

Transparente KI Methoden in der Medizin

Carsten Eickhoff

EBERHARD KARLS
UNIVERSITÄT
TÜBINGEN



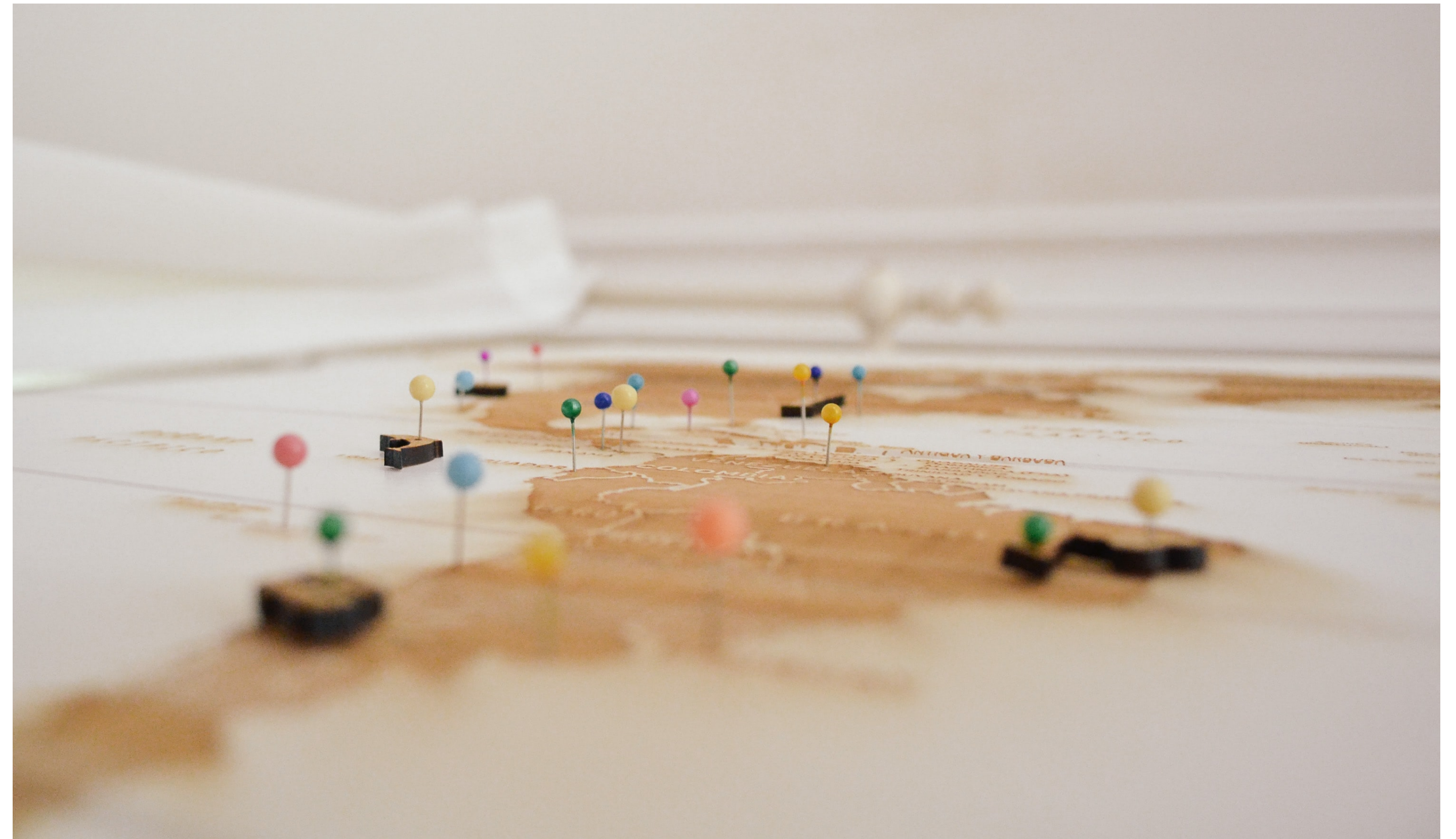
BROWN

Overview

- 1. Background, Research Interests, Lab**
- 2. Some Teasers**
- 3. Zero-shot [Text Classification | Diagnostic Decision Support]**
- 4. Discussion**

