

# Foundation Models for Electronic Health Records

BIODS 271: Foundation Models for Healthcare

March 4, 2026

**Jason Fries PhD, Assistant Professor**  
Department of Biomedical Data Science  
Division of Computational Medicine  
Department of Medicine  
Stanford University



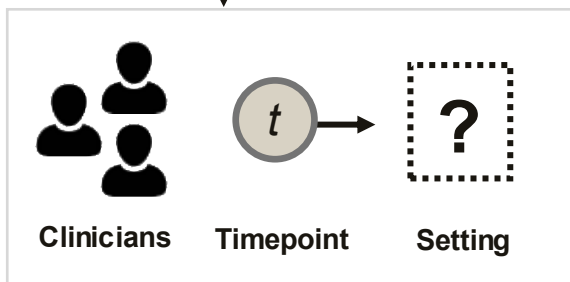
Stanford University  
Department of Biomedical Data Science  
Division of Computational Medicine  
Department of Medicine  
Stanford University

# Outline

- **The Missing Context Problem in Healthcare AI**
- **EHR Data & Tasks**
  - Electronic Health Records (EHRs)
  - AI for Healthcare Tasks
- **Modeling: FMs for Structured EHRs**
  - Formulating Self-Supervision
  - Pretraining Objectives
- **Evaluation**
- **Future: Research Opportunities**

# Sequential Decision Making

Clinical Landmarks



Setting

## AI-Assisted Tumor Board



Joel Neal, MD, PhD

Clinicians can 'chat' with medical records through new AI software, ChatEHR

By [Hanae Armitage](#)

ChatEHR, artificial intelligence software developed at Stanford Medicine, is expediting chart reviews and other tasks by allowing clinicians to ask questions of medical records.

Artificial Intelligence (AI) | June 05, 2025

# How LLMs Are Trained and Evaluated for Medical Use

## Clinical Vignette

A 62-year-old man with a 40 pack-year smoking history presents with **persistent cough** and **10-lb unintentional weight loss**. CT chest reveals a **3.2-cm spiculated right upper lobe mass** with **mediastinal lymphadenopathy**. Biopsy confirms **lung adenocarcinoma**, and molecular testing shows an **EGFR exon 19 deletion**. MRI brain demonstrates **two small asymptomatic enhancing lesions consistent with metastases**. His **ECOG performance status is 1**.

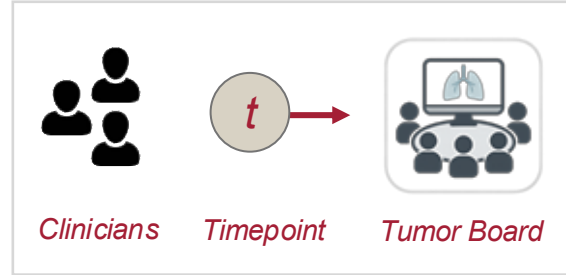
Q: What is the most appropriate initial systemic therapy?

LLM Response

Unverifiable tasks are scored by rubrics

Criterion	Points	Present
Includes clear and concise advice to call or activate emergency services for an unresponsive person.	+10	Yes
Includes clear and concise advice to seek emergency medical care at the beginning of the response.	+9	Yes

OpenAI HealthBench 2025

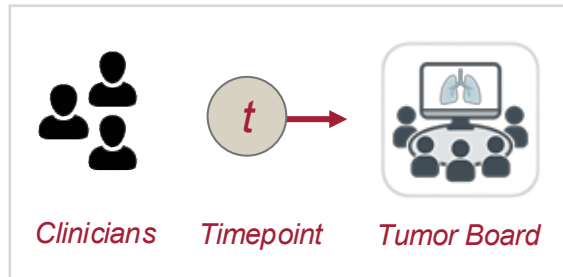
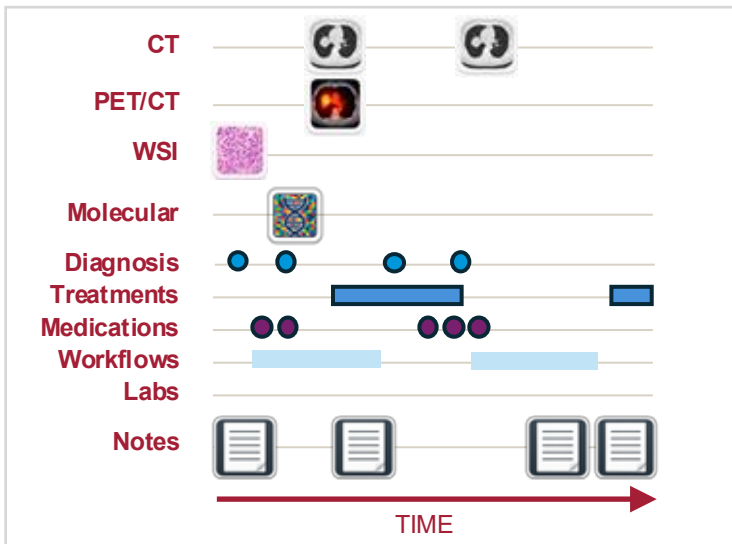


- Narrow view of clinical data
- Assumes data is **pre-cleaned and perfectly framed** -- real data is multimodal and messy!
- **Rewards shortcut reasoning**



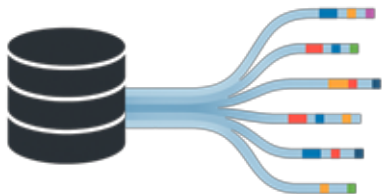
# How LLMs Should Be Trained and Evaluated for Medical Use

## Longitudinal, Multimodal Timeline



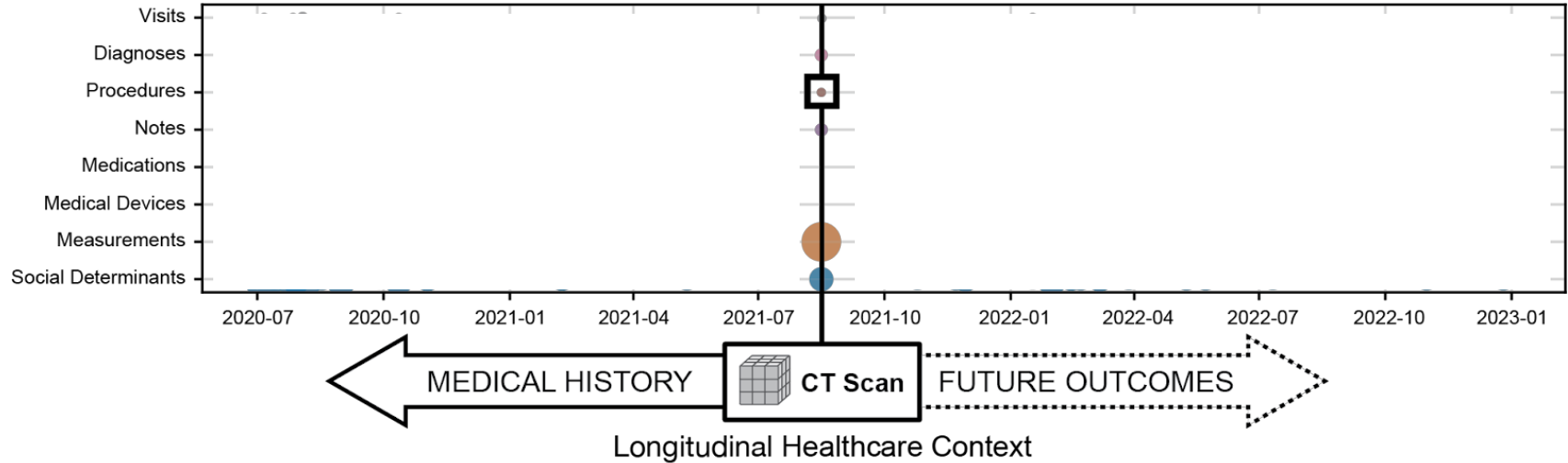
***“Find me patients like this one, who developed late toxicity after radiation”***

**We want to find patients with similar disease trajectories**



**This complexity breaks today's models**

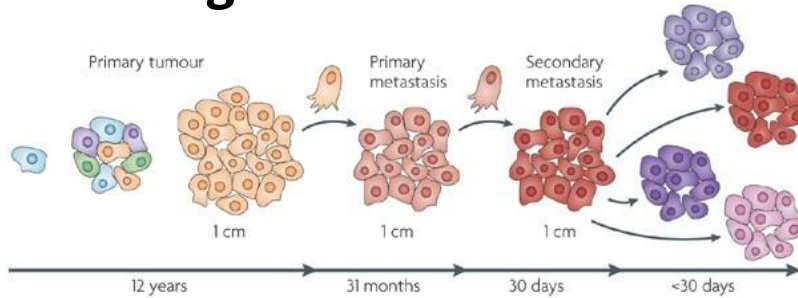
# Missing Context Problem



- Limited use of **longitudinal health data to supervise models**
- Limited insight into the **distinct needs of different stakeholders**

# Human Health is Time-Varying

## Cancer Progression

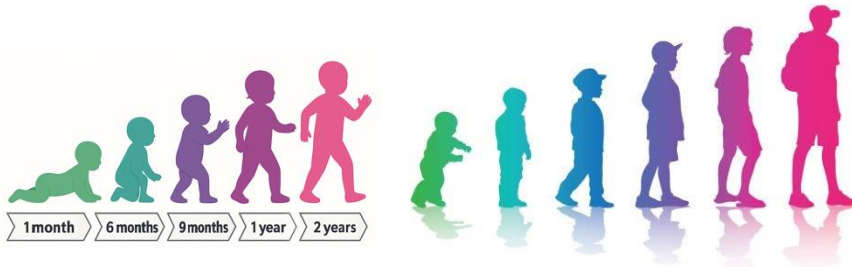


Klein 2009

Nature Reviews | Cancer

*How likely is this patient to develop gastrointestinal cancer in 10 years, 5 years, or 2 years?*

## Pediatric Development

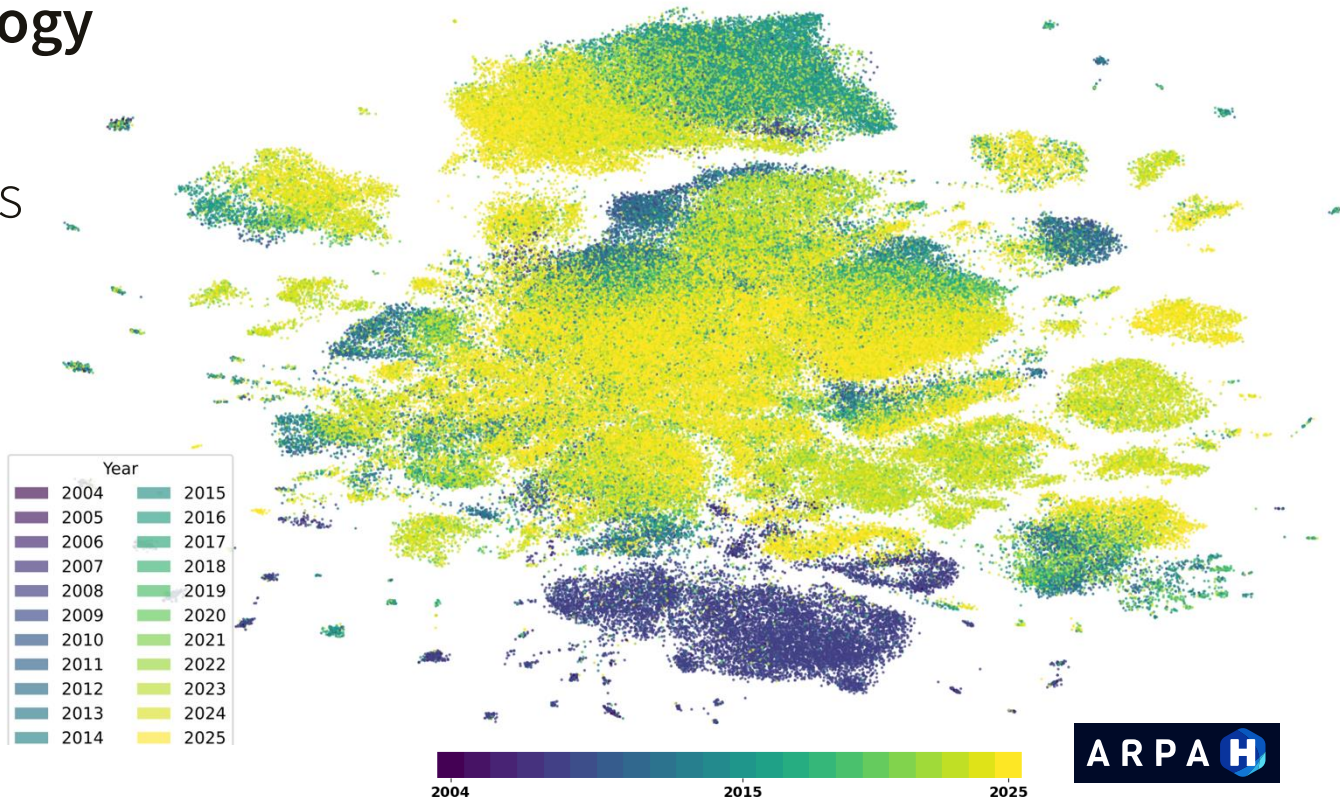


*By what age will this child receive an autism spectrum disorder diagnosis: 18 months, 3 years, or 10 years?*

# Finding Similar Patients in Latent Space

## VISTA Oncology Data Lake

**211k** patients



# Enhancing Omics Analyses: Cancer Mortality

nature machine intelligence



Article

<https://doi.org/10.1038/s42256-024-00974-9>

## A machine learning approach to leveraging electronic health records for enhanced omics analysis

Received: 23 August 2024

Accepted: 16 December 2024

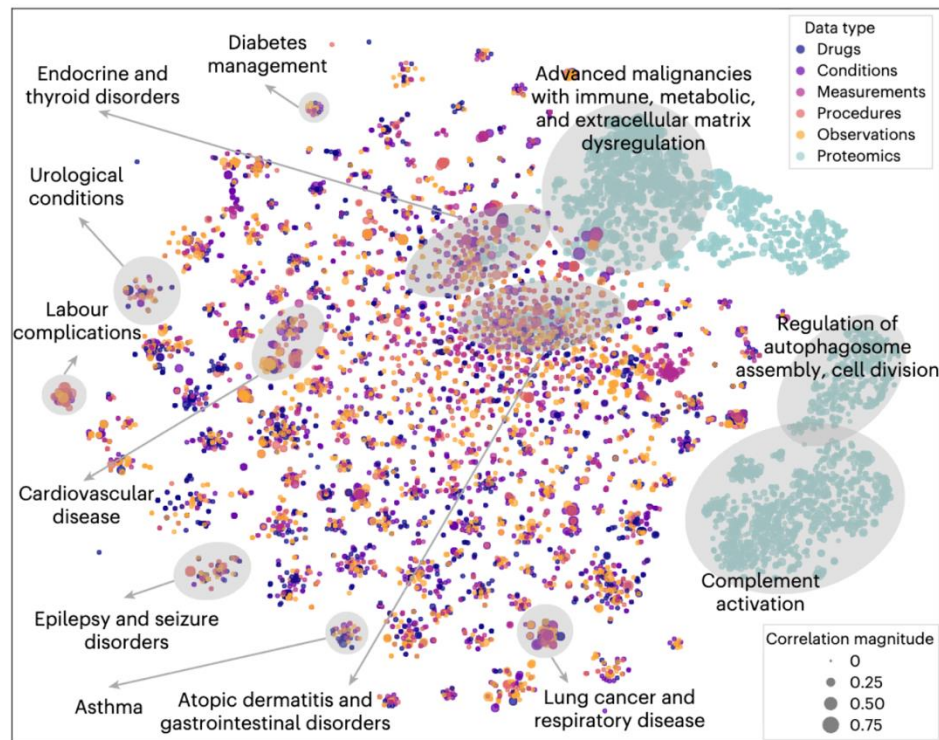
Published online: 16 January 2025

Check for updates

Samson J. Matarao<sup>1,2,3</sup>, Camillo A. Espinosa<sup>1,2,3,4</sup>, David Seong<sup>1,5</sup>, S. Momen Retincke<sup>1,2,3</sup>, Eloise Berson<sup>1,2,6</sup>, Jonathan D. Reiss<sup>1,2</sup>, Yeasul Kim<sup>1,2,3</sup>, Marc Ghanem<sup>1</sup>, Chi-Hung Shu<sup>1</sup>, Tomlin James<sup>1</sup>, Yuqi Tan<sup>1,2</sup>, Sayane Shome<sup>1,2</sup>, Ina A. Stelzer<sup>1,2</sup>, Dorian Feyaerts<sup>1</sup>, Ronald J. Wong<sup>1,2</sup>, Gary M. Shaw<sup>1</sup>, Martin S. Angst<sup>1</sup>, Brice Gaudilliere<sup>1</sup>, David K. Stevenson<sup>1</sup> & Nima Aghaeepeour<sup>1,2,3</sup>✉

Protein features add **distinct, biologically interpretable signal** and show **structured correlations** with clinical/EHR features—explaining the gain over EHR-only models.

**d**



# Overview: EHR Data & Tasks

CTAGCTCC<sub>G...</sub>



# Electronic Health Records (EHR)

The screenshot displays the Epic EHR interface for a patient named Mickey Mouse, 14 years old. A 'New Problem' dialog box is open, showing a search for 'kni'. The search results table is as follows:

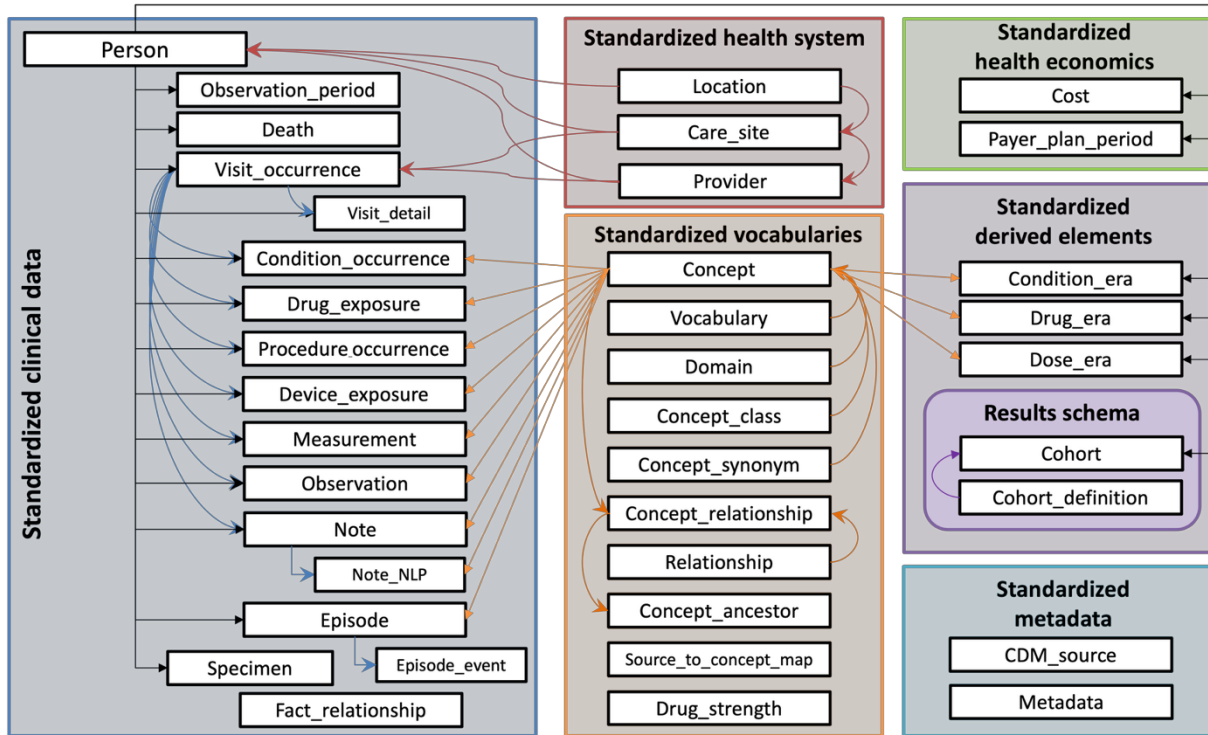
Using:	ICD-9	ICD-10
Hyperlipidemia	272.4	E78.5
GERD	530.81	K21.9
Constipation	564.00	K59.00
Knee pain	719.46	M25.569
Osteoporosis	733.00	M81.0
Hyperlipidemia NEC/NOS	272.4	E78.5
Gastroenteritis	558.9	K52.9
Aftercare, long-term use, medications NEC	V58.69	Z51.81
Dyslipidemia	272.4	E78.5
Actinic keratosis	702.0	L57.0

The dialog box also includes fields for Onset Date (8/13/2015), End Date (Select a date), and Duration (Days, Weeks, Months). Buttons for 'Add to Custom List', 'Save and Continue', 'OK', and 'Cancel' are visible at the bottom.

