	64 0.7162		
.2 PP.2	43 0.2230 23 0.1737	_	_	-	-	-	-	-	-	-	Ī	_	_	95 0.1370		63 0.2751		II				14 0.2017				52 0.1708		47 0.4233	
1 P.R.2	44 0.7143 83 0.6423	-	_	_	13 0.6739	_	_	_	_	22 0.6679	07.09.0 86	0.6756	01 0.6672	75 0.5695	-	23 0.7563	_		61 0.5000	_	_	_	-	Ĭ	_	21 0.6552		63 0.7747	
:1 P.P.1	73 0.2444	_	_	Ĭ	Ĭ	_	_	_	_	Ĭ	Ĭ	_	_	05 0.1175	ı	02 0.4623		9880:0	00 0.0661	-	_		_		_	21 0.1421		03 0.3063	
20 P-R.I	68 0.7973 48 0.7932		Ξ.	_	_	_	-	-	Ī	_			_	96 0.6305	306 0.8992		52 0.7655	II.,	121 0.5000	_		-	60 0.7145	_	_	41 0.6621		170 0.8603	
@0 PP.0	-	0	_	Ĭ	65 0.3770	_	_	28 0.2999	0	76 0.3816	33 0.3408	81020 59	65 0.3172	12 0.2996	9059.0		37 0.4052	13 0.1876	00 0.2121	_	_	_	_	-	_	72 0.3141	_	_	33 0.3846
PR.@0	3 0.7784 9 0.7179	Ĭ	Ĭ	_	Ĭ	Ĭ	Ĭ	_	_	_	Ĭ	2	_	9 0.6112	9 0.8763		6 0.7237	ľ	0.5000	_	_	2 0.7222	-	_	_	2 0.6472		_	
. IP.	9 0.4423										_	_	_	8 0.3109	0.6069		2 0.4016		3 0.1519			3 0.4162	-	2 0.3111	_	3 0.2192	l	_	0.3660
g IR.	0.8379	0.739	0.872	0.840	0.761	0.801	0.719	0.6235	0.827	0.805	0.672	0.867	0.8125	0.704	0.889	0.8694	0.801	0.617	0.5773	0.496	0.8586	0.8173	0.755	0.7592	0.706	0.699	0.8760	0.864	0.775
Value sharing	××	×	×	×	×	×	×	×	×	×	×	×	×	×	×	×	×	0	0	0	0	0	0	0	0	0	0	0	С
Methods	HEART	HEART	MOTOR	MOTOR	MOTOR	EHRSHOT	EHRSHOT	EHRSHOT	TRADE	TRADE	TRADE	EHR mamba	EHR mamba	EHR mamba	Ours (No value share)	Ours (No value share)	Ours (No value share)	FM4EHR	FM4EHR	FM4EHR	ETHOS	ETHOS	ETHOS	STraTS	STraTS	STraTS	Ours (Value share)	Ours (Value share)	Ours (Value share)

Table 10: Performance Results (Part 1); I., and P. denote In-Hospital-Mortality, and Phenotyping respectively. R., and P. denote AUROC and AUPRC respectively.

R., and P. denote

Table 11: Performance Results (Part 2); P., L., H., and V. denote Phenotyping, Length-of-Stay, Oliguria/Anuria, and Vasopressor respectively.

AUROC and AUPRC respectively. Value Value Value Ours (No va Ours (No va Ours (No va

# G LLM USAGE CLARIFICATION

In addition to the uses of LLMs described in the main text, we employed them for summarizing content, translation, grammar correction, and sentence refinement during the writing of the manuscript. In the early stages of the study, we used LLMs to search for related work, and the retrieved papers were then read and verified by the researchers.