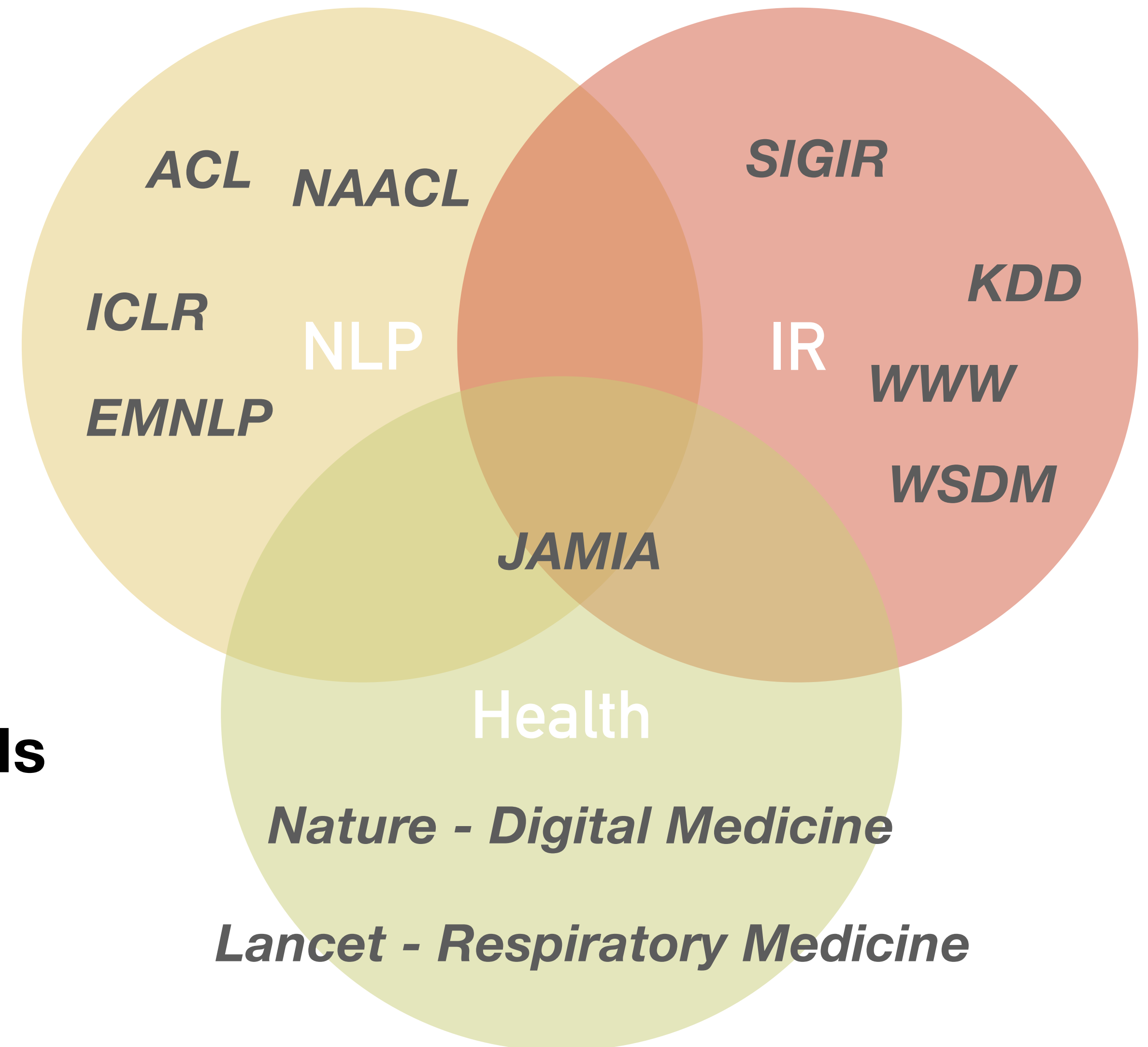
About me

- **Hannover**
- **University of Edinburgh**
- **TU Delft**
- **Microsoft**
- **ETH Zurich**
- **Harvard University**
- **Brown University**
- **University of Tübingen**