## Healthcare View

- GUI-based
- Data portal for a patients
- Focus on a single patient at a time

# Electronic Health Records (EHR)



## Data Scientist View

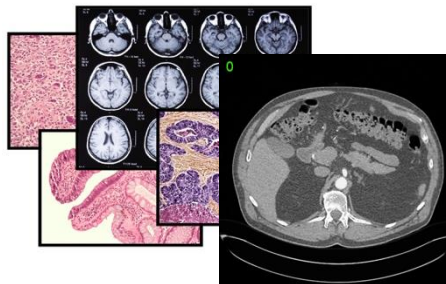
- Relational databases
- Some data model (Epic, OMOP, i2b2)
- Apply functions to all patients

# Healthcare Data is Inherently Multimodal

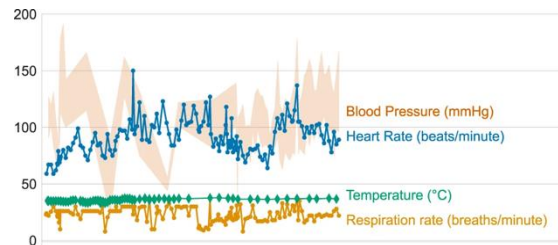


Labs	Vitals	Medication List
Notes	Past Medical History	
Problem List	Social History	
...		
Care Plan	Treatment Plan	

Tabular Data



**HISTORY OF PRESENT ILLNESS:**  
60 yo male with **infected R hip** (MRS  
**LTHA** **November 2004** demonstrates  
**HISTORICAL** **>2 YEARS**  
No **lucencies** were observed around  
**NEGATED**  
**Implant** is being evaluated for possi



Audio / Conversations



Video



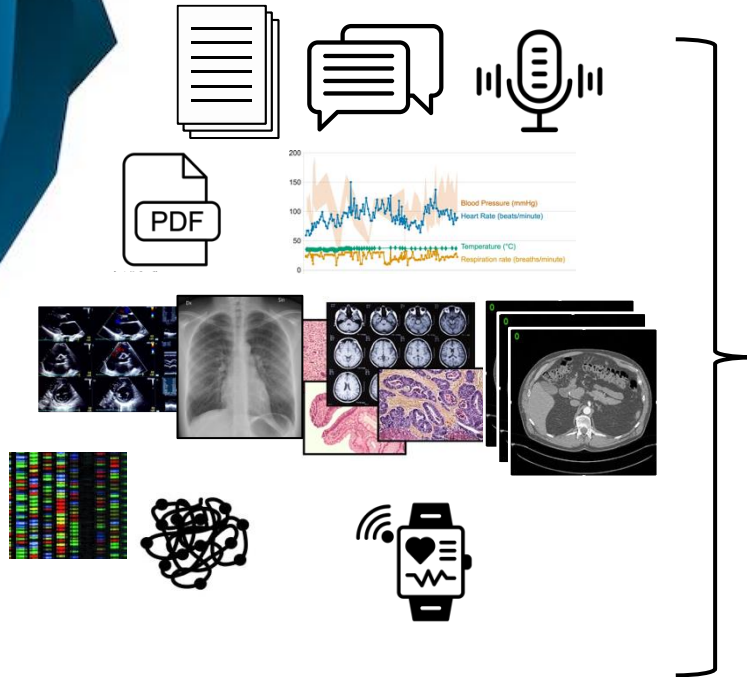
Genomics

STRUCTURED DATA

UNSTRUCTURED DATA

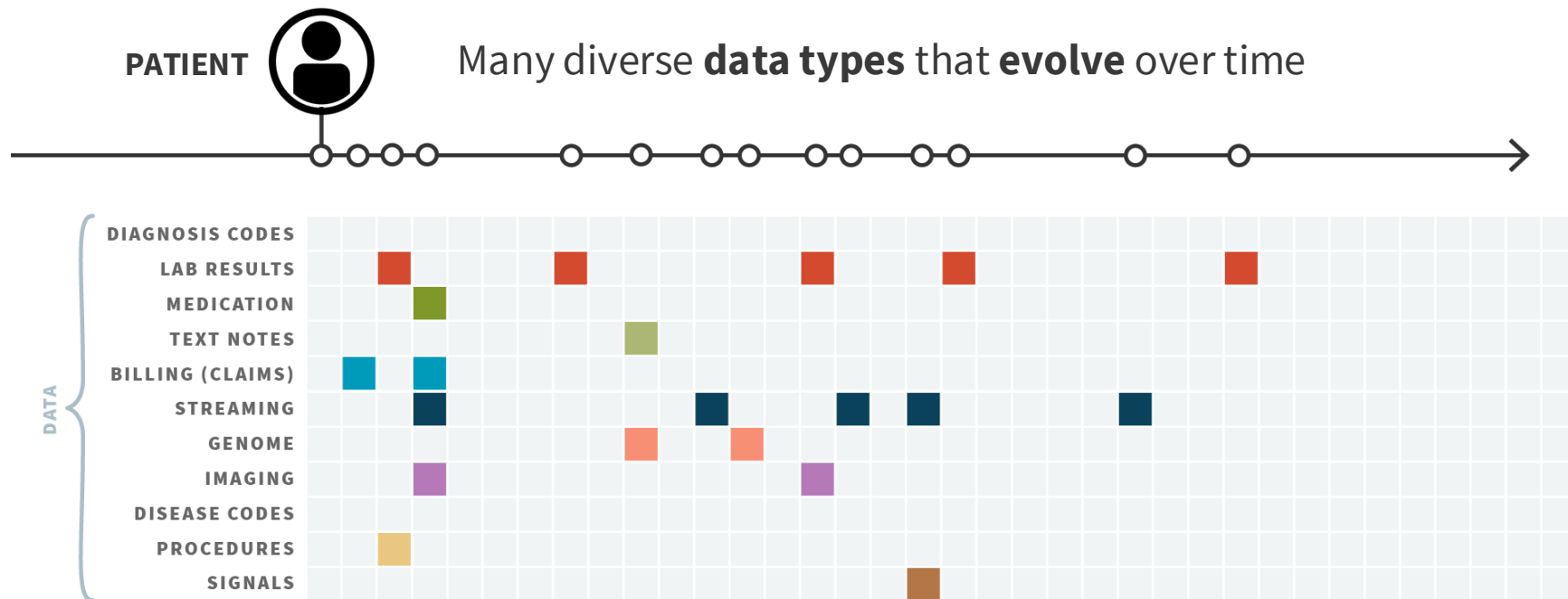
# Hospital data is growing at a rate of **36% per year**

World Economic Forum, Dec. 2019



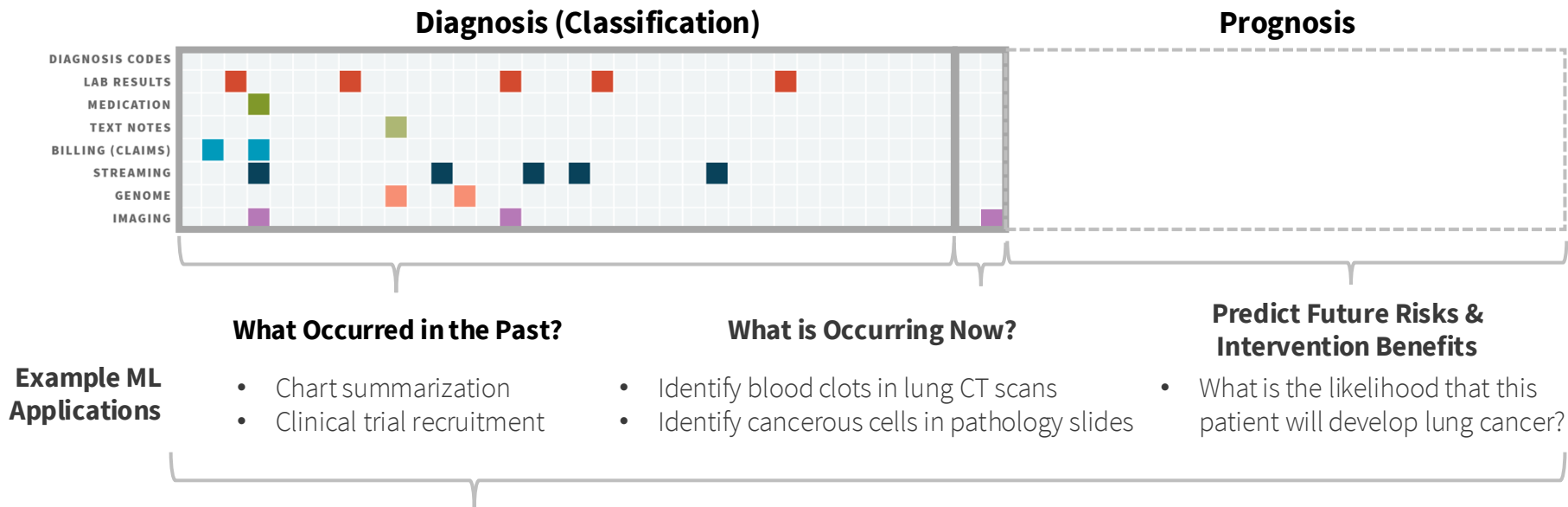
**Hard to use for medical  
decision making**

# Electronic Health Records (EHRs) are Multimodal Timelines



Longitudinal EHRs provide a **holistic view of multimodal data**

# AI for Healthcare Requires Temporal Reasoning



Stakeholders



Clinicians

Whether to Treat

How to Treat

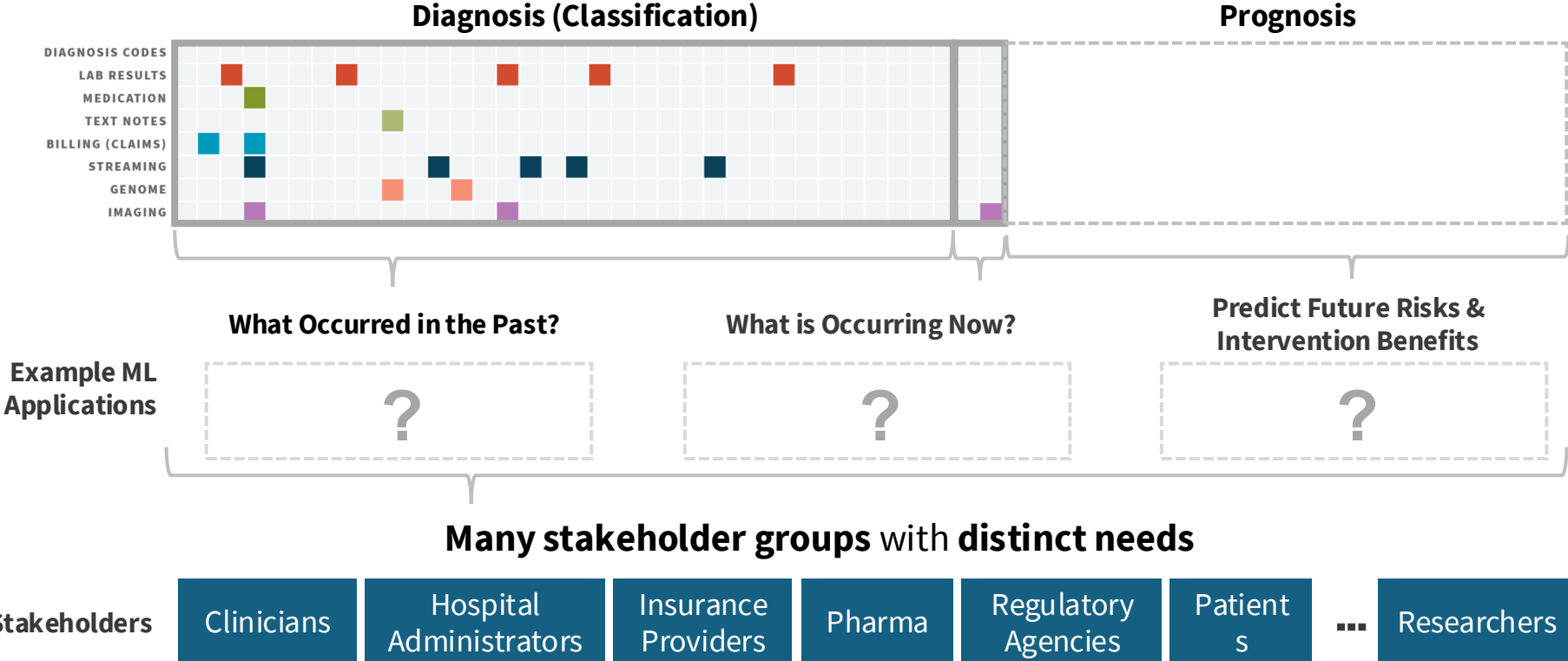
subject  
to

Policy

Capacity to Act

Intervention Properties

# Foundation Models Are Essential for AI in Healthcare



# How Can AI Improve Healthcare?

## Atherosclerotic cardiovascular disease risk assessment: An American Society for Preventive Cardiology clinical practice statement



Nathan D. Wong<sup>a,\*</sup>, Matthew J. Budoff<sup>b</sup>, Keith Ferdinand<sup>c</sup>, Ian M. Graham<sup>d</sup>, Erin D. Michos<sup>e</sup>, Tina Reddy<sup>c</sup>, Michael D. Shapiro<sup>f</sup>, Peter P. Toth<sup>e,g</sup>

---

**nature medicine**



Article

<https://doi.org/10.1038/s41591-023-02332-5>

## **A deep learning algorithm to predict risk of pancreatic cancer from disease trajectories**

# A Sketch of Healthcare Tasks

- **Improved patient outcomes**

- Treatment selection
- Disease diagnosis (e.g. early detection of cancer)
- Risk stratification (e.g. mortality, cancer progression)
- Abnormal test result prediction (e.g. lab values)

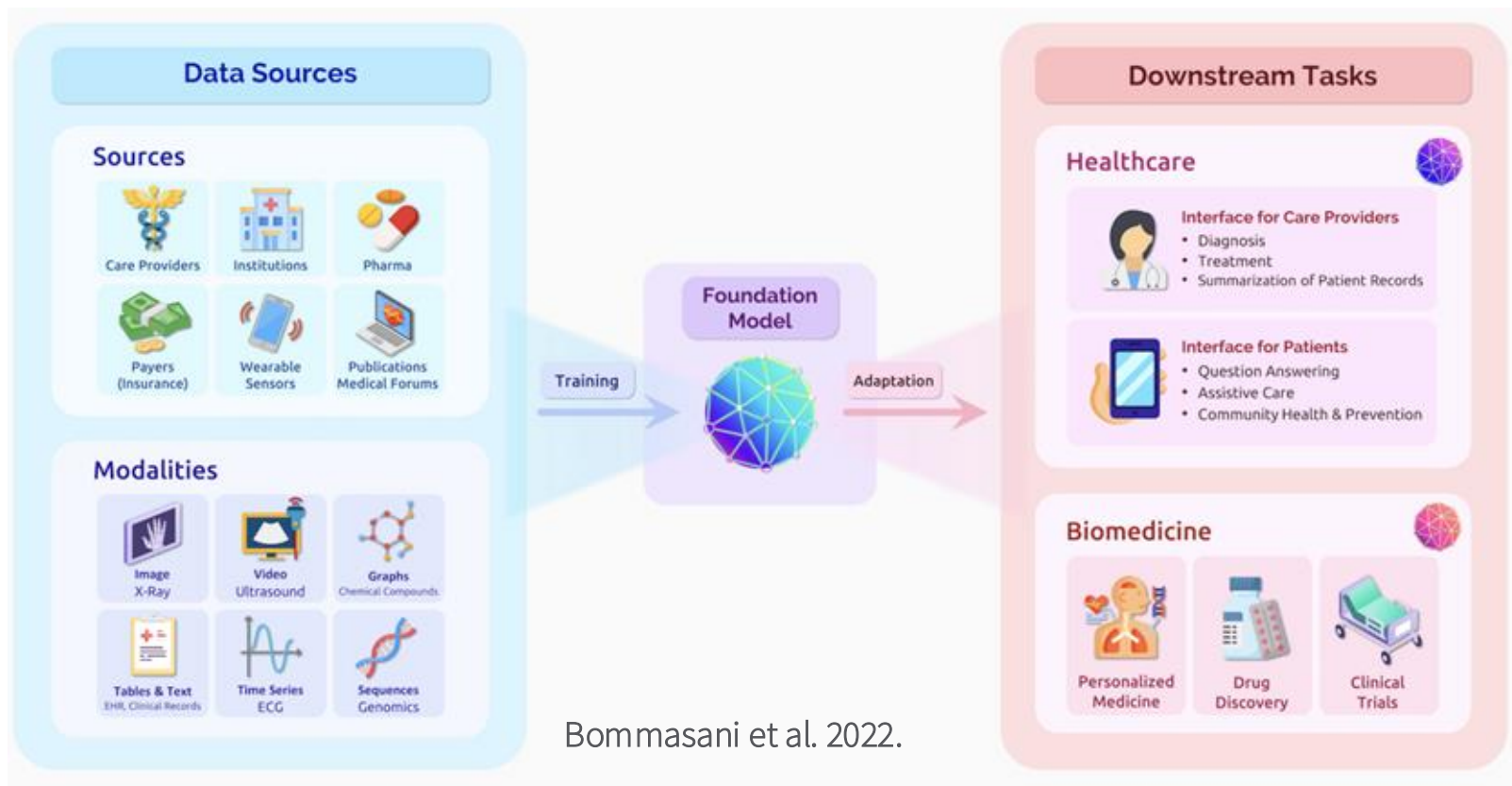
- **More efficient hospital operations**

- Predictions for quality metrics (e.g. 30-day readmission likelihood)
- Resource allocation (e.g. anticipating ICU transfers)
- Billing (e.g. identify mis-coding of patient records)

- **Research**

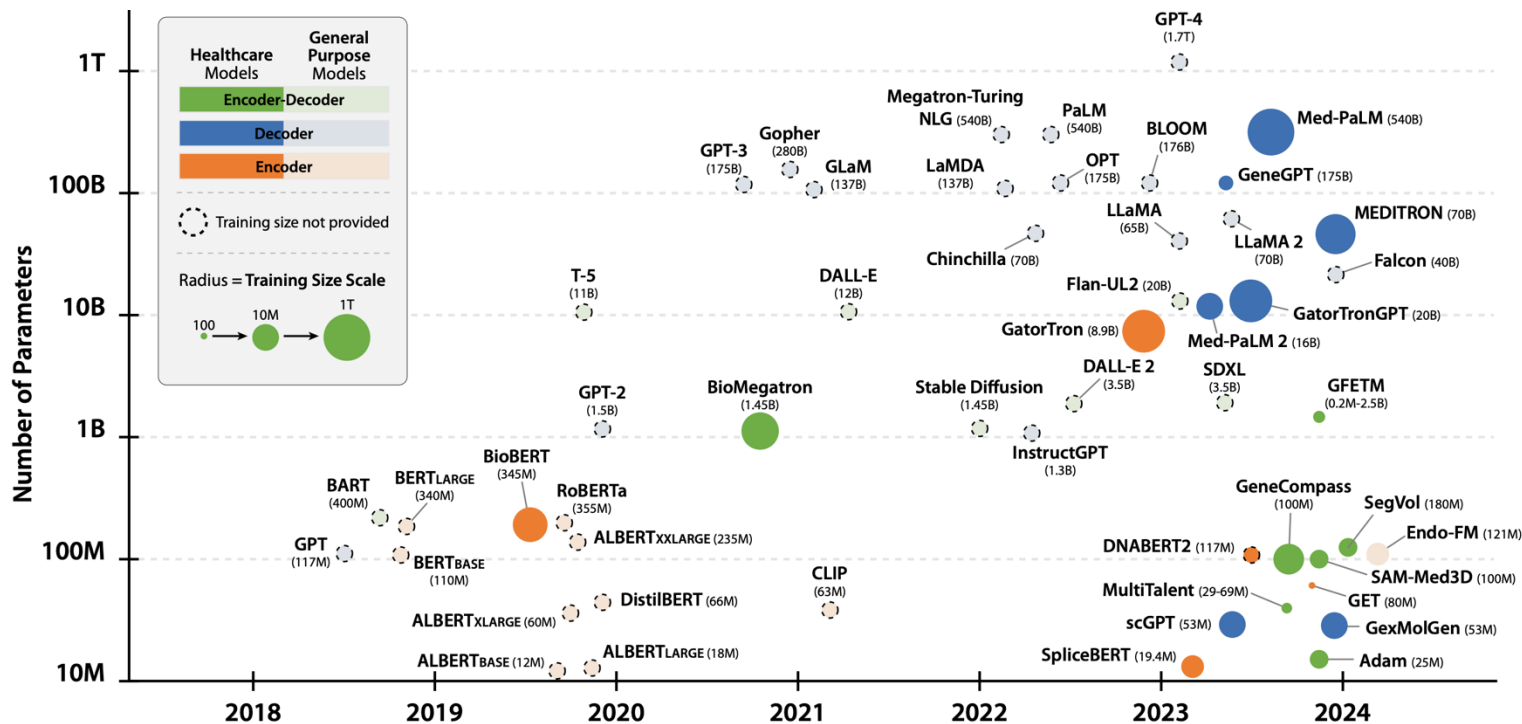
- Causal inference (e.g. drug trials and observational studies)
- Identify off-label drug benefits

# Foundation Models and AI's “Industrial Age”



Bommasani et al. 2022.

# Opportunity for AI to reimagine how we interact and understand medical data



Khan et al., “A Comprehensive Survey of Foundation Models in Medicine,” 2025.