1778 1779

1780

```
1729
1730
1731
1732
1733
1734
               Birth: 1845. 05. 11
                                      hemoglobin: 31.0
                                                            - GCS: 15
                                                                                   - HR: 116
1735
                                                            - GCS-V: 5
               Sex: Female
                                        Mean corpuscular
                                                                                   - MBP: 72
1736
               Ethnicity: WHITE
                                      hemoglobin concen-
                                                            - HR: 117
                                                                                   - O2 saturation: 79
1737
                                      tration: 34.4
                                                             - MBP: 73
                                                                                   - RR: 19
               Age: 300
                                      - Mean corpuscular
                                                            - O2 saturation: 92
                                                                                   - SBP: 116
1738
               2145-05-11 17:12:55
                                      volume: 90.0
                                                             - RR: 19
                                                                                   2145-05-12 06:00:00
1739
                - ICU transfer
                                      - Platelets: 200.0
                                                            - SBP: 115
                                                                                   - DBP: 47
1740
               2145-05-11 17:30:00
                                                                                   - HR: 114
                                      - Prothrombin time :
                                                            - Temperature : 36.9
1741
               - DBP : 62
                                                            2145-05-12 00:00:00
                                                                                   - MBP: 70
                                      13.3
1742
               - SBP: 103
                                      - Partial thromboplas-
                                                            - SBP: 116
                                                                                   - O2 saturation: 92
               - HR: 123
1743
                                      tin time: 32.4
                                                             - RR: 19
                                                                                   - RR: 20
                                                                                   - SBP: 107
               - MBP: 83
                                      - Red Cell Distribu-
                                                            - O2 saturation: 91
1744
                                      tion Width: 13.3
                - RR:4
                                                             - HR: 118
                                                                                   2145-05-12 07:00:00
1745
               2145-05-11 18:00:00
                                      - Red blood cell count
                                                            - DBP: 53
                                                                                   - HR: 111
1746
                                                             - MBP: 74
               - DBP: 51
                                      : 3.74
                                                                                   - DBP: 49
1747
               - HR: 110
                                      - White blood cell
                                                            2145-05-12 01:00:00
                                                                                   - GCS-E: 4
1748
               - MBP: 73
                                      count: 12.20
                                                             - DBP: 50
                                                                                   - GCS-M: 6
               - RR: 17
                                      - RR: 18
                                                            - HR: 116
                                                                                   - GCS: 14
1749
                                      2145-05-11 19:00:00
                                                            - MBP: 76
               - SBP: 111
                                                                                   - GCS-V: 5
1750
                  -# Gen Start #-
                                      - DBP: 56
                                                            - O2 saturation: 93
                                                                                   - MBP: 73
1751
                                      - HR: 119
               - O2 saturation: 90
                                                             - RR: 18
                                                                                   - O2 saturation: 93
1752
                                                                                   - RR: 19
               - Weight: 54.2
                                      - MBP: 73
                                                             - SBP: 113
               - Temperature: 36.7
                                      - O2 saturation: 92
                                                                                   - SBP: 114
1753
                                                            2145-05-12 02:00:00
               - GCS: 15
                                      - RR: 19
                                                            - HR: 111
                                                                                   - Temperature : 36.6
1754
               - GCS-M: 6
                                      - SBP: 111
                                                            - DBP: 49
                                                                                   2145-05-12 08:00:00
1755
               - GCS-E: 4
                                      2145-05-11 20:00:00
                                                            - MBP: 74
                                                                                   - DBP: 52
1756
               - GCS-V: 5
                                      - DBP: 50
                                                            - RR: 20
                                                                                   - HR: 112
1757
               - O2 saturation: 92
                                      - MBP: 76
                                                            - SBP: 116
                                                                                   - MBP: 71
1758
                - HR: 117
                                      - HR: 110
                                                             - O2 saturation: 92
                                                                                   - O2 saturation: 89
               - SBP: 108
                                      - RR: 19
                                                                                   - RR: 19
                                                            2145-05-12 03:00:00
1759
               - Glucose: 261
                                      - SBP: 111
                                                            - DBP: 52
                                                                                   - SBP: 114
1760
                                      - O2 saturation: 90
               - MBP: 76
                                                             - GCS-E: 4
                                                                                   2145-05-12 09:00:00
1761
               - DBP: 50
                                      2145-05-11 21:00:00
                                                            - GCS-M: 6
                                                                                   - DBP: 52
1762
               - Anion gap: 15.0
                                      - DBP: 51
                                                            - GCS: 15
                                                                                   - HR: 116
1763
                                      - HR: 113
                                                            - GCS-V: 5
                                                                                   - MBP: 76
               - Bicarbonate : 24.0
               - Calcium Total: 8.6
                                      - MBP: 71
                                                            - HR: 119
                                                                                   - O2 saturation: 88
1764
               - Chloride: 105.0
                                      - O2 saturation: 93
                                                             - MBP: 72
                                                                                   - RR: 20
1765
                - Creatinine: 0.9
                                      - RR: 19
                                                             - O2 saturation: 91
                                                                                   - SBP: 110
1766
               - Magnesium: 1.8
                                      - SBP: 115
                                                            - RR: 20
                                                                                   2145-05-12 10:00:00
1767
                                                            - SBP: 109
               - Phosphate: 2.7
                                      2145-05-11 22:00:00
                                                                                   - DBP: 53
1768
               - Potassium: 4.0
                                      - DBP: 50
                                                             - Temperature: 36.9
                                                                                   - HR: 113
               - Sodium: 142.0
                                      - HR: 114
                                                            2145-05-12 04:00:00
                                                                                   - MBP: 75
1769
               - Blood urea nitrogen
                                     - MBP: 74
                                                             - DBP: 50
                                                                                   - O2 saturation: 92
1770
                                      - O2 saturation: 89
               : 12.0
                                                            - HR: 117
                                                                                   - RR: 20
1771
               - Hematocrit: 33.0
                                      - RR: 19
                                                            - MBP: 73
                                                                                   - SBP: 111
1772
                                                             - RR: 18
                                                                                   2145-05-12 11:00:00
                - Hemoglobin chem-
                                      - SBP: 116
                                                            - SBP: 113
1773
               istry: 11.3
                                      2145-05-11 23:00:00
                                                                                   - DBP: 52
                                      - DBP: 51
                 Prothrombin time
                                                             - O2 saturation: 85
                                                                                   - GCS-E: 4
1774
                                      - GCS-E: 4
                                                            2145-05-12 05:00:00
                                                                                   - GCS-M: 6
               INR: 1.10
1775
                 Mean corpuscular
                                      - GCS-M: 6
                                                             - DBP: 49
1776
1777
```