Interests

- **Dense Retrieval**
- **XIR**
- **Uncertainty Aware Models**
- **Grounded Language Modeling**
- **Manifold Learning for Neural LMs**
- **Clinical Decision Support**



Team



ADEEL ABBASI
ASSISTANT PROFESSOR



FLORIAN ROTTACH
PHD STUDENT



MARÍA VENEGAS-CARRO
SCIENTIFIC GRANT WRITER



ALI BAHRAINIAN
POSTDOC



GEORGE ZERVEAS
PHD STUDENT



MICHAL GOLOVANEVSKY
PHD STUDENT



AMINA ABDULLAHI
PHD STUDENT



JACK MERULLO
PHD STUDENT



RUOCHEN ZHANG
PHD STUDENT



CARSTEN EICKHOFF
PROFESSOR, DIRECTOR



LARS NETH
ADMINISTRATIVE ASSISTANT



SHAKILA MOSTAAN
SCIENTIFIC GRANT WRITER



CATHERINE CHEN
PHD STUDENT

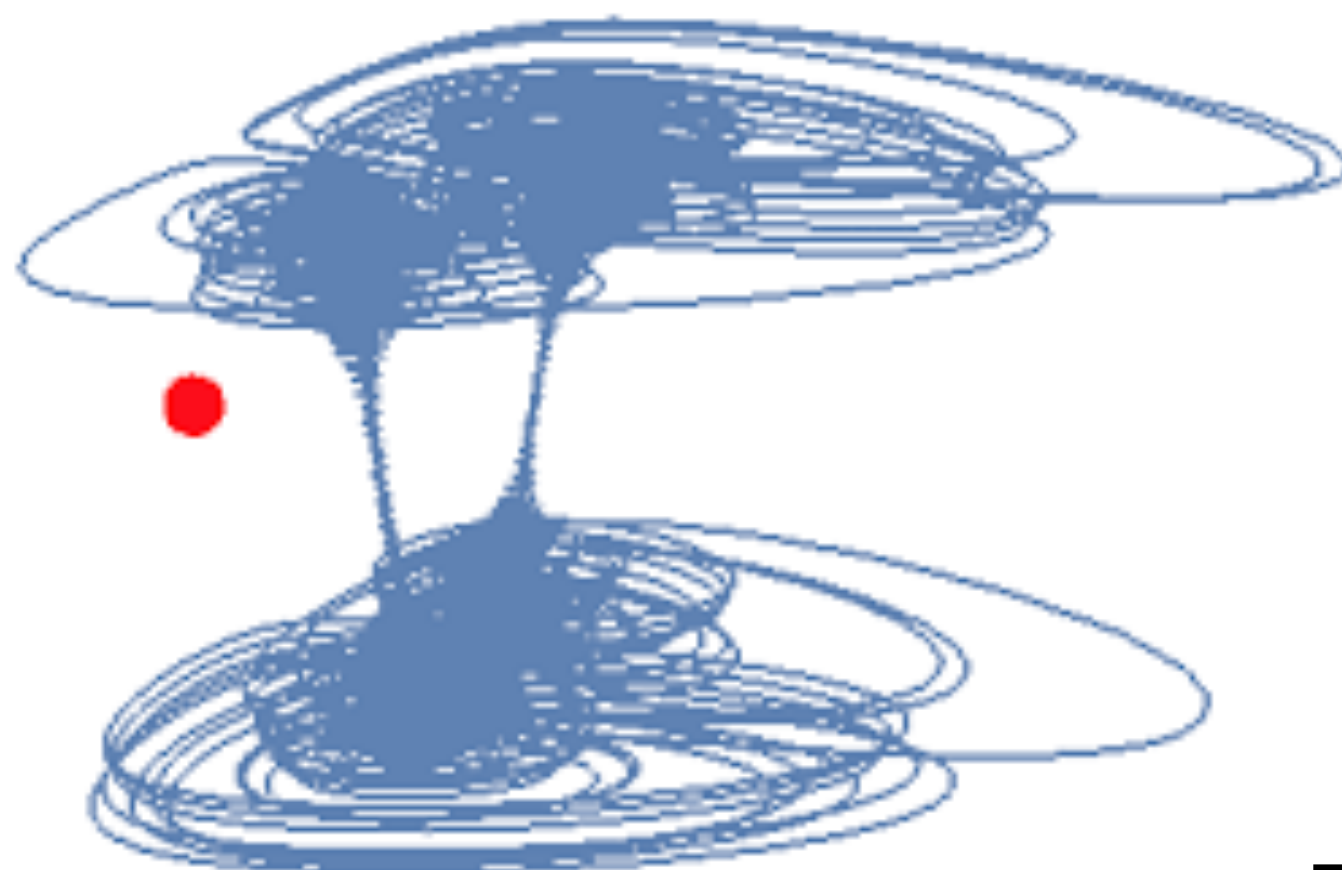
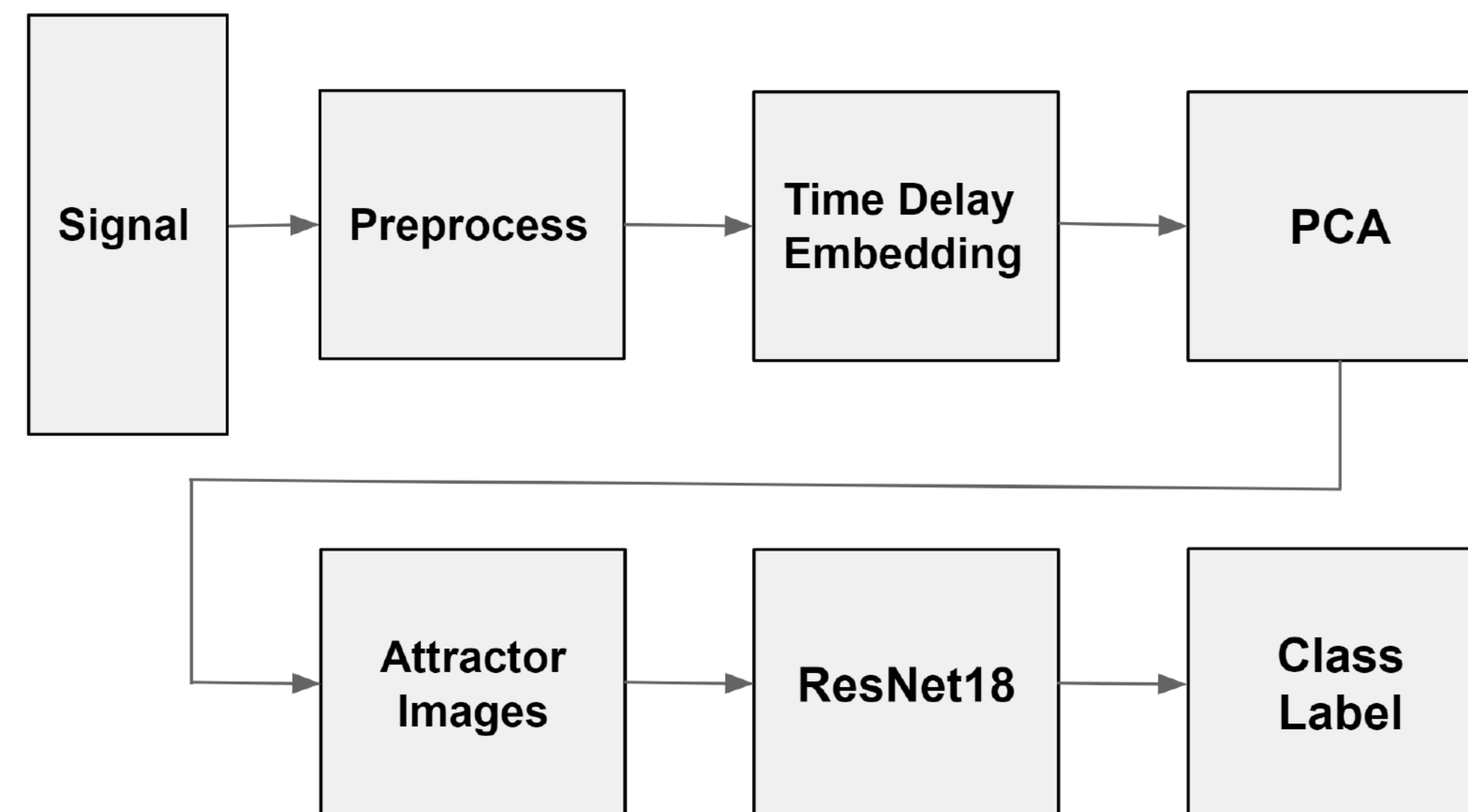
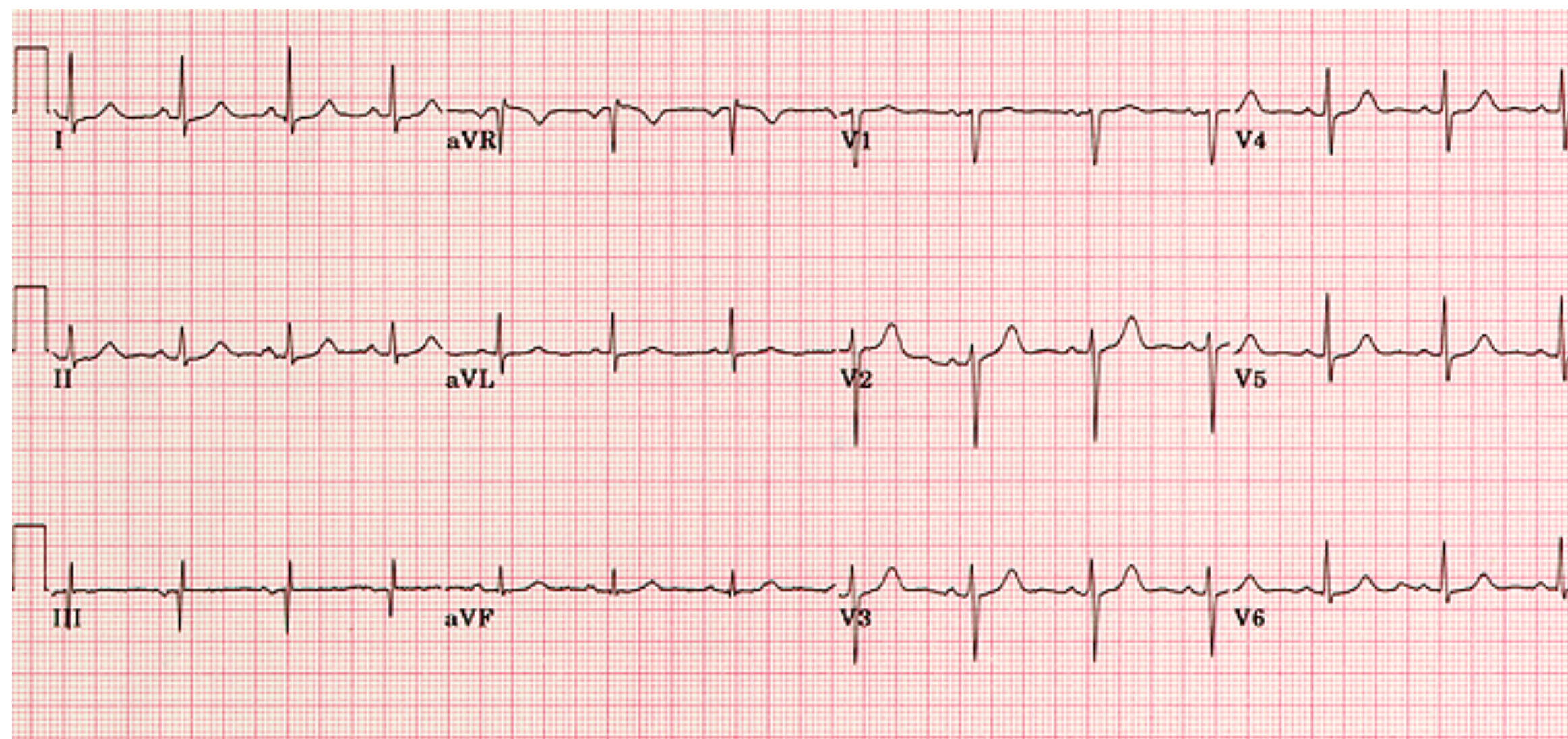