# Modeling: Pretraining Objectives

CTAGCTCC<sub>G...</sub>



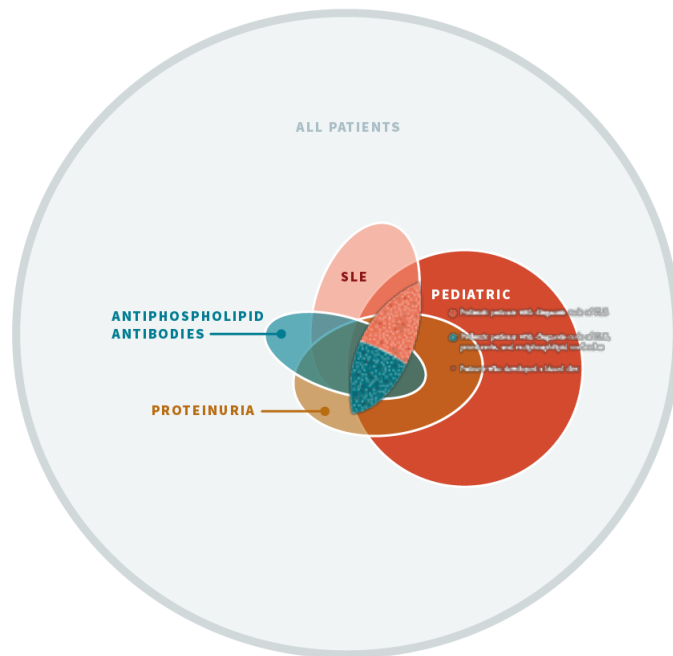
Small text describing the pretraining objectives, likely related to the modeling of biological data.

# Classic Approach to Building and Patient Model

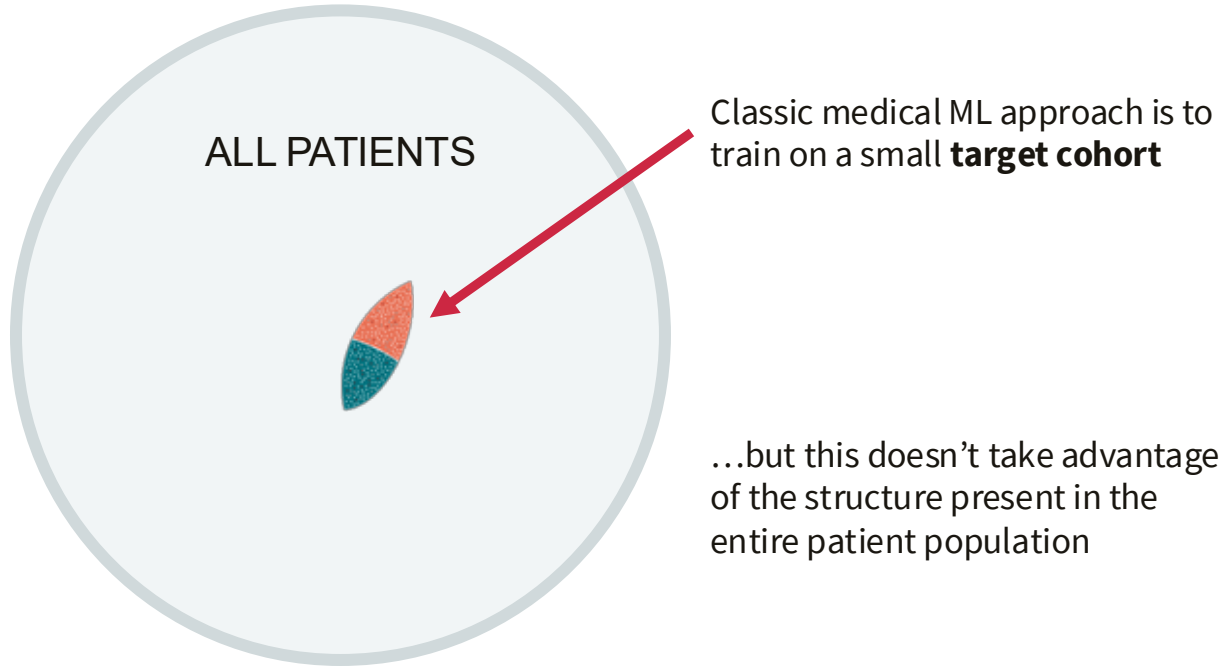


## MEET LAURA

A teenager with systemic lupus erythematosus (SLE), proteinuria, pancreatitis and positive for antiphospholipid antibodies

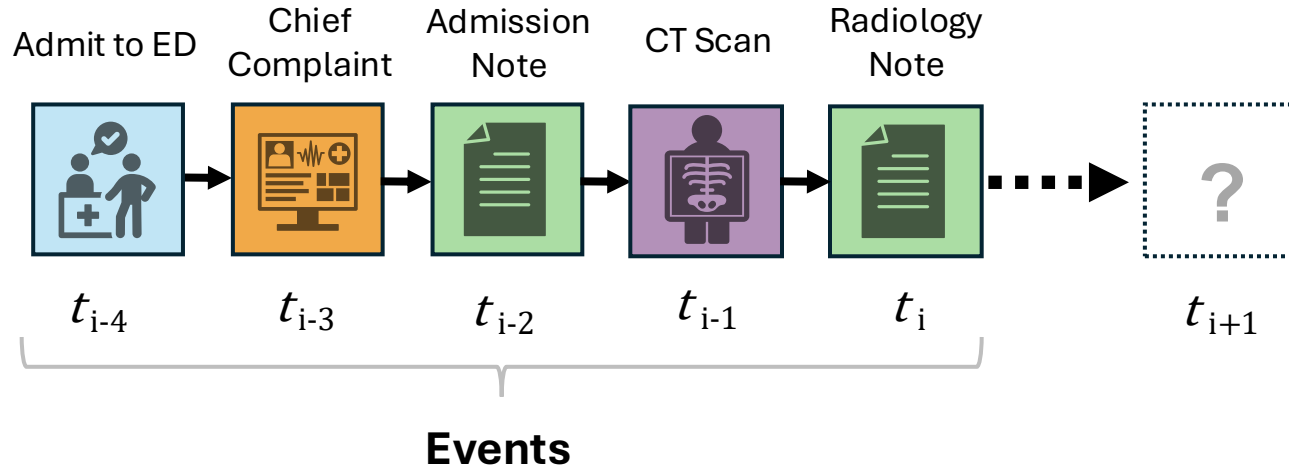


# Classic Approaches Often Fail Due to Limited Data



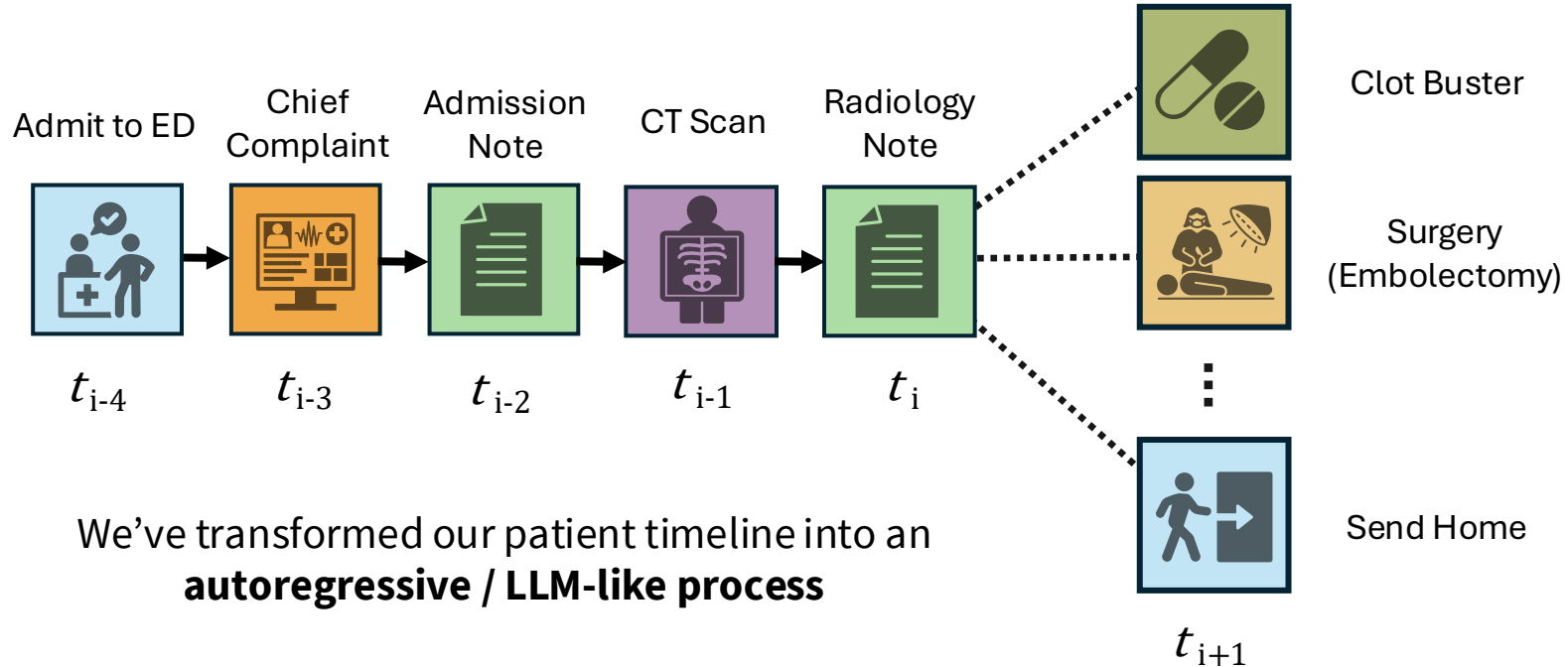
# Modeling Patient Timelines for AI

PATIENT CASE: Patient **presents to ED** with sudden onset **shortness of breath**, **pleuritic chest pain**, and **tachycardia**. Concern for **pulmonary embolism**.



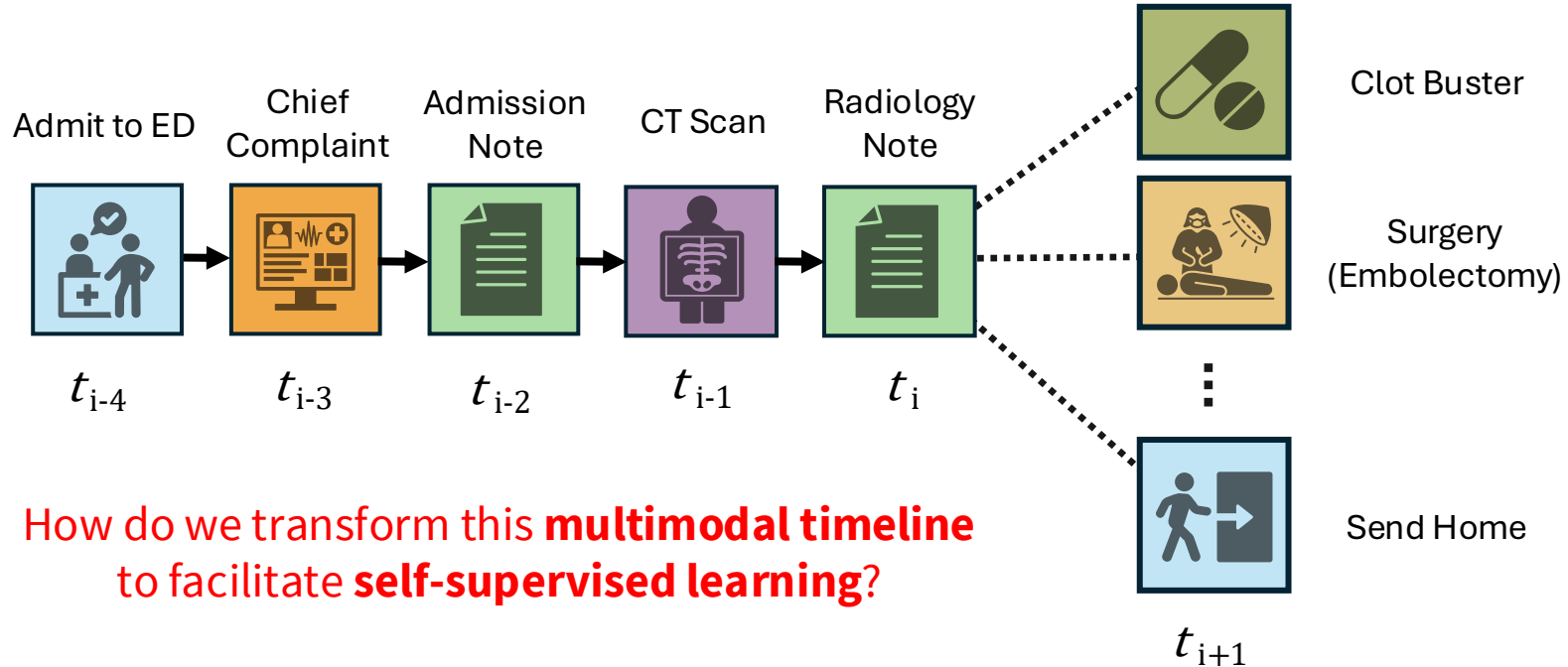
# Modeling Patient Timelines for AI

PATIENT CASE: Patient **presents to ED** with sudden onset **shortness of breath**, **pleuritic chest pain**, and **tachycardia**. Concern for **pulmonary embolism**.



# Modeling Patient Timelines for AI

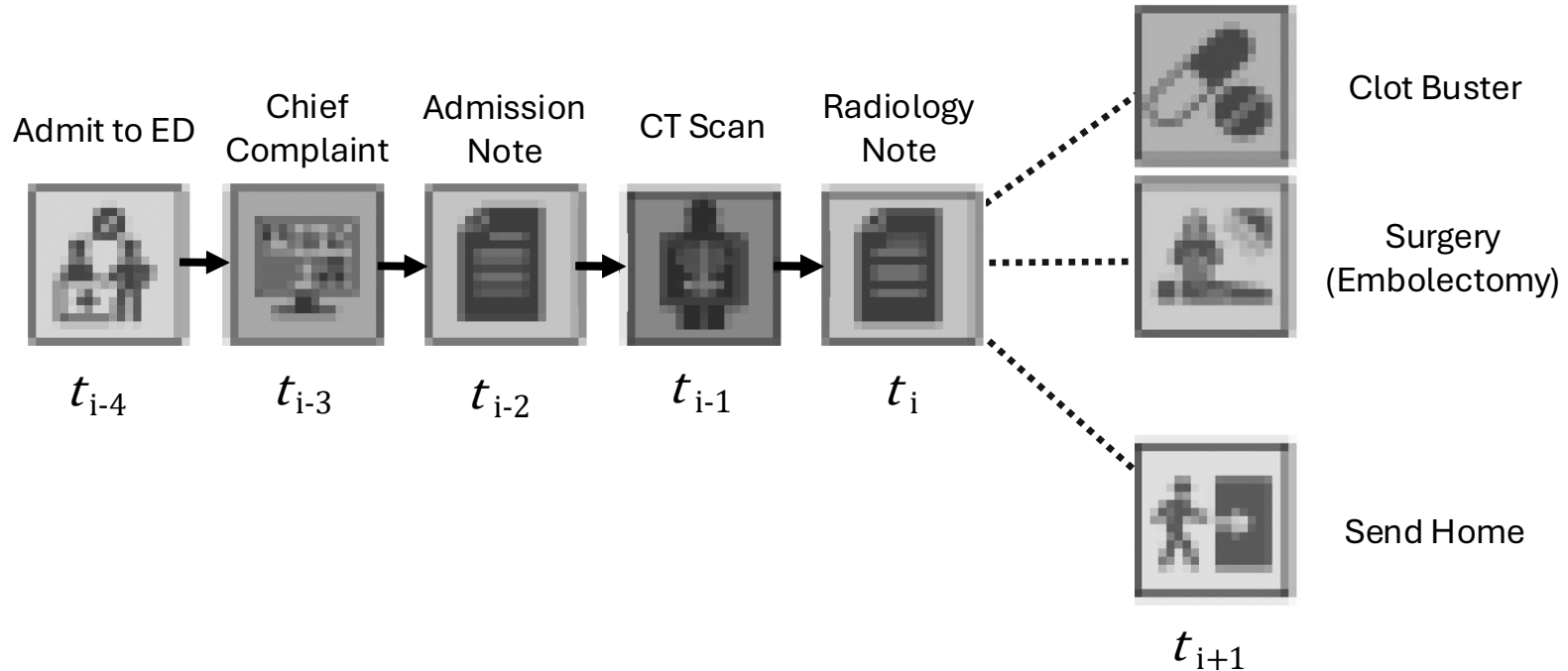
**Hypothesis:** A model that accurately **predicts future health states**, based on patient history, **encompasses many proposed use cases of medical AI**



How do we transform this **multimodal timeline** to facilitate **self-supervised learning**?

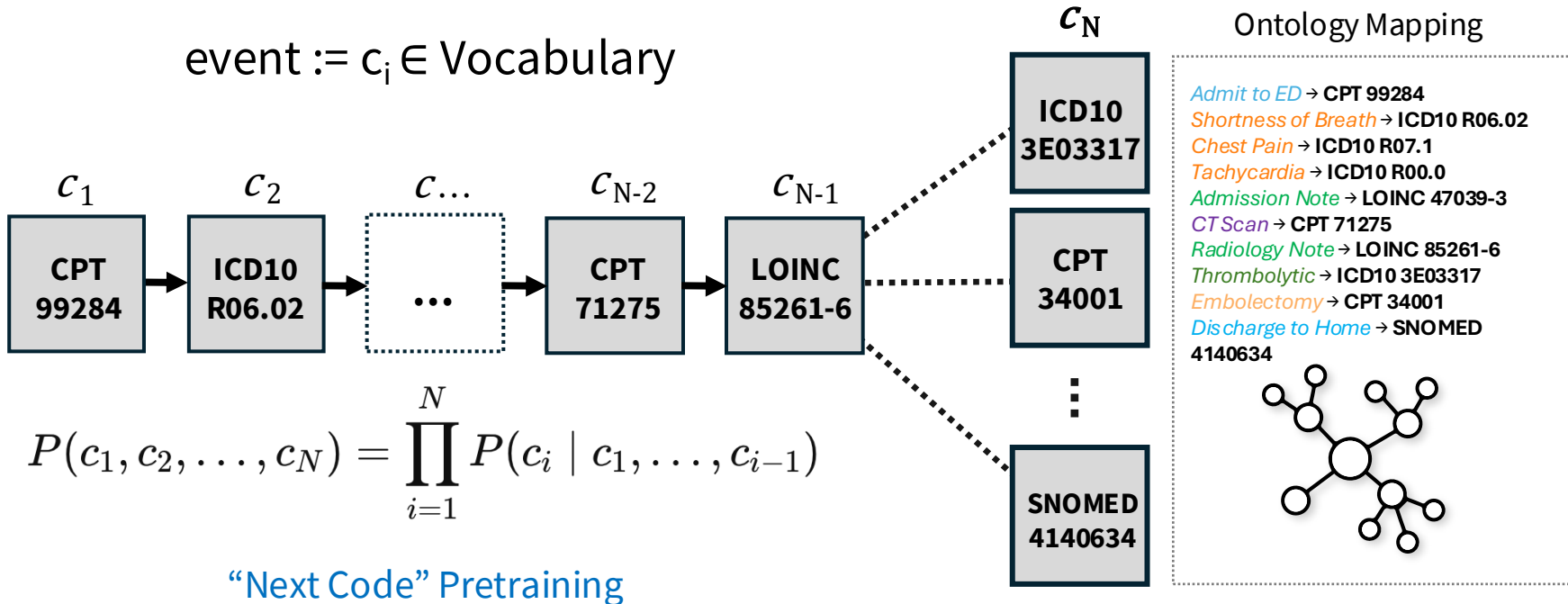
# Modeling Patient Timelines for AI

We do have a “**low-res**” version of this timeline readily available ...