Figure 8: Temporal EHR history generated by our model (used value share version for fair comparison). The data before this "—# Gen Start #—" marker is given, and the data after it is generated.

1831

1832

```
1784
1785
1786
1787
1788
1789
               Birth: 1845. 05. 11
                                      - DBP : 42
                                                            - O2 saturation: 85
                                                                                   2145-05-12 05:05:00
1790
                                                            - HR: 119
                                      - MBP: 60
                                                                                   - HR: 115
               Sex: Female
1791
               Ethnicity: WHITE
                                      2145-05-11 20:45:00
                                                            - DBP: 43
                                                                                   - O2 saturation: 93
1792
                                                            - MBP: 60
               Age: 300
                                      - RR: 30
                                                                                   - RR: 29
1793
                                      - SBP: 84
                                                            2145-05-12 00:30:00
                                                                                   2145-05-12 05:25:00
1794
               2145-05-11 17:12:55
                                     - O2 saturation: 93
                                                            - RR: 32
                                                                                   - DBP: 31
               - ICU transfer
                                      - HR: 115
                                                                                   - MBP: 59
                                                            - HR: 114
1795
               2145-05-11 17:30:00
                                      - DBP: 44
                                                            - O2 saturation: 93
                                                                                   - SBP: 94
1796
               - DBP: 62
                                      - MBP: 49
                                                            - SBP: 68
                                                                                   2145-05-12 06:25:00
1797
               - SBP: 103
                                      2145-05-11 21:25:00
                                                            - MBP: 43
                                                                                   - RR: 33
1798
               - HR: 123
                                                            - DBP: 43
                                      - HR: 122
                                                                                   - O2 saturation: 92
1799
               - MBP: 83
                                      - MBP: 58
                                                            2145-05-12 01:00:00
                                                                                   - HR: 120
                                                            - DBP: 42
                                                                                   2145-05-12 06:45:00
               - RR: 4
                                      - O2 saturation: 93
               2145-05-11 18:00:00
                                     - RR: 29
                                                            - HR: 117
                                                                                   - DBP: 27
1801
               - DBP: 51
                                      - SBP: 88
                                                            - MBP: 53
                                                                                   - MBP: 55
                                                            - O2 saturation: 92
               - HR: 110
                                      - DBP: 40
                                                                                   - SBP: 83
1803
               - MBP: 73
                                      2145-05-11 22:05:00
                                                            - RR: 31
                                                                                   2145-05-12 07:35:00
               - RR: 17
                                      - SBP: 81
                                                            - SBP: 92
                                                                                   - HR: 119
               - SBP: 111
                                      - RR: 29
                                                            2145-05-12 01:50:00
1805
                                                                                   - O2 saturation: 91
                  -# Gen Start #-
                                      - O2 saturation: 87
                                                                                   - RR: 29
                                                            - HR: 112
1806
               - O2 saturation: 97
                                      - HR: 134
                                                            - O2 saturation: 90
                                                                                   2145-05-12 08:05:00
1807
               2145-05-11 18:30:00
                                     - DBP: 43
                                                                                   - DBP: 42
                                                            - RR : 33
1808
                                                            2145-05-12 02:40:00
               - DBP: 51
                                      - MBP: 59
                                                                                   - MBP: 42
1809
               - HR: 124
                                      2145-05-11 23:00:00
                                                            - DBP : 39
                                                                                   - SBP: 86
               - MBP: 67
                                      - SBP: 88
                                                            - MBP: 48
                                                                                   2145-05-12 09:05:00