LAURA MERCURIO
ASSISTANT PROFESSOR



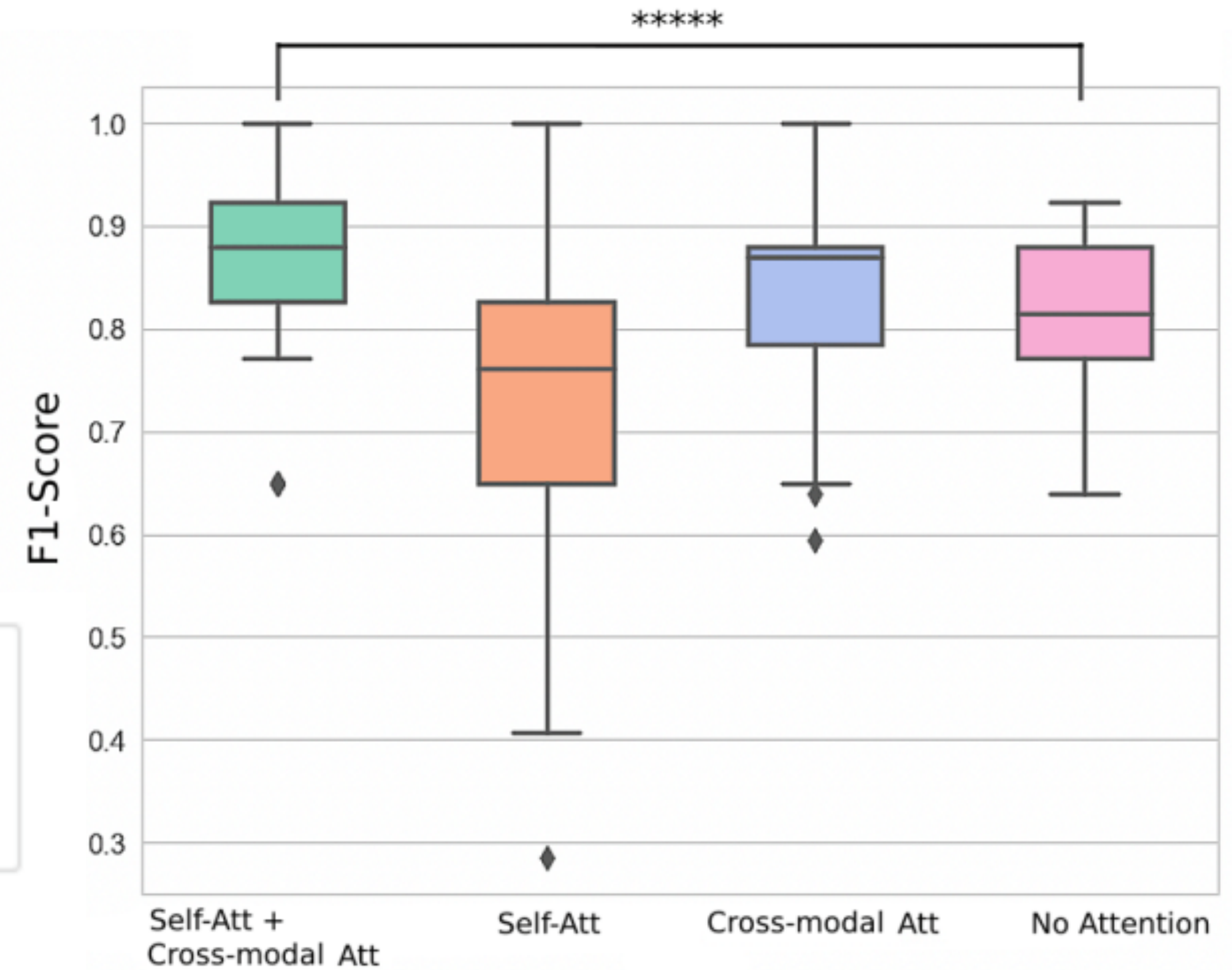