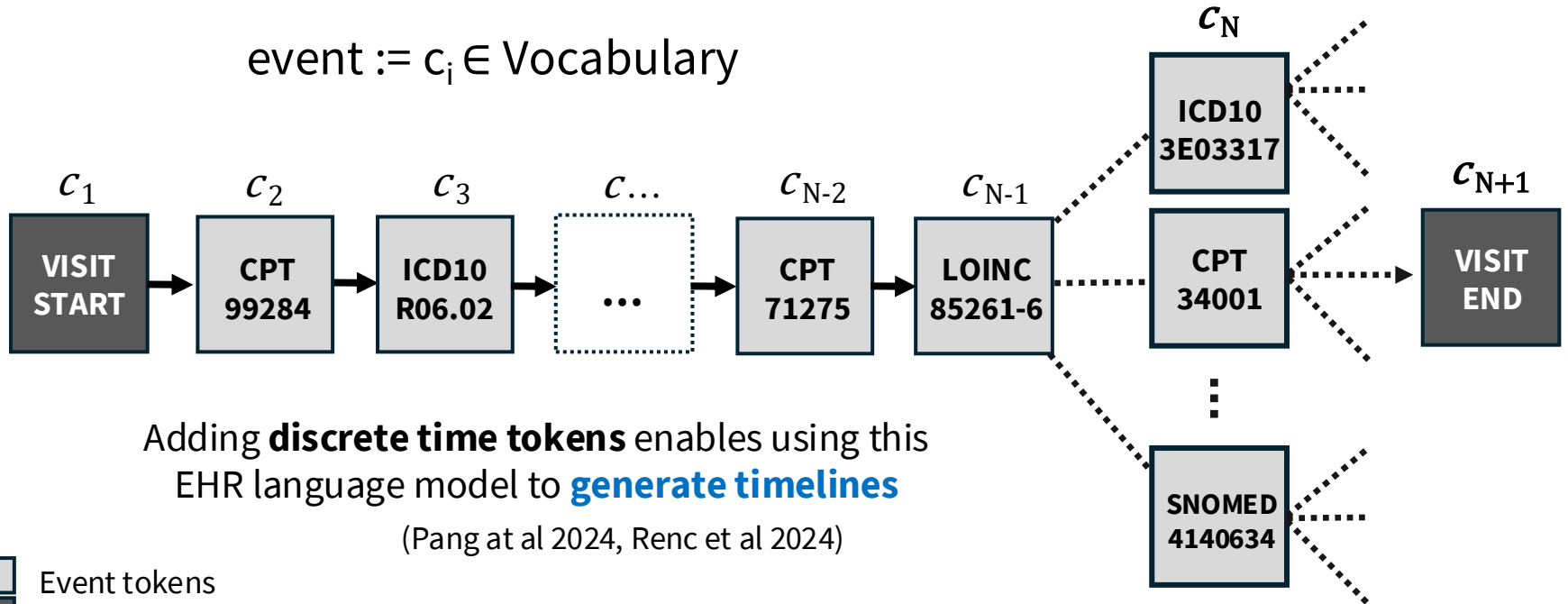
# Modeling Structured EHR Timelines

**Map events to ontologies** to define a “language” based on medical codes



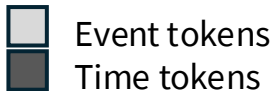
# Modeling Structured EHR Timelines

**Map events to ontologies** to define a “language” based on medical codes



Adding **discrete time tokens** enables using this EHR language model to **generate timelines**

(Pang et al 2024, Renc et al 2024)



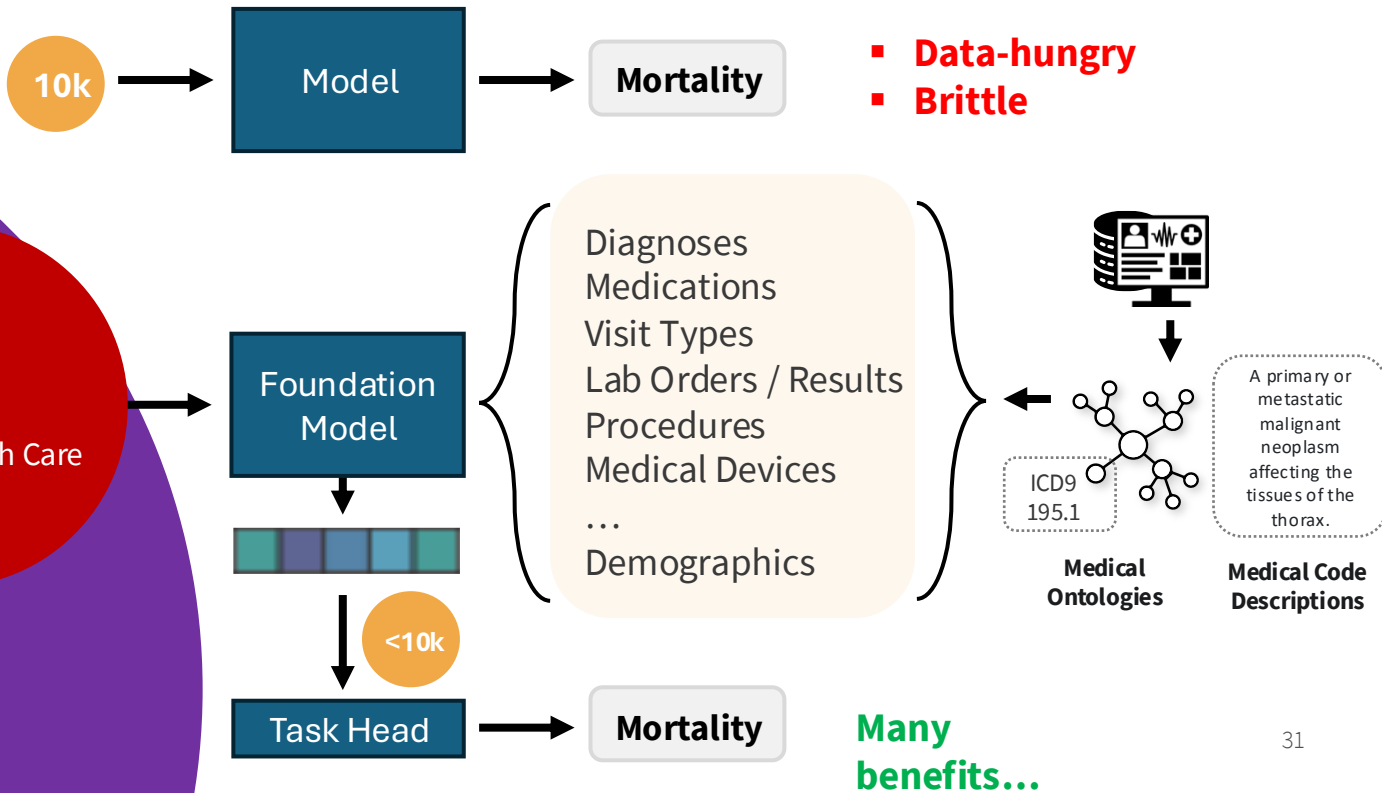
# Self-Supervised Training of an EHR Foundation Model

## PATIENT POPULATION

1,699 Hospitals  
**296M**  
Epic Cosmos

**2.57M**  
Stanford Health Care

## TASKS



# Self-Supervised Pretraining Objectives for Structured Event Data

## BERT-Style (Masked Language Modeling)

- BEHRT (Li et al. 2020)
- MedBERT (Rasmy et al. 2021)
- CEHR-BERT (Pang et al 2021)
- ClaimPT (Zeng et al. 2022)
- *et alia*

## GPT-Style (Autoregressive)

- CLMBR (Steinberg et al. 2020)
- TransformEHR (Yang et al. 2023)
- CEHR-GPT (Pang et al 2024)
- ETHOS (Renc et al. 2024)

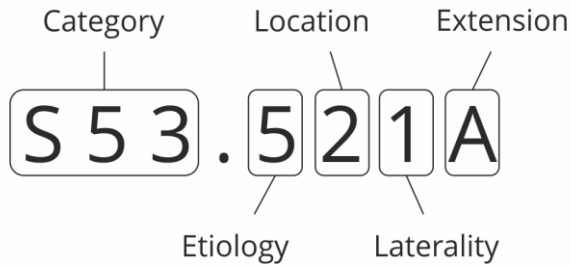
## Time-to-Event

- MOTOR (Steinberg et al. 2024)

**Won't talk about masked language modeling**  
**Will focus on structured (medical code) models**

# Structured Data: Medical Vocabularies

## ANATOMY OF AN ICD-10 CODE



ICD-10 code for torus fracture of lower right end of right radius, initial encounter for closed fracture

<https://blogs.halodoc.io/>

- Controlled Vocabularies
- **Knowledge Graphs**

code<sub>i</sub> ∈ Vocabulary

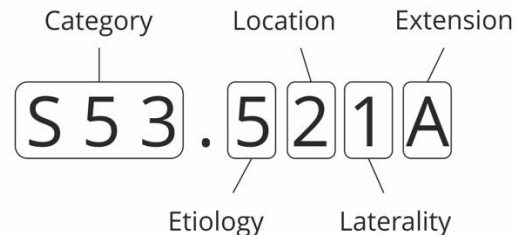


Category or Name		
- {component} 103832		
- Laboratory 63121		
+ Microbiology and Antimicrobial susceptibility 5731		
+ Skin challenge 47		
- Chemistry and Chemistry - challenge 14248		
+ Chemistry - non-challenge 10420		
- Chemistry - routine challenge 27		
+ 17-Hydroxypregnenolone 2		
+ Cortisol 7		
+ Dehydroepiandrosterone 1		
- Glucose 17		
- Glucose   Blood   Chemistry - routine challenge 3		
	Glucose p meal Bld-mCnc	Glucose^post meal
	Deprecated Glucose pre-meal Bld-mCnc	Glucose^pre-meal
	Glucose pre-meal Bld-mCnc	Glucose^pre-meal

# More Like NLP Now, but Key Differences!

## Tokenization / Vocabulary

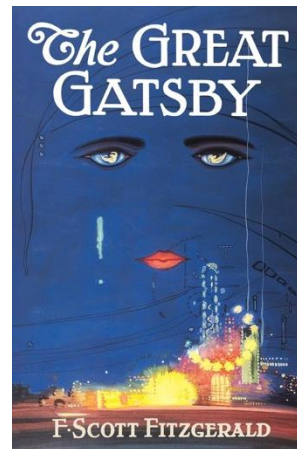
	<b>NLP</b>	<b>EHR</b>
Vocabulary Size	50k	<b>250k+</b>
Subwords	Yes	<b>No</b>
Tokens Semantics	Flat	



***Hierarchical, Complex Dependencies***

## Sequence Properties

	<b>NLP</b>	<b>EHR</b>
Sequence Length	32k	<b>250k+</b>
Ordering	Total	<b><i>Partial</i></b>
Time Intervals	None	<b><i>Discontinuous</i></b>
Sampling Fidelity	All	<b><i>Sparse/Errors</i></b>

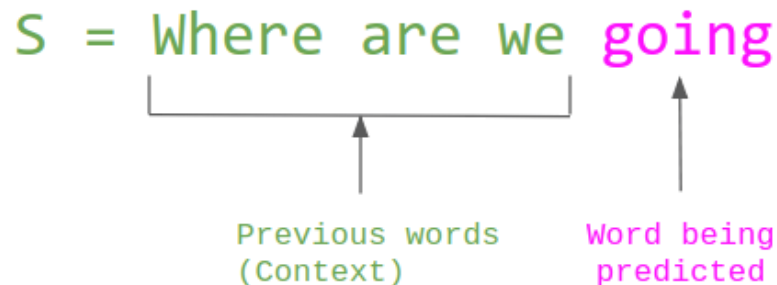


**50% Patients**  
**>= 68k tokens**

# GPT-Style (Autoregressive)

- CLMBR (Steinberg et al. 2020)
- TransformEHR (Yang et al. 2023)
- CEHR-GPT (Pang et al 2024)
- ETHOS (Renc et al. 2024)

# Self-Supervised Pretraining in Natural Language

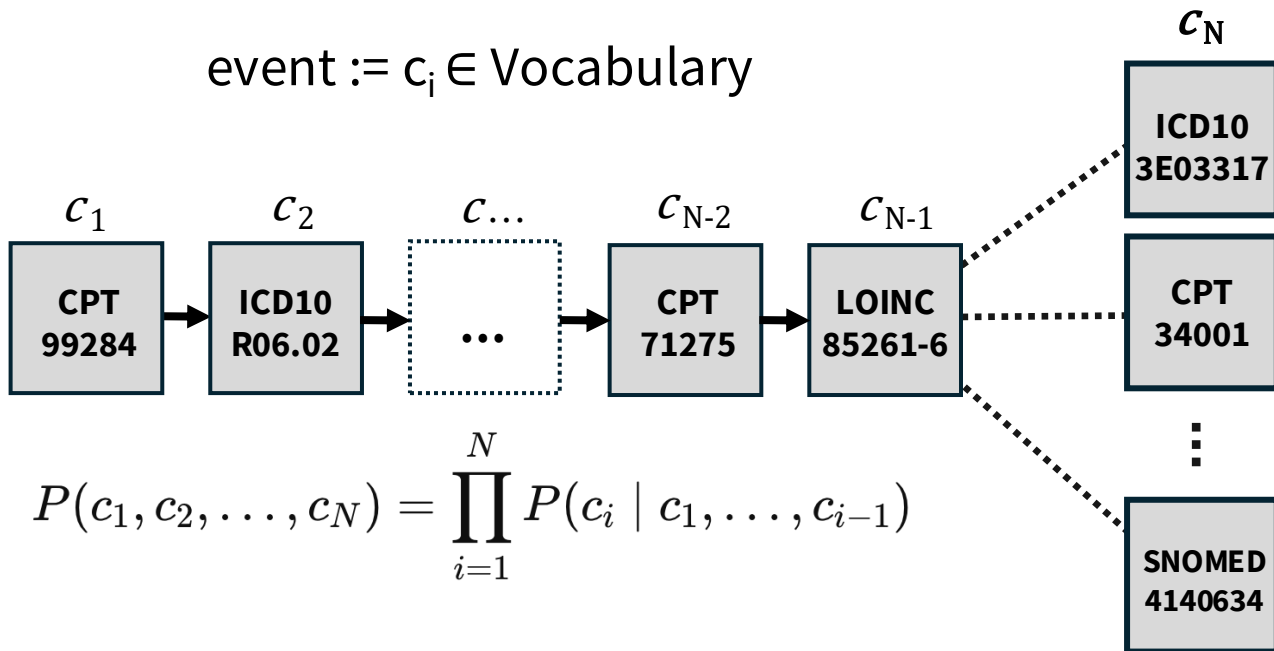


$$P(S) = P(\text{Where}) \times P(\text{are} \mid \text{Where}) \times P(\text{we} \mid \text{Where are}) \times P(\text{going} \mid \text{Where are we})$$

$$\begin{aligned} P_{(w_1, w_2, \dots, w_n)} &= p(w_1)p(w_2|w_1)p(w_3|w_1, w_2)\dots p(w_n|w_1, w_2, \dots, w_{n-1}) \\ &= \prod_{i=1}^n p(w_i|w_1, \dots, w_{i-1}) \end{aligned}$$

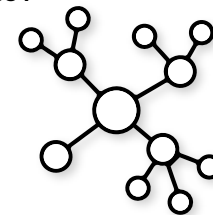
# Next Code Pretraining

event :=  $c_i \in \text{Vocabulary}$

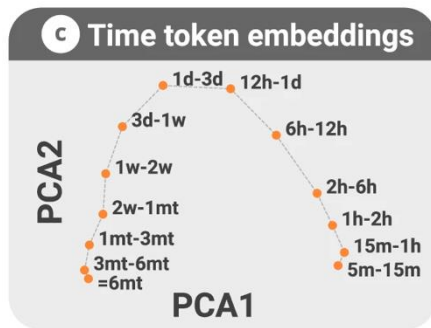
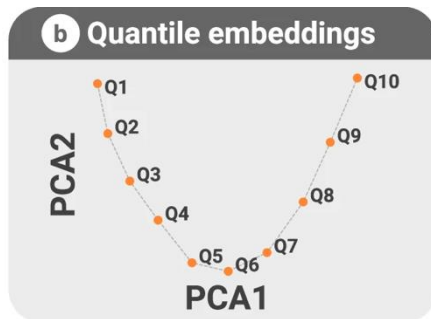
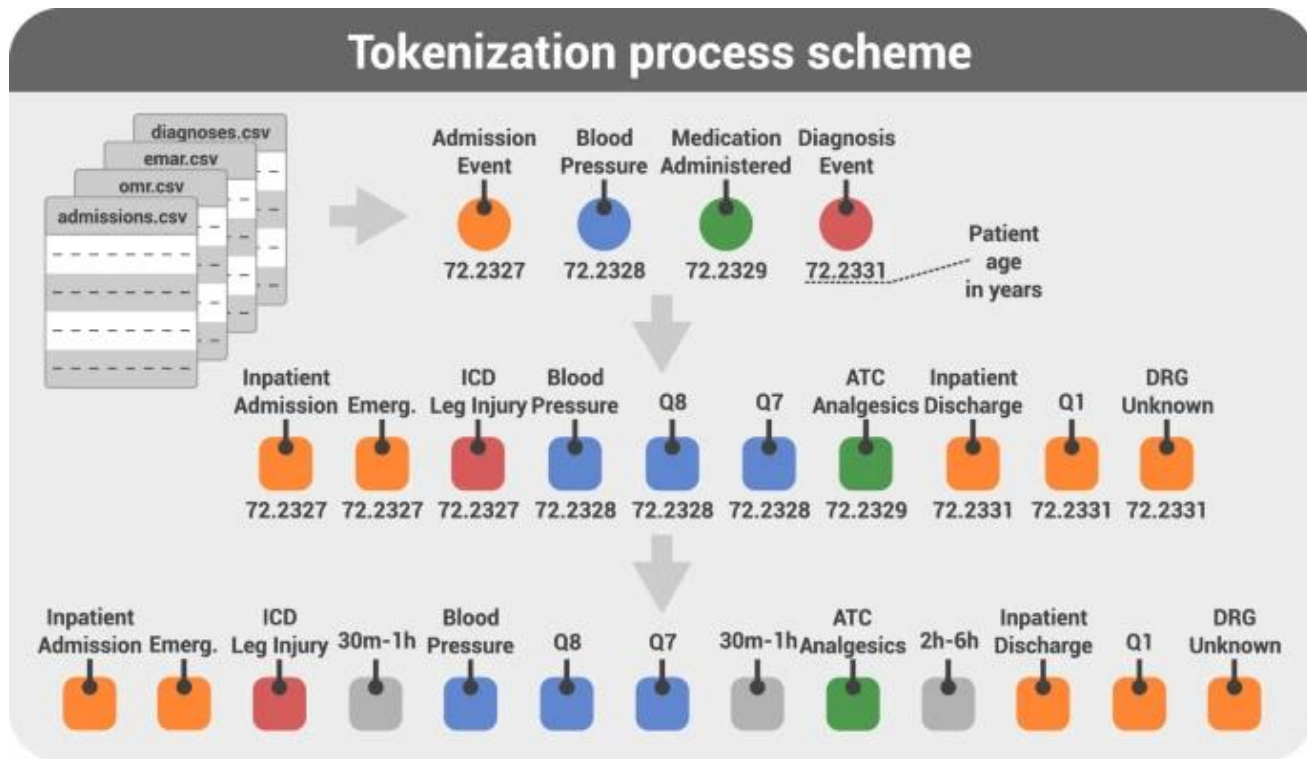


Ontology Mapping

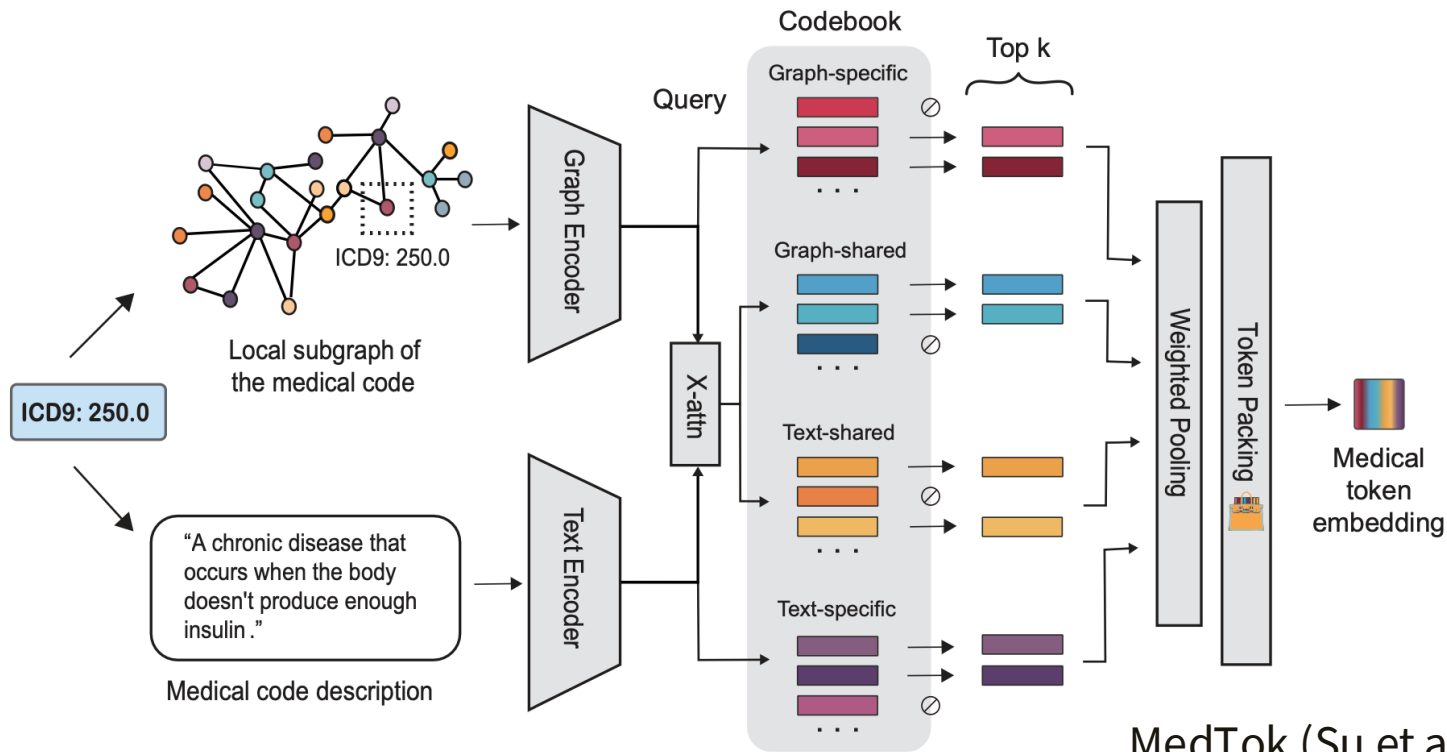
*Admit to ED* → CPT 99284  
*Shortness of Breath* → ICD10 R06.02  
*Chest Pain* → ICD10 R07.1  
*Tachycardia* → ICD10 R00.0  
*Admission Note* → LOINC 47039-3  
*CT Scan* → CPT 71275  
*Radiology Note* → LOINC 85261-6  
*Thrombolytic* → ICD10 3E03317  
*Embolectomy* → CPT 34001  
*Discharge to Home* → SNOMED 4140634