1810
               - O2 saturation: 97
                                      - RR: 36
                                                            - SBP: 91
                                                                                   - HR: 116
1811
               - RR: 32
                                      - O2 saturation: 91
                                                            2145-05-12 03:15:00
                                                                                   - MBP: 50
1812
                                      - HR: 115
               - SBP: 89
                                                            - DBP: 43
                                                                                   - O2 saturation: 92
1813
               2145-05-11 19:10:00
                                     - DBP: 22
                                                            - HR: 120
                                                                                   - RR: 37
                                      - MBP: 58
1814
               - SBP: 80
                                                            - MBP: 43
                                                                                   - SBP: 81
               - RR: 35
                                      2145-05-11 23:30:00
                                                            - O2 saturation: 91
                                                                                   - DBP: 43
1815
               - O2 saturation: 96
                                      - SBP: 84
                                                            - RR: 38
                                                                                   2145-05-12 10:00:00
1816
                                                                                   - HR: 132
               - HR: 119
                                      - RR: 32
                                                            - SBP: 94
1817
               - DBP: 42
                                      - O2 saturation: 92
                                                            2145-05-12 04:05:00
                                                                                   - O2 saturation: 92
1818
               - MBP: 59
                                                                                   - RR: 33
                                      - HR: 137
                                                            - RR: 29
1819
               2145-05-11 19:45:00
                                     - DBP: 38
                                                            - O2 saturation: 91
                                                                                   2145-05-12 10:15:00
               - DBP: 44
                                      - MBP: 49
                                                            - HR: 116
                                                                                   - DBP: 35
1820
               - HR: 129
                                      2145-05-11 23:45:00
                                                            - MBP: 51
                                                                                   - MBP: 55
1821
               - MBP: 60
                                                            - DBP: 38
                                                                                   - SBP: 94
                                      - DBP: 39
               - O2 saturation: 90
                                      - HR: 115
                                                            - SBP: 88
                                                                                   2145-05-12 11:10:00
                                                            2145-05-12 04:20:00
               - RR : 30
                                      - MBP : 52
                                                                                   - DBP : 43
1824
               - SBP: 94
                                      - O2 saturation: 91
                                                            - RR: 30
                                                                                   - HR: 136
               2145-05-11 20:15:00
                                                            - O2 saturation: 87
                                     - RR: 34
                                                                                   - MBP: 34
1825
               - RR: 32
                                      - SBP: 86
                                                            - HR: 123
                                                                                   ...
1826
                                                            - DBP: 42
               - SBP: 78
                                      2145-05-12 00:15:00
1827
                                                            - MBP: 54
               - O2 saturation: 93
                                      - SBP: 88
1828
               - HR: 140
                                      - RR: 36
                                                            - SBP: 89
1829
1830
```

Figure 9: Temporal EHR history generated by ETHOS. The data before this "—# Gen Start #—" marker is given, and the data after it is generated.

You are a physician with extensive ICU experience and an AI researcher familiar with text generation models, such as LLMs. In this survey, you will compare the quality of EHR texts generated by two different models from the same initial patient history. The quality of an EHR depends on whether the right clinical events occur at the right times. Please consider both the timing of events and the appropriateness of the events themselves. First, you will see a few sample ICU EHR texts. Then, for each pair of generated candidates (A and B), you will be asked to decide which one appears more realistic. <Sample EHR texts> 1. ## Sample 1 ## 2. ## Sample 2 ## 3. ## Sample 3 ## <end of EHR samples> <Evaluation candidate A> ## ETHOS generated Sample (Random order; Ours can be candidate A) ## <Evaluation candidate B> ## Ours generated Sample (Random order; ETHOS can be candidate B) ## <Compare two candidates> 

Figure 10: LLM input prompt for generated EHR evaluation. We compared the generative performance of our model and ETHOS on LLMs with this prompt.