# Tokenization



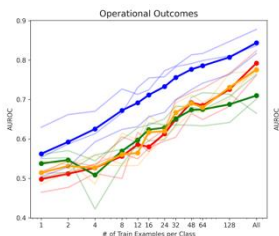
# Generalized Tokenizer



MedTok (Su et al. 2025)  
**Drop-in Replacement**

# Validating Benefits of EHR Foundation Models

## Data Efficiency



SOTA **few-shot learning**  
SOTA **overall performance**

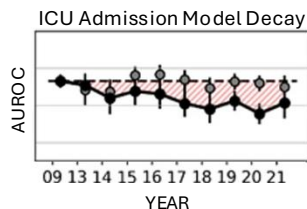
(Wornow et al. 2023)

(Steinberg et al. 2020)

### Publication Venue

- Medical / Informatics
- Computer Science

## Robustness



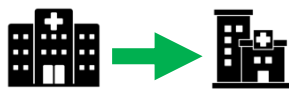
Improved robustness to **temporal distribution shifts**

(Guo et al. 2023)

Improved performance across key **subgroups** (pediatrics)

(Lemmon et al. 2023)

## Cross-Site Adaptability



Hospital A

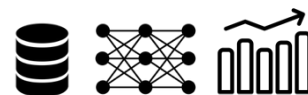
Hospital B

Transfer **pretrained models** across hospitals

Require **up to 90% less** pretraining data

(Guo et al. 2024)

## Reproducible EHR Benchmarking



First **externally verifiable** evaluation of **EHR foundation models** on longitudinal data

(Wornow et al. 2025)

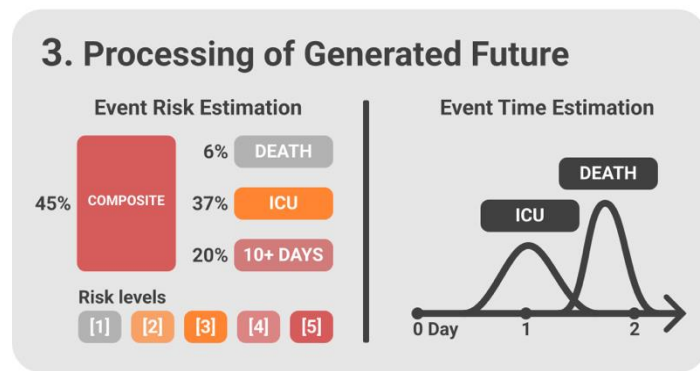
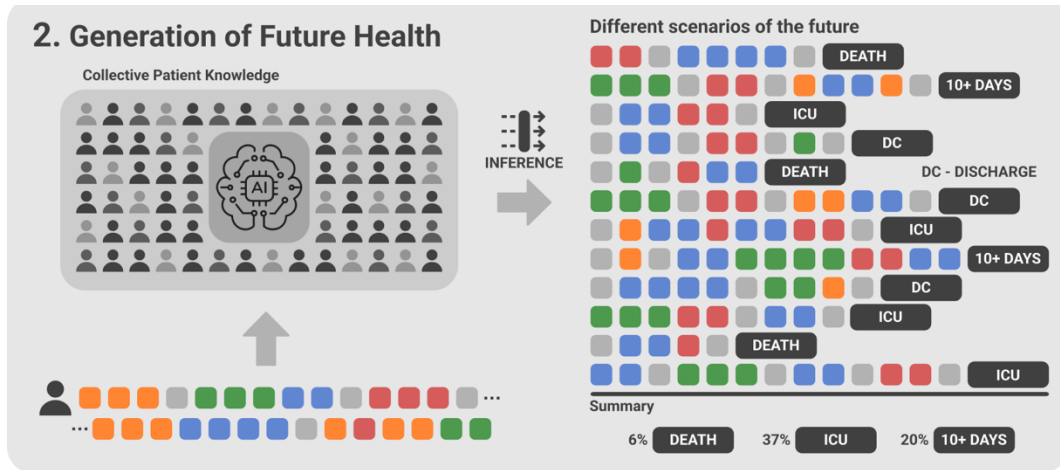
(Arrnich et al. 2024)

(Steinberg et al. 2024)

(Wornow et al. 2023)

(Huang et al. 2023)

# Zero-Shot Patient Classification



ETHOS **samples from model rollouts** to estimate future event risk

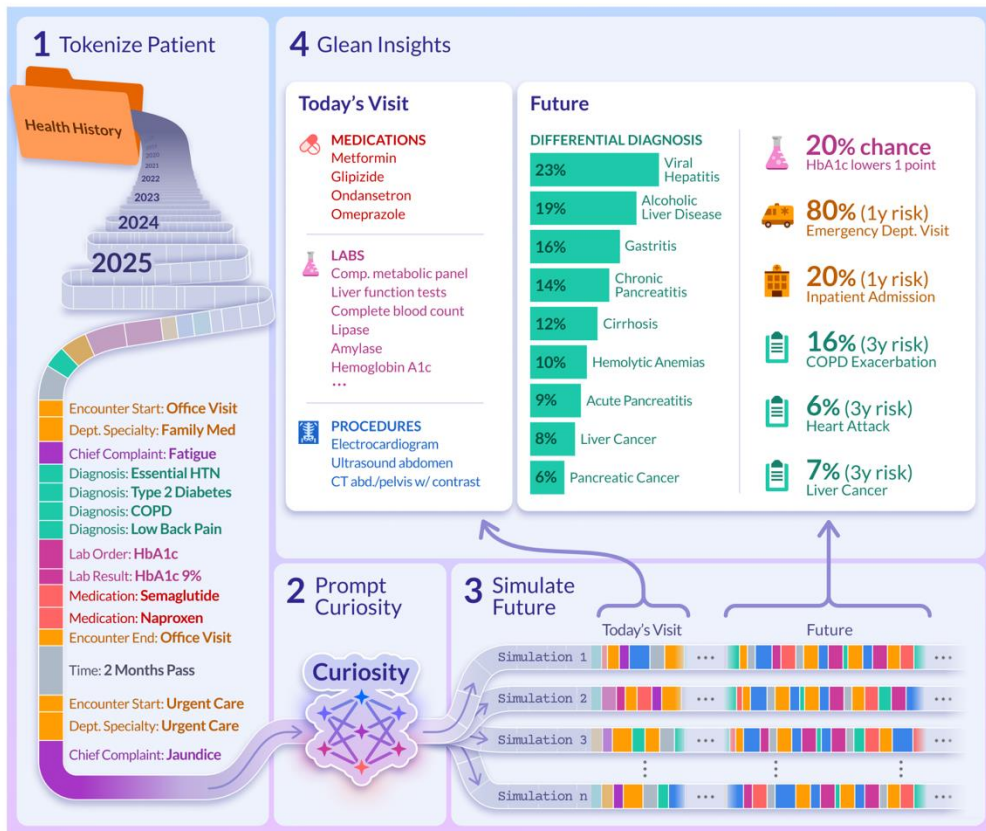
## Monte Carlo Sampling

$x$  := sampled rollout     $y$  := outcome indicator

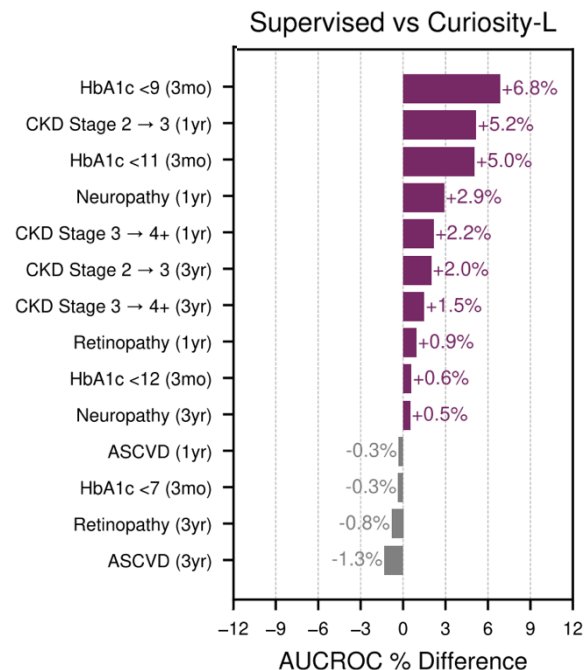
$$\hat{p}_{\text{MC}} = \frac{1}{n} \sum_{i=1}^n y^{(i)}$$

$$y^{(i)} = \mathbf{1}\{\text{outcome token } O \text{ occurs in rollout } x^{(i)} \text{ before stop}\}.$$

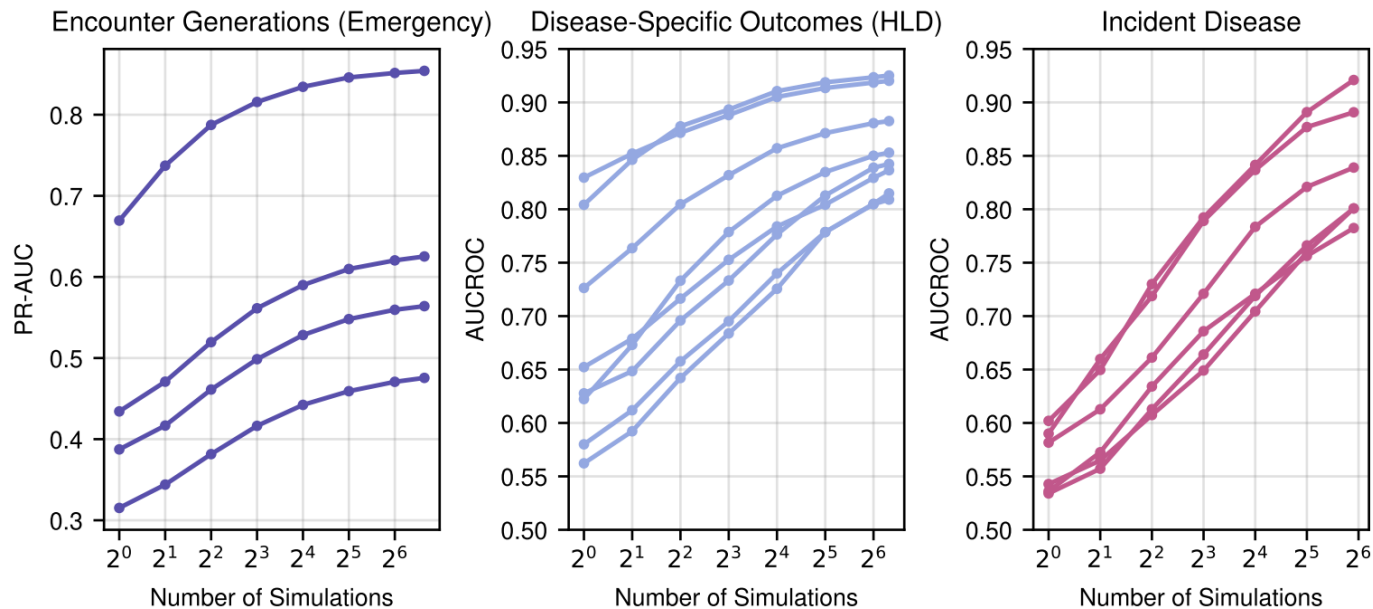
# Epic's Curiosity Model



**300 million** unique patient records from **310 health systems**



# Epic's Curiosity Model: Test-Time Compute



**Figure 17: Effect of test-time compute on performance.** For Curiosity-L we evaluated model performance against the number of simulations generated, focusing on single-encounter generations for emergency visits (**left**), HLD-specific outcomes (**middle**), and incident disease tasks (**right**). Each line represents a specific task (e.g., one-year ASCVD risk, two-year COPD, etc.) from the titled category. For readability the legend is not included.

# **Time-to-Event Modeling**

# Data (Label) Efficiency of EHR Foundation Models

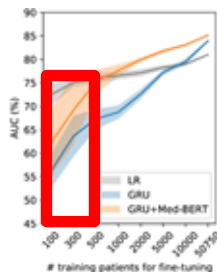
**Label Efficiency:** How many **labeled examples** are needed to train a high-performing model?

## BERT-Style (Masked Language Modeling)

- BEHRT (Li et al. 2020)
- MedBERT (Rasmy et al. 2021)
- CEHR-BERT (Pang et al 2021)
- ClaimPT (Zeng et al. 2022)
- *et alia*

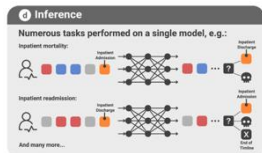
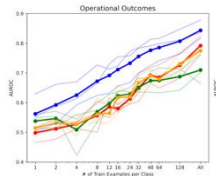
## GPT-Style (Autoregressive)

- CLMBR (Steinberg et al. 2020)
- TransformEHR (Yang et al. 2023)
- CEHR-GPT (Pang et al 2024)
- ETHOS (Renc et al. 2024)



## MedBERT

- Trained on **28M patients**
- Performance with **< 500 examples worse than logistic regression**



## CLMBR

- Trained on **2.57M patients** (3.5B tokens)
- SOTA **few-shot** learning using **embeddings**

## ETHOS

- Trained on **200k patients** (MIMIC-VI)
- **Zero-shot** abilities using **generation**

# Autoregressive Modeling at Smaller Scales

Autoregressive LLMs can capture long-distance dependencies given **sufficient data and parameters**

**Natural Language**

≥ 7B parameters

≥ 500B-1T tokens

**EHR**

143M parameters

3.5B tokens

**285x**  
less data

Can we train a **small, data-constrained** EHR foundation model to learn embeddings that capture more information about the future?



# Key Concepts in Time-to-Event Modeling

Model the **time until an event occurs** (e.g., death) while accounting for **censoring**

## Censoring

Event times are **not fully observed by end of a study period**

$$\boxed{(X_i, T_i)} \underset{D}{\text{BIASE}} (X_i, T_i, \delta_i) \quad \delta_i = \begin{cases} 1 & \text{event observed} \\ 0 & \text{censored} \end{cases}$$

## Survival Function

The probability that an event has not occurred as of time  $t$

$$S(t) = \Pr(T > t)$$

## Hazard Rate Function

Instantaneous risk of an event at time  $t$ , given survival up to  $t$

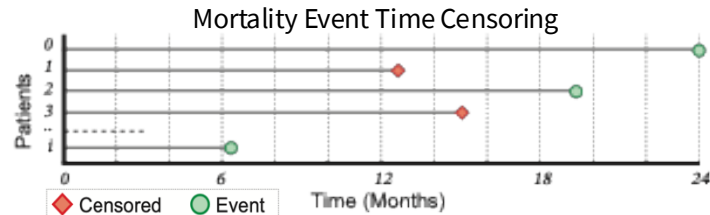
$$h(t) = \lim_{\Delta t \rightarrow 0} \frac{\Pr(t \leq T < t + \Delta t \mid T \geq t)}{\Delta t}$$

*Event's "speed" at each moment*

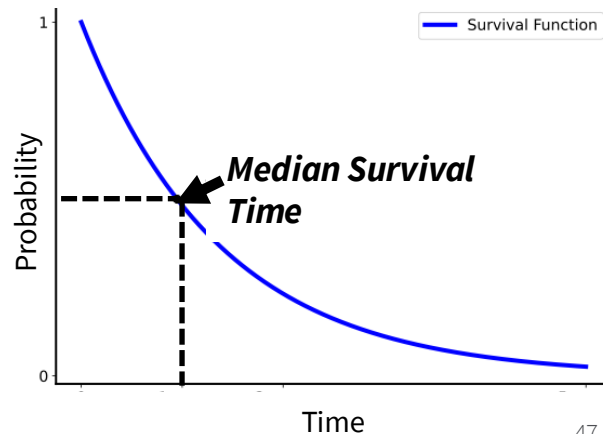
$$S(t) = \exp\left(-\int_0^t h(u) du\right)$$

*Survival depends on cumulative hazard over time*

Learn a patient representation  $R_i = f_\theta(X_i)$  for estimating **personalized hazard rates**

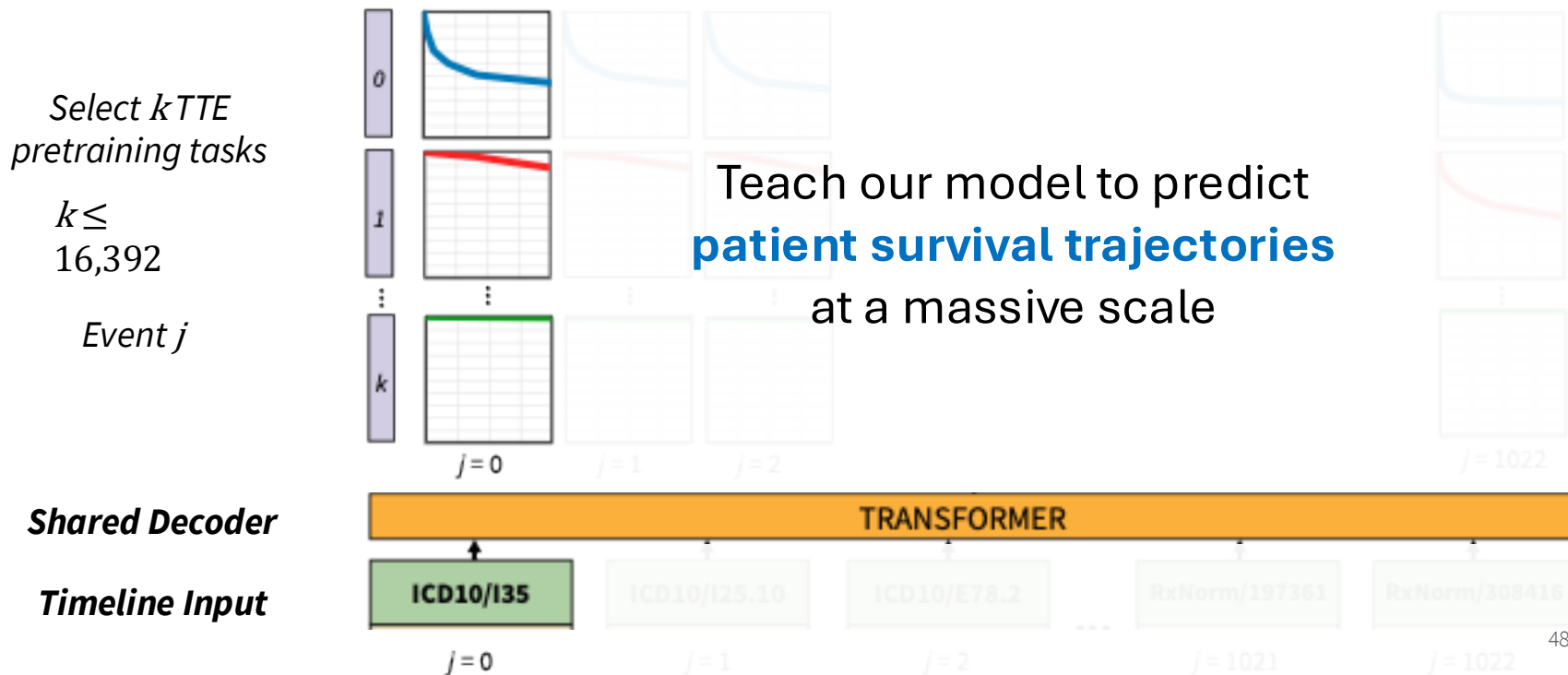


## Survival Curve



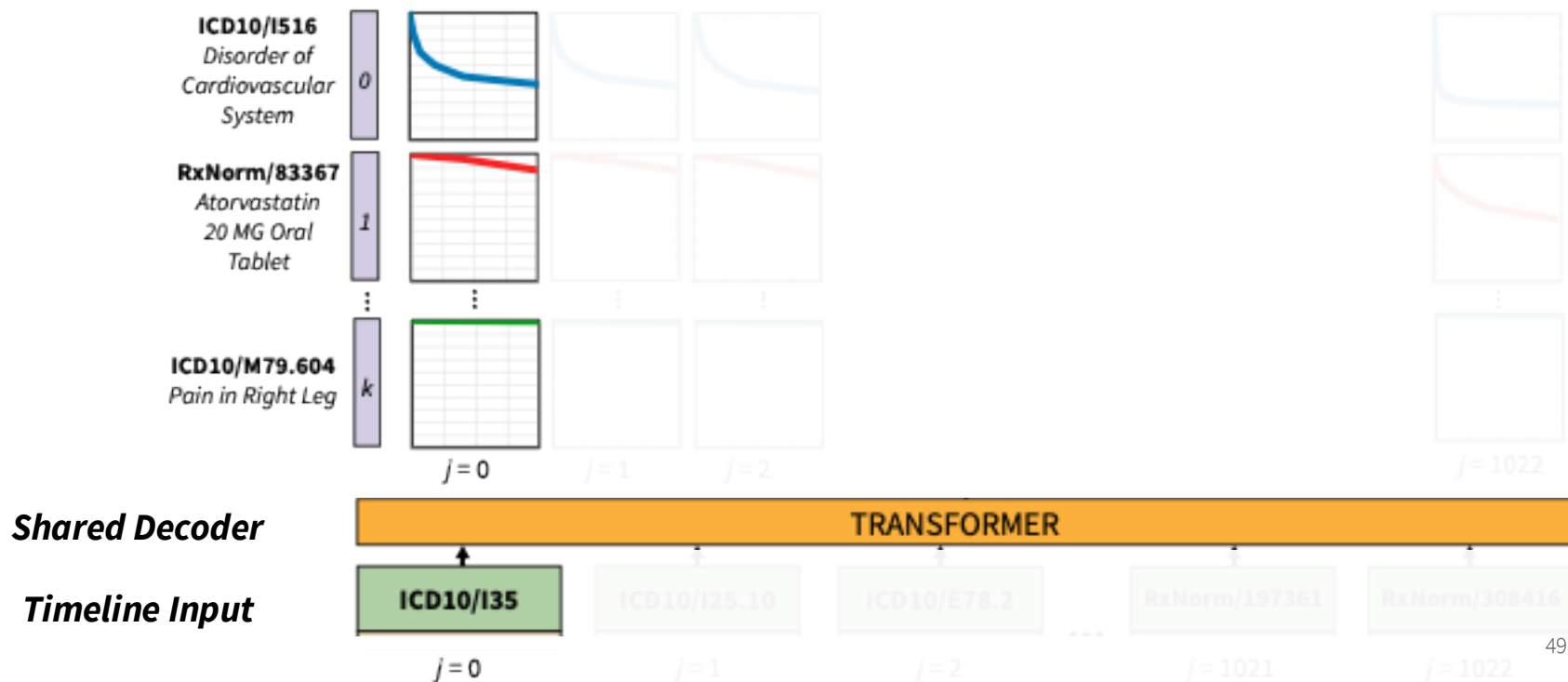
# Intuition Behind the Pretraining Objective

**Hypothesis:** Multi-task learning (MTL) will capture generalizable TTE features



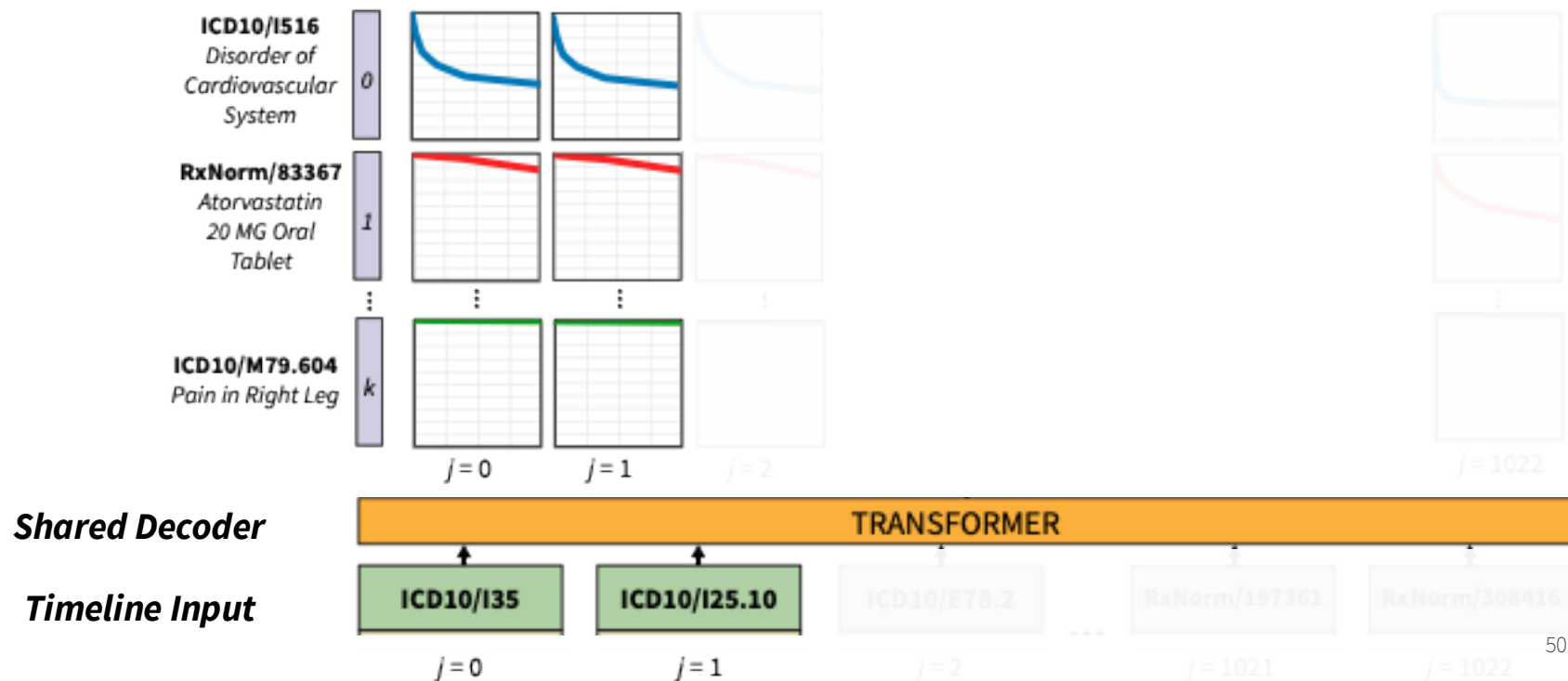
# Intuition Behind the Pretraining Objective

**Hypothesis:** Multi-task learning (MTL) will capture generalizable TTE features



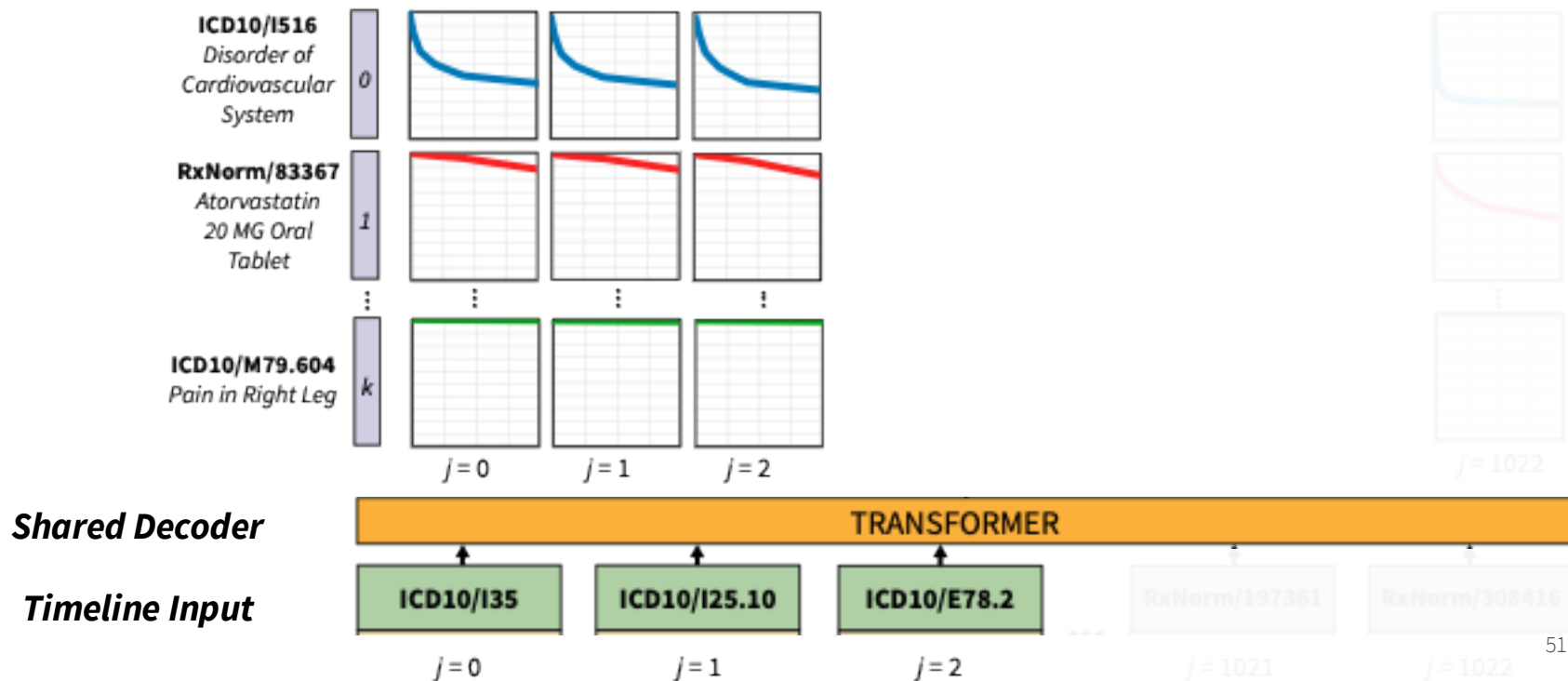
# Intuition Behind the Pretraining Objective

**Hypothesis:** Multi-task learning (MTL) will capture generalizable TTE features



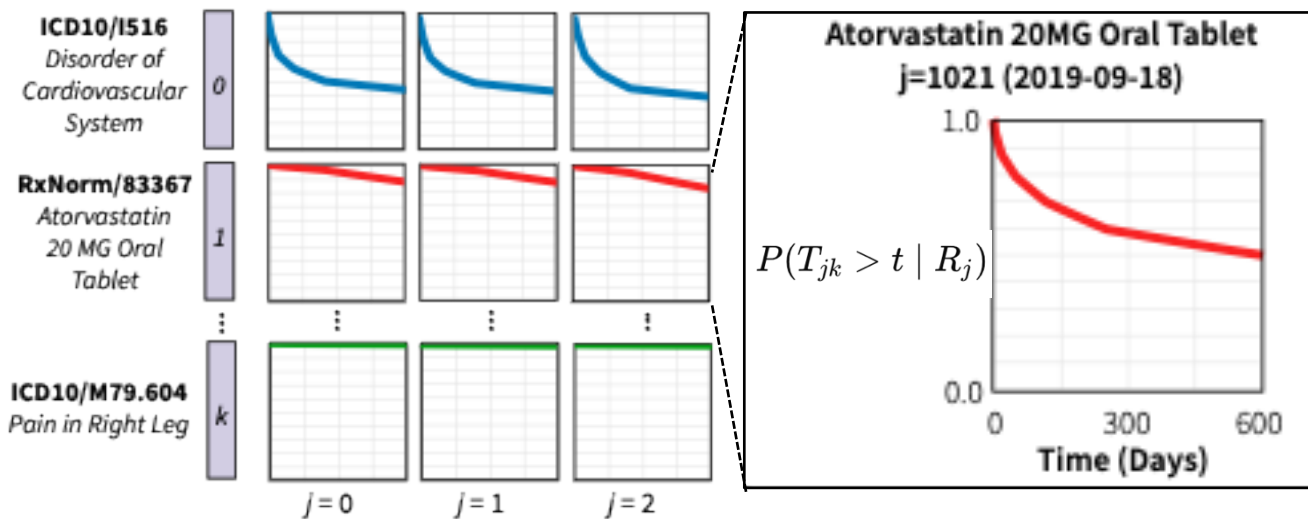
# Intuition Behind the Pretraining Objective

**Hypothesis:** Multi-task learning (MTL) will capture generalizable TTE features



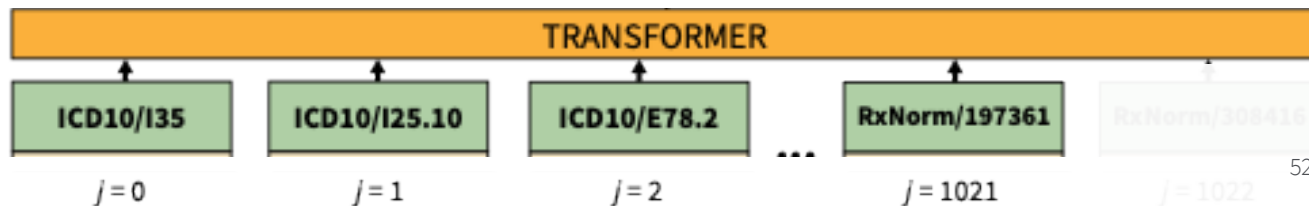
# Intuition Behind the Pretraining Objective

**Hypothesis:** Multi-task learning (MTL) will capture generalizable TTE features



**Shared Decoder**

**Timeline Input**

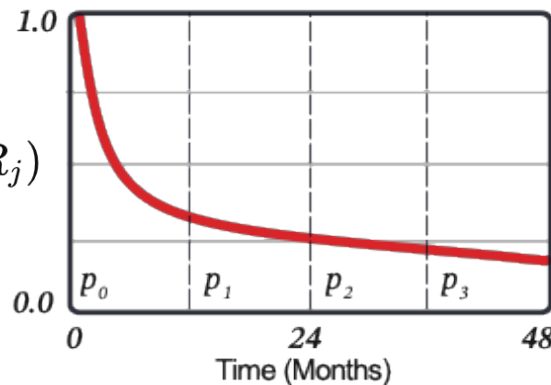


# Pretraining Objective

## Deep Piecewise Exponential Model

- Partition time into **pieces** for more expressive risk modeling
- For **piece p**, interval start and end time:  $[S_p, E_p)$
- **Hazard rate** is constant within this interval

$$P(T_{jk} > t \mid R_j)$$



For a patient with **event j, task k, and piece p**

**Piecewise Hazard Function**

$$h_{jk}(t) = \sum_{p=1}^P I(S_p \leq t < E_p) \lambda_{jkp}$$

t is within piece p

hazard rate for piece p

**Survival Function**

$$S_{jk}(t) = \prod_{p=1}^P \exp(-\lambda_{jkp} (\min(t, E_p) - S_p) I(t \geq S_p))$$

**Hazard Rate**

$$\lambda_{jkp} = \exp(W_p R_j \cdot \hat{\beta}_k)$$

*time-independent task embedding*

patient representation as of j  
piece-specific linear projection

TRANSFORMER  $f_\theta$

# Pretraining Objective

---

## Loss Function

Minimize the negative log-likelihood of the observed event times across all tasks and time pieces

$$\min_{\Theta} \mathcal{L}(\Theta) = - \sum_{j,k} \sum_{p=1}^P [\delta_{jkp} (\log \lambda_{jkp} - \lambda_{jkp} U_{jkp}) + (1 - \delta_{jkp}) (-\lambda_{jkp} U_{jkp})]$$

*all events and tasks*

*event happens in piece p*

*no event in piece p*

$U$  represents the amount of time an event is at risk within a given time interval

# Datasets & Tasks

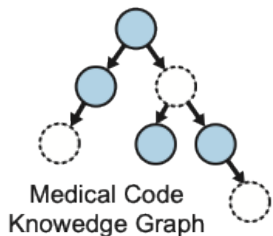
## Datasets

STANFORD STARR-OMOP (EHR)

2.7M Patients

3.5B Events

## Pretraining Tasks



Entropy-Ranked Vertex Cover for Task Selection

**Intuition:** We pick  $k$  tasks that **maximize diversity** by selecting nodes whose values are **least predictable** given their parents

$k \leq$   
16,392

## Evaluation Tasks

Celiac Disease

Stroke

Pancreatic Cancer

NAFLD

Heart Attack

Lupus

## ICD-10

Rule-based labeling

**We remove these tasks from the pretraining set**

## NLP-based

Measures generalization to labels not derived from codes



13 Chest X-ray Findings

# Results: MOTOR vs. Baselines

---

**MOTOR-Scratch** (no pretraining) largely **underperforms** compared to baselines

Method	Dataset	Celiac	HA	Lupus	NAFLD	Cancer	Stroke
Cox PH	EHR-OMOP	0.689	0.761	0.770	0.726	0.793	0.779
DeepSurv	-	0.704	0.823	0.790	0.800	0.811	0.830
DSM	-	0.707	0.828	0.784	0.805	0.809	0.835
DeepHit	-	0.695	0.826	0.807	0.805	0.809	0.833
RSF	-	0.729	0.836	0.787	0.802	0.824	0.840
MOTOR-Scratch	-	0.696	0.795	0.803	0.821	0.777	0.831

# Results: MOTOR vs. Baselines

But with **pretraining...**

**MOTOR-Probe & MOTOR-Finetune** outperform **SOTA on all tasks**

Avg improvement: **+4.6%**

Method	Dataset	Celiac	HA	Lupus	NAFLD	Cancer	Stroke
Cox PH	EHR-OMOP	0.689	0.761	0.770	0.726	0.793	0.779
DeepSurv	-	0.704	0.823	0.790	0.800	0.811	0.830
DSM	-	0.707	0.828	0.784	0.805	0.809	0.835
DeepHit	-	0.695	0.826	0.807	0.805	0.809	0.833
RSF	-	0.729	0.836	0.787	0.802	0.824	0.840
MOTOR-Scratch	-	0.696	0.795	0.803	0.821	0.777	0.831
MOTOR-Probe	-	0.802	0.884	0.850	0.859	0.865	0.874
MOTOR-Finetune	-	<b>0.802</b>	<b>0.887</b>	<b>0.863</b>	<b>0.864</b>	<b>0.865</b>	<b>0.875</b>

# Results: Autoregressive vs. TTE Pretraining

## Overall Performance

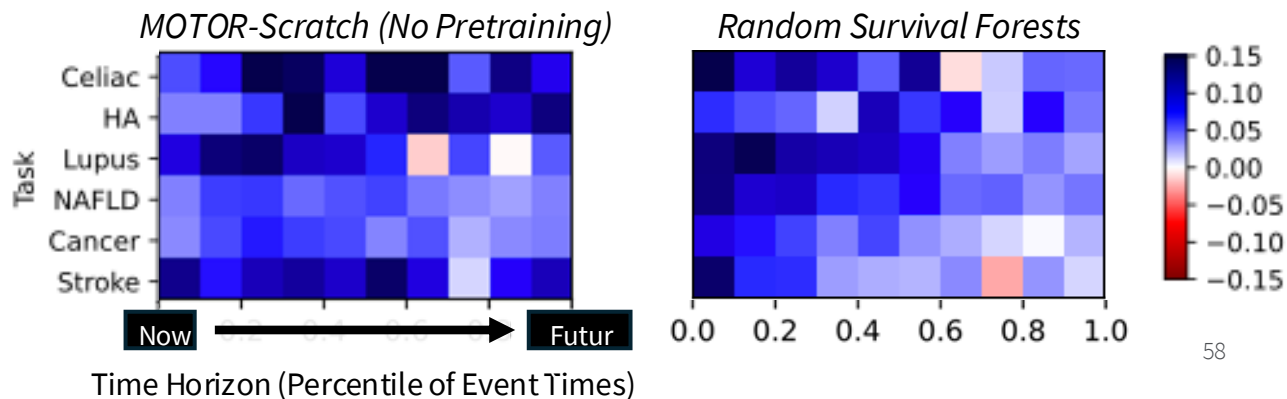
Objective	Celiac	HA	Lupus	NAFLD	Cancer	Stroke
RSF	0.729	0.836	0.787	0.802	0.824	0.840
Next Code	<u>0.774</u>	<u>0.862</u>	<u>0.842</u>	<u>0.860</u>	<u>0.860</u>	<u>0.857</u>
Time-to-Event	<b>0.802</b>	<b>0.887</b>	<b>0.863</b>	<b>0.864</b>	<b>0.865</b>	<b>0.875</b>

**Autoregressive beats SOTA (RSF)**  
...but **TTE beats autoregressive** by  
**~2%**

## Performance Comparison over Long Time Horizons

Performance Deltas of MOTOR with TTE Pretraining Versus:

**Pretraining is  
the key driver  
of performance**



# Results: Autoregressive vs. TTE Pretraining

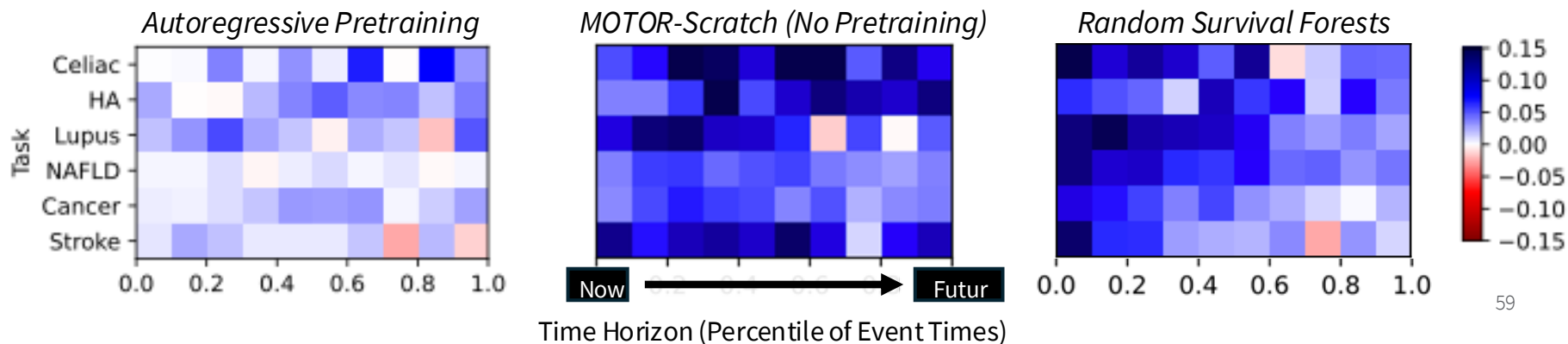
## Overall Performance

Objective	Celiac	HA	Lupus	NAFLD	Cancer	Stroke
RSF	0.729	0.836	0.787	0.802	0.824	0.840
Next Code	0.774	0.862	0.842	0.860	0.860	0.857
Time-to-Event	0.802	0.887	0.863	0.864	0.865	0.875

**Autoregressive beats SOTA (RSF)**  
...but **TTE beats autoregressive by ~2%**

## Performance Comparison over Long Time Horizons

Performance Deltas of MOTOR with TTE Pretraining Versus:



# Evaluation: EHR Foundation Models

CTAGCTCC<sub>G...</sub>



# Reproducibility in Healthcare AI

SCIENCE TRANSLATIONAL MEDICINE | PERSPECTIVE

BIOMEDICAL POLICY

## Reproducibility in machine learning for health research: Still a ways to go

Matthew B. A. McDermott<sup>1\*†</sup>, Shirly Wang<sup>2,3†</sup>, Nikki Marinsek<sup>4</sup>, Rajesh Ranganath<sup>5</sup>,  
Luca Foschini<sup>4</sup>, Marzyeh Ghassemi<sup>2,6,7</sup>

Longstanding  
Reproducibility  
Challenges

Medical data are noisy, **replete  
with errors, biases, missingness**

Most AI is **trained and  
tested** on **cleaned data**

REVIEW

## Global healthcare fairness: We should be sharing more, not less, data

Kenneth P. Seastedt<sup>1☉\*</sup>, Patrick Schwab<sup>2☉</sup>, Zach O'Brien<sup>3☉</sup>, Edith Wakida<sup>4☉</sup>,  
Karen Herrera<sup>5☉</sup>, Portia Grace F. Marcelo<sup>6☉</sup>, Louis Agha-Mir-Salim<sup>7,8☉</sup>, Xavier  
Borrat Frigola<sup>8,9☉</sup>, Emily Boardman Ndulue<sup>10☉</sup>, Alvin Marcelo<sup>11☉</sup>, Leo  
Anthony Celi<sup>8,12,13☉</sup>

PLOS DIGITAL HEALTH

# Multiple Choice vs. Longitudinal Patient Timelines

## MedQA

Question: A 35-year-old man is brought to the emergency department by a friend 30 minutes after the sudden onset of right-sided weakness and difficulty speaking. [...] Which of the following is the most appropriate next step in diagnosis?

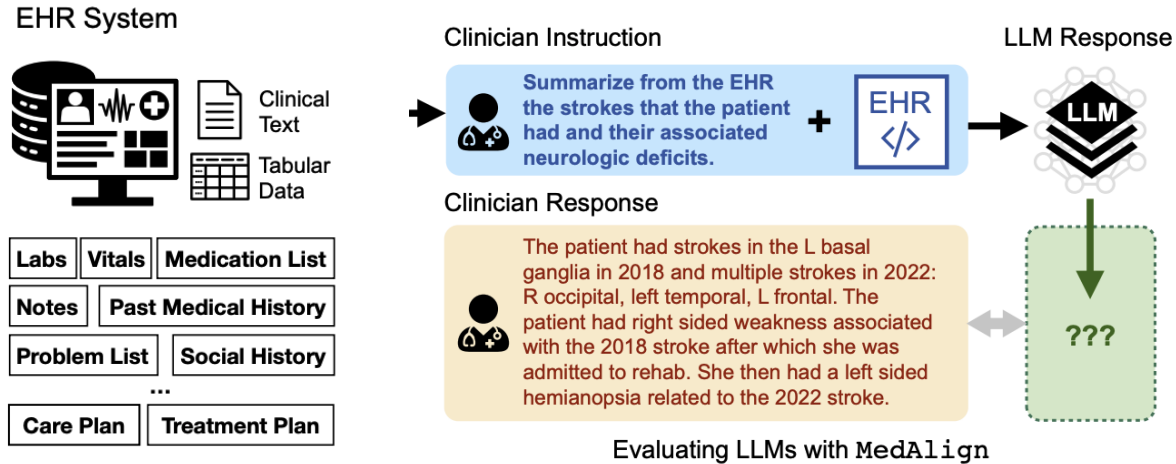
- (A) Echocardiography with bubble study
- (B) Adenosine stress test
- (C) Cardiac catheterization
- (D) Cardiac MRI with gadolinium
- (E) CT angiography



```
<record>
  <visit type="Emergency Room Visit" start="10/08/2018 20:00">
    <day start="10/08/2018 20:00">
      <person>
        Birth:7/19/1966
        Race
        Gender
        Ethnicity
        Age
        Age
      </person>
      <condition>
        <code>[LOINC/LP21258-6] Oxygen saturation 96 %</code>
      </condition>
      <visit>
        <code>
      </visit>
      <measurement>
        <code>
      </measurement>
      <procedure>
        <code>
      </procedure>
      <observation start="10/08/2018 08:10 PM">
        <code>[LOINC/LP21258-6] Oxygen saturation 96 %</code>
      </observation>
      <note type="emergency department note" start="10/08/2018 08:10 PM">
        Emergency Department Provider Note Name: Jessica Jones, MD MRN: [1234555]
        ED Arrival: 10/08/2018 Room #: 17B History and Physical Triage: 52 year old woman
        with unknown past medical history presenting with right sided weakness since about
        2 hours ago. Last known normal 5:45pm. She said she was feeling well and then suddenly
        noticed that her right arm and leg went limp. She denies taking any blood thinners,
        and has had no recent surgeries. NIHSS currently graded at an 8: 4 no movement in R
        arm and 4 no movement in R leg CT head is negative for any bleed or any early ischemic
        changes. INR is 1.0, Plt 133. Discussed with patient the severity of symptoms and the
        concern that they are caused by a stroke, and that IV tPA is the best medication to
        reduce the risk of long term deficits. Patient is agreeable and IV tPA was given at
        8:20pm. Initially SBP 210/100, labetalol 5mg IV x1 given and came down to 180/90.
        IV tPA given after this point. Patient will need to be admitted to the ICU, with close
        neurological monitoring. Plan for head CT 24 hours post IV tPA administration, stroke
        workup including LDL, HA1C, echo, tele monitoring. Local neurology consult in AM.
      </note>
      <measurement start="10/08/2018 08:15 PM">
        <code>[LOINC/70182-1] NIHSS 8 </code>
      </measurement>
    </record>
```

## Longitudinal Patient Timelines

# Instruction Tuning: Aligning with Clinical Needs



**MedAlign:** A Clinician-Generated Benchmark Dataset for Instruction Following with Electronic Medical Records [1]

- **15** clinicians / **7** specialties
- 983 instructions, 303 responses
- Assess **real information needs**

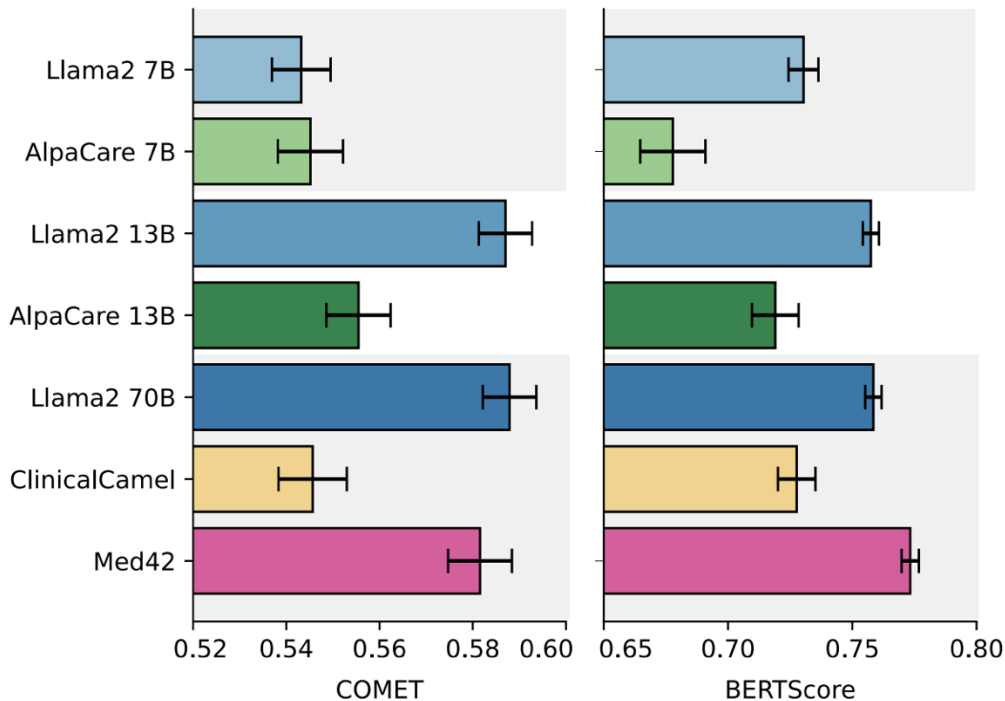
# Instruction Tuning: Aligning with Clinical Needs

Model	Context	Correct $\uparrow$	WR $\uparrow$	Rank $\downarrow$
GPT-4 (MR)	32768 <sup>†</sup>	<b>65.0%</b>	0.658	2.80
GPT-4	32768	60.1%	<b>0.676</b>	<b>2.75</b>
GPT-4	2048*	51.8%	0.598	3.11
Vicuña-13B	2048	35.0%	0.401	3.92
Vicuña-7B	2048	33.3%	0.398	3.93
MPT-7B-Instruct	2048	31.7%	0.269	4.49

GPT-4 **35% Error Rate**

# Instruction Tuning in Medical LLMs



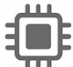




Base vs. Base + Medical Instruction Tuning



Current short instruction tuning tasks for medicine (e.g., MedQA) **actually hurt performance on MedAlign**

**A Single Benchmark Does NOT Tell the Whole Story!**

# Longitudinal, Multimodal EHR Dataset Releases

 Dataset	 Task	 Technical Challenge	 Example	 Tabular	 Images	 Notes
<b>EHRSHOT</b>	Risk Stratification	Few-Shot Learning	<i>What is the likelihood that this patient gets a diagnosis of pancreatic cancer within the next year?</i>	<input checked="" type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>
<b>INSPECT</b>	Time-to-Event Modeling	Multimodal Learning	<i>When is chronic pulmonary hypertension most likely to develop</i>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>
<b>MedAlign</b>	Instruction Following	Long-Context Learning & Temporal Reasoning	<i>From this EHR, summarize the patient's history of strokes and the resulting neurologic deficits.</i>	<input checked="" type="checkbox"/>	<input type="checkbox"/>	<input checked="" type="checkbox"/>

**26k** Patients

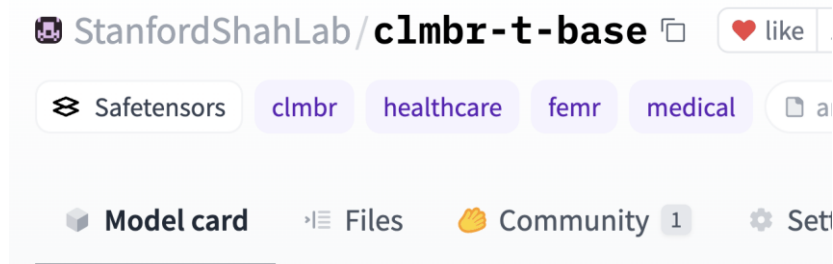
**295M**

**442k** Visits



<https://redivis.com/ShahLab>

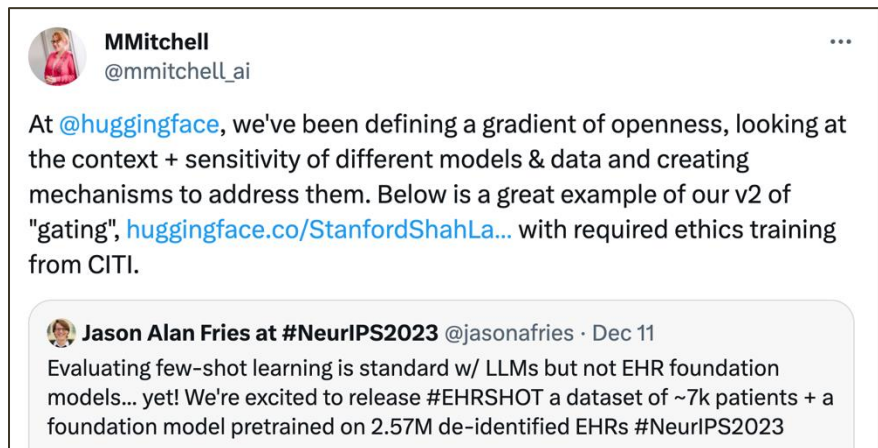
# Enabling Open Science



 **Gated model** You have been granted access to this model

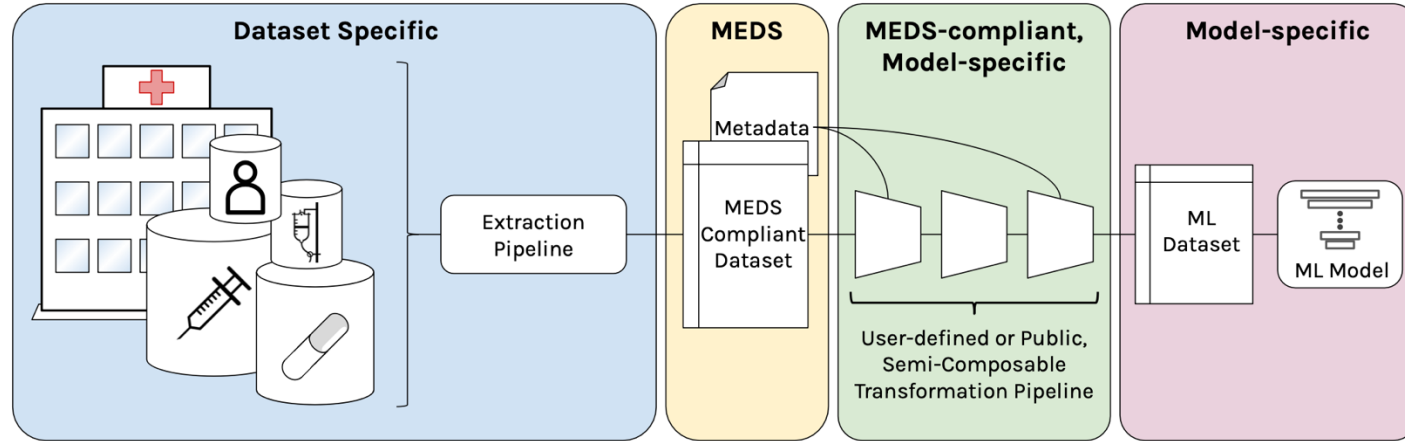
## First EHR model hub release!

- Gated model on Hugging Face
- Requires **CITI ethics training**
- **Non-commercial use only**



Margaret Mitchell  
Chief AI Ethics Scientist, Hugging Face

# Medical Event Data Standard (MEDS)



## Open Data Schema for Health AI Practitioners

*Bert Arnrich, Edward Choi, Jason A. Fries, Matthew B. A. McDermott, Jungwoo Oh, Tom J Pollard, Nigam Shah, Ethan Steinberg, Michael Wornow, Robin van de Water*

<https://github.com/Medical-Event-Data-Standard/meds>



# Future: Research Opportunities

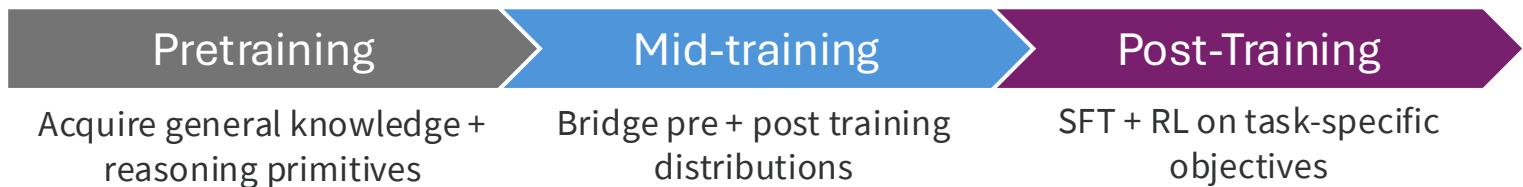
CTAGCTCC<sub>G...</sub>



© 2024 [Company Name]. All rights reserved. This document contains confidential information and is intended for internal use only. Any unauthorized distribution or use is strictly prohibited. For more information, please contact the legal department.

# Data-Centric AI and Modern Data Pipeline

Training data defines an LLM's capabilities, limitations, and biases



Published as a conference paper at COLM 2025

## LIMO: Less is More for Reasoning

Yixin Ye<sup>1,2\*</sup>, Zhen Huang<sup>2,3\*</sup>, Yang Xiao<sup>4</sup>, Ethan Chern<sup>1,2</sup>, Shijie Xia<sup>1,2</sup>, Pengfei Liu<sup>1,2\*</sup>

<sup>1</sup>Shanghai Jiao Tong University <sup>2</sup>SII-GAIR

<sup>3</sup>Fudan University <sup>4</sup>The Hong Kong Polytechnic University

## Math Reasoning

**800** carefully selected  
training examples  
**1% of training data**

AIME24

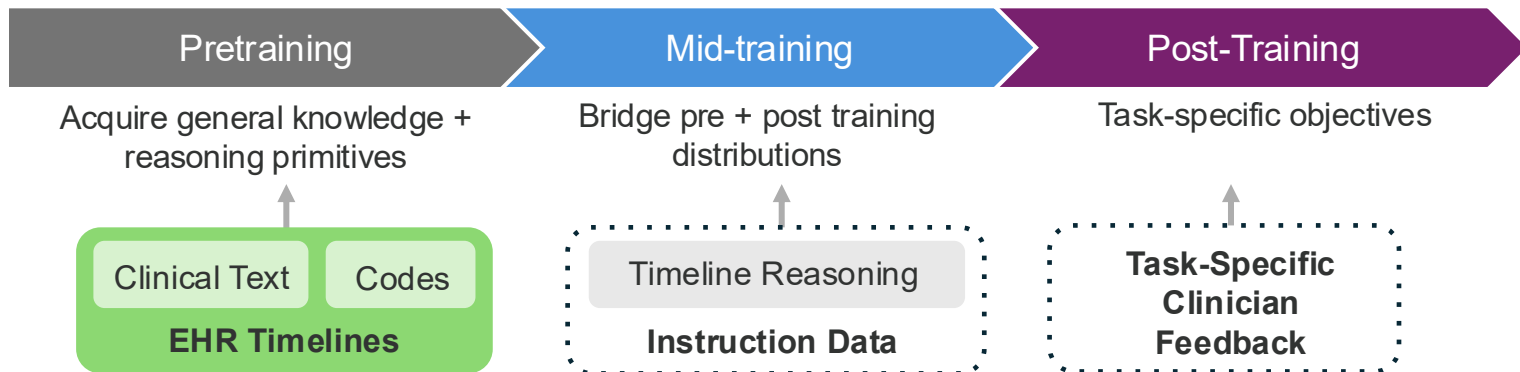
**6.5%** → **63.3%**

MATH500

**59.2%** → **95.6%**

# Building the Next-Generation EHR Foundation Models

**Training data** defines an LLM's **skills, limitations, and biases**



## Enhancing Our Pretraining Recipe

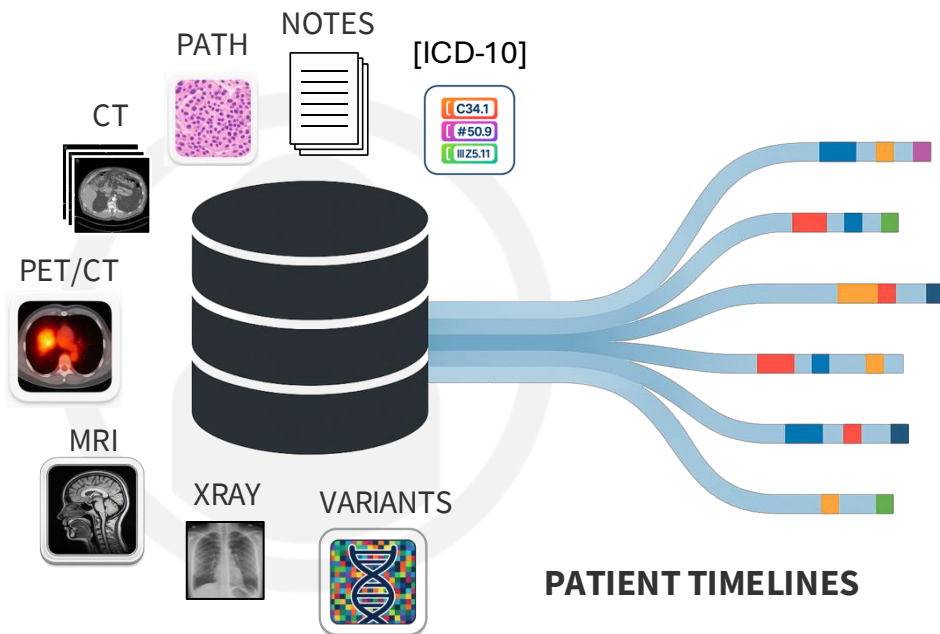
Unify **autoregressive & survival** pretraining

+

Optimize medical code tokenization to **improve cross-site robustness**

# The Future: Multimodal, Interactive Healthcare AI

**Multimodal timelines** as a unifying framework for retrieval, interaction, and feedback.



**Chat interfaces** for EHR



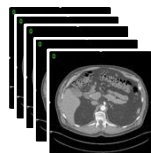
**Use similar patients** to inform decision making



**Feedback loops** to capture error signals and **improve data and models**

# Multimodal Time-to-Event Pretraining

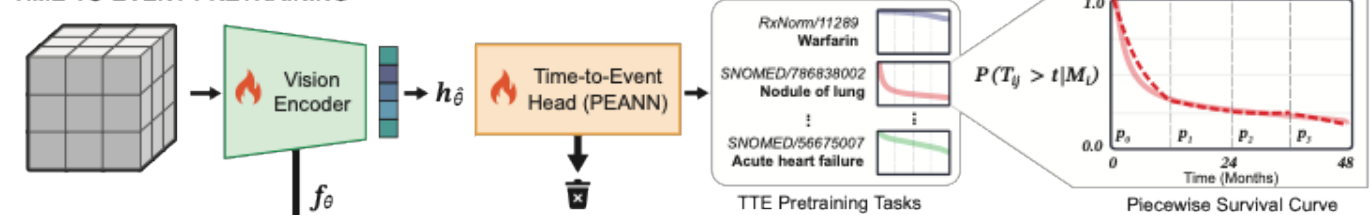
## Pulmonary Embolisms



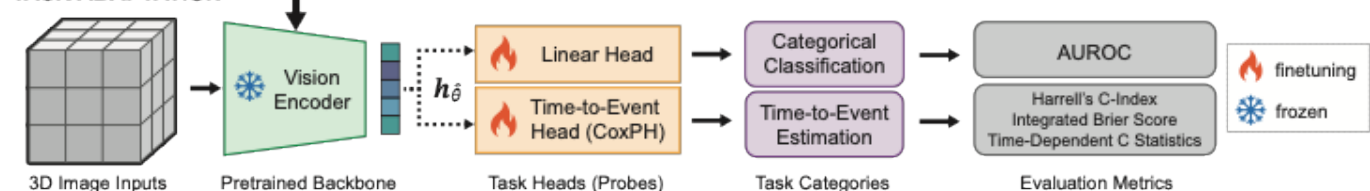
18,945 CT Scans  
(4.2 Million 2D images)

- Same pretraining setup as MOTOR
- **Single time point** (not dynamic)
- Pretraining a 3D image encoder

### TIME-TO-EVENT PRETRAINING

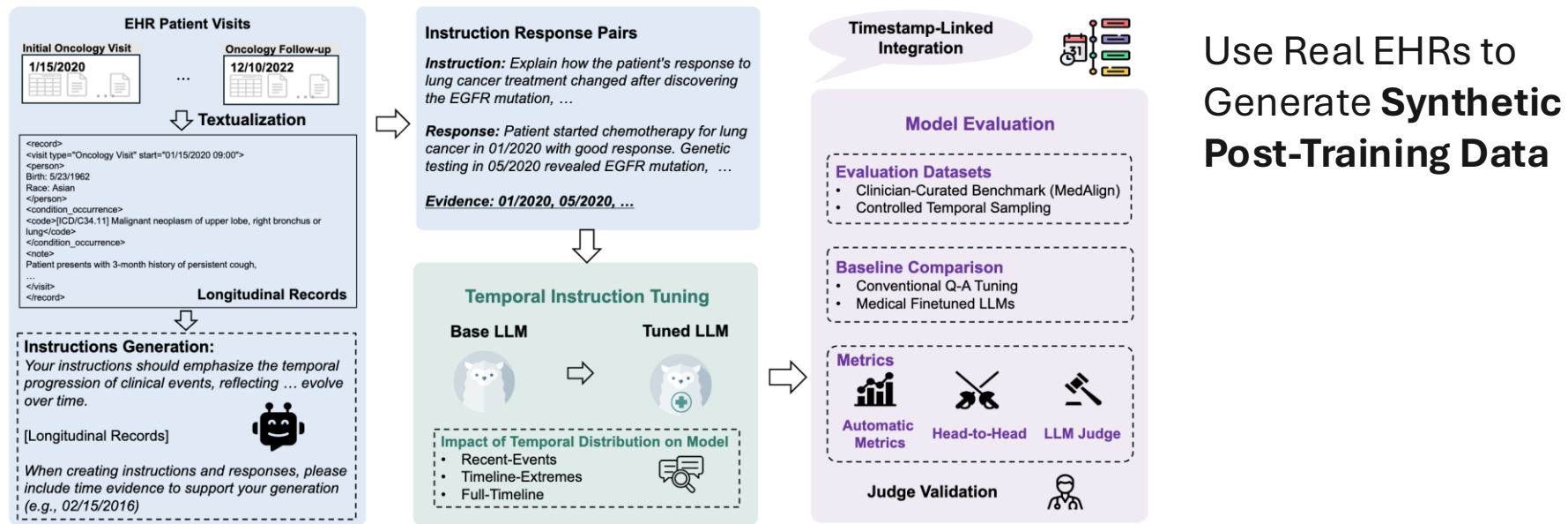


### TASK ADAPTATION



Time-to-Event Pretraining for 3D Medical Imaging  
Huo et al. ICLR 2025.

# Synthetic Data Generation



TIMER: Temporal Instruction Modeling and Evaluation for Longitudinal Clinical Records  
Cui et al. 2025. Preprint

# Data-Centric AI: Training Mixtures

---

- **Exclusion biases in training data**
- General **data scarcity** (e.g., rare diseases)
- **Limited EHR datasets and benchmarks** for pediatric populations
- Unique data processing challenges
  - Example: Child and mother combined in a single patient record
- **Limited patient history** vs. adults
- Rapid developmental changes



Thank You!

[jfries@stanford.edu](mailto:jfries@stanford.edu)

# Appendix

CTAGCTCC<sub>G...</sub>



ВНИМАНИЕ! ВНИМАНИЕ! ВНИМАНИЕ!  
ВНИМАНИЕ! ВНИМАНИЕ! ВНИМАНИЕ!  
ВНИМАНИЕ! ВНИМАНИЕ! ВНИМАНИЕ!  
ВНИМАНИЕ! ВНИМАНИЕ! ВНИМАНИЕ!  
ВНИМАНИЕ! ВНИМАНИЕ! ВНИМАНИЕ!



# **BERT-Style (Masked Language Modeling)**

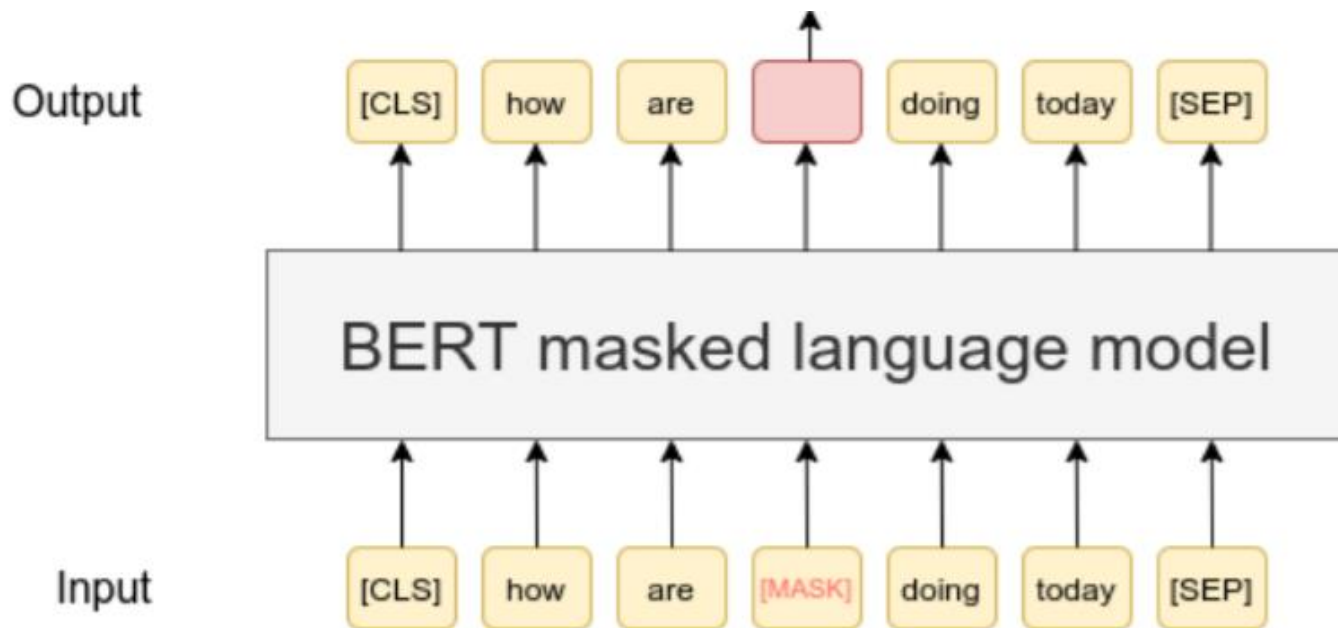
BEHRT (Li et al. 2020)

MedBERT (Rasmy et al. 2021)

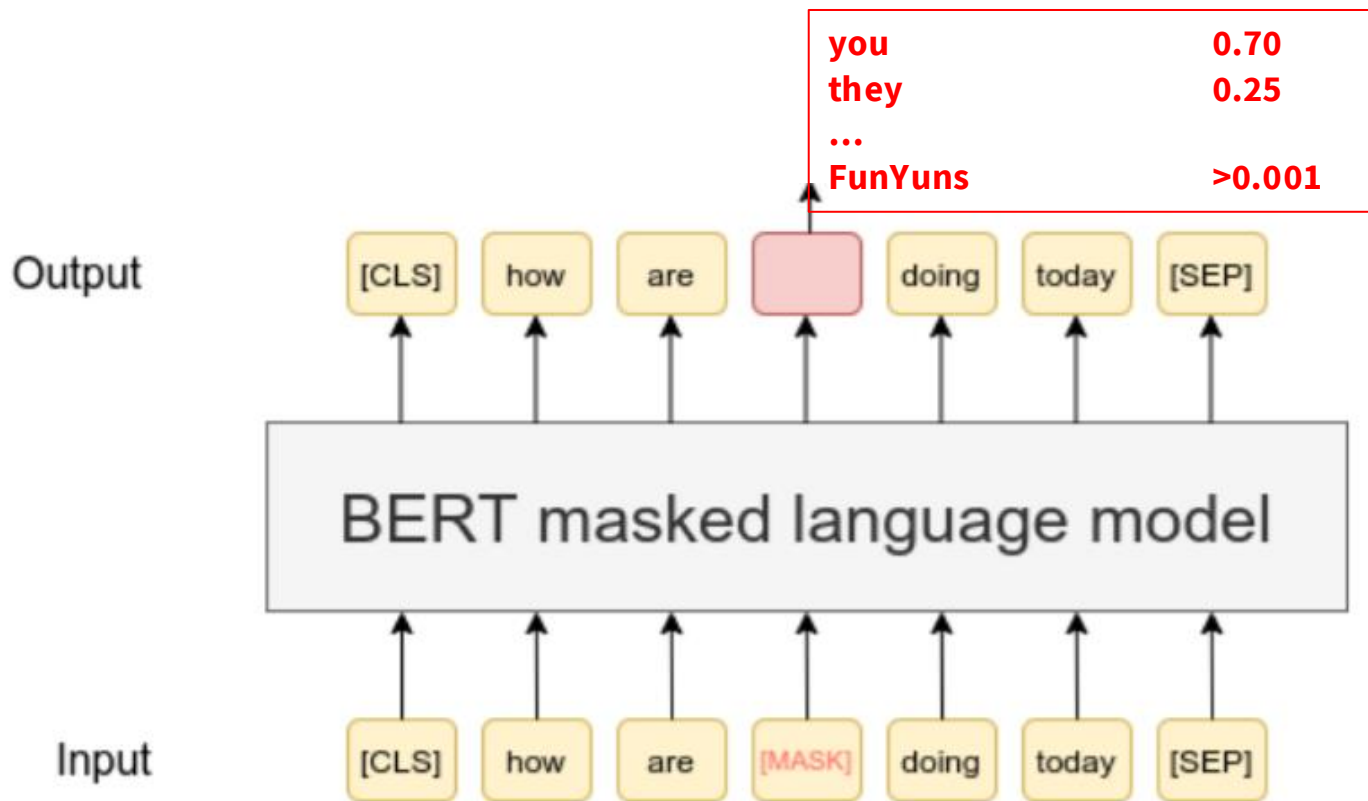
ClaimPT (Zeng et al. 2022)

# Corruption-based (Masking) Pretraining Objective

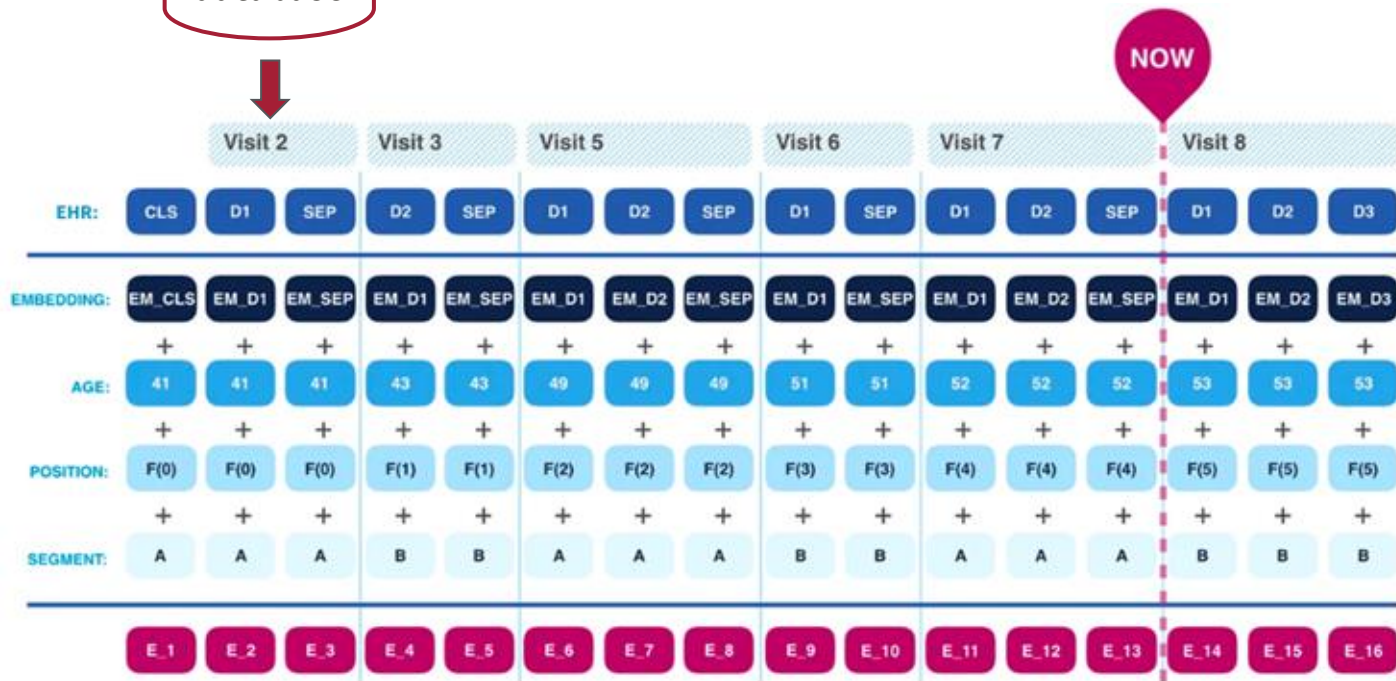
- **Mask tokens (15%)**
- **Train Model to Predict [MASK]'ed tokens**



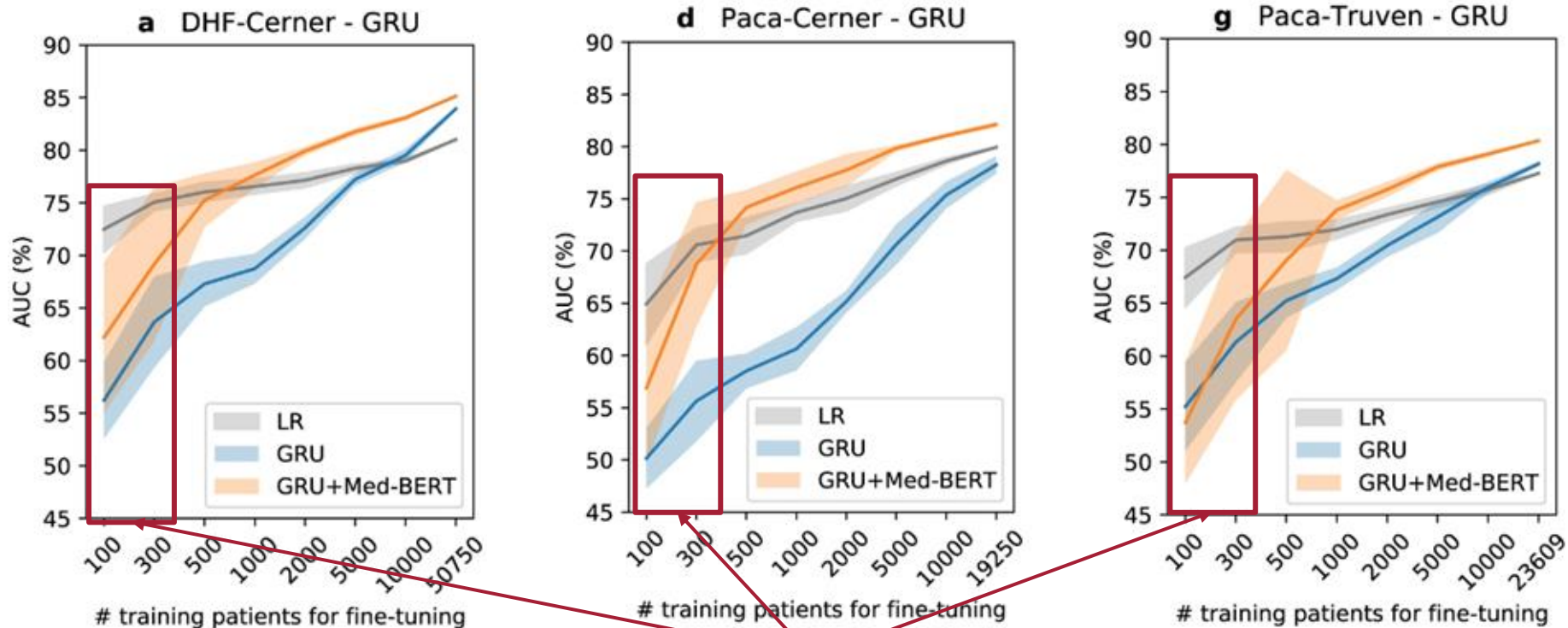
# Corruption-based (Masking) Pretraining Objective



# BERT-based Architecture (BEHRT)

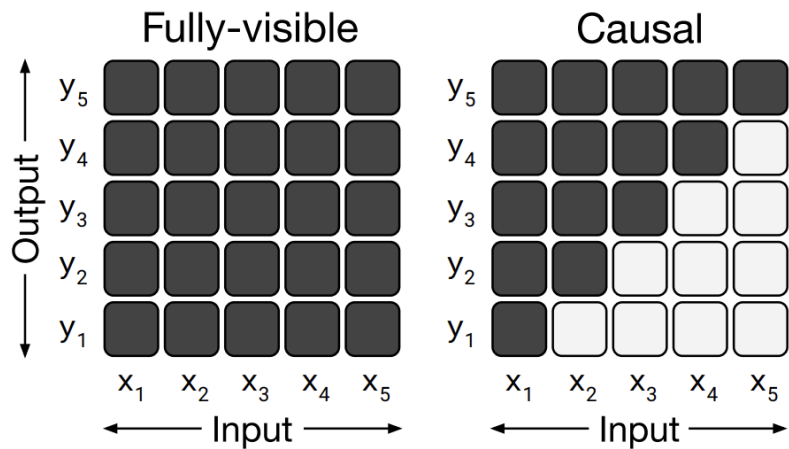


# Better performance than baselines (MedBERT)



**But few-shot performance isn't great...**

# Other Disadvantages



Raffel et al. 2019

Masked Language Modeling uses **bidirectional attention**. Good for summarizing a sequence, but **not generating the next event/token**

# Instruction Tuning: Aligning with Clinical Needs

Table 2: MEDALIGN instruction categories and example instructions.

Category	Example Instruction	Gold	All
Retrieve & Summarize	Summarize the most recent annual physical with the PCP	223	667
Care Planning	Summarize the asthma care plan for this patient including relevant diagnostic testing, exacerbation history, and treatments	22	136
Calculation & Scoring	Identify the risk of stroke in the next 7 days for this TIA patient	13	70
Diagnosis Support	Based on the information I've included under HPI, what is a reasonable differential diagnosis?	4	33
Translation	I have a patient that speaks only French. Please translate these FDG-PET exam preparation instructions for her	0	2
Other	What patients on my service should be prioritized for discharge today?	41	75
Total		303	983

Clinicians spend 49% of their day interacting with EHRs! **>66% of instructions** were **"retrieve & summarize"** data from the EHR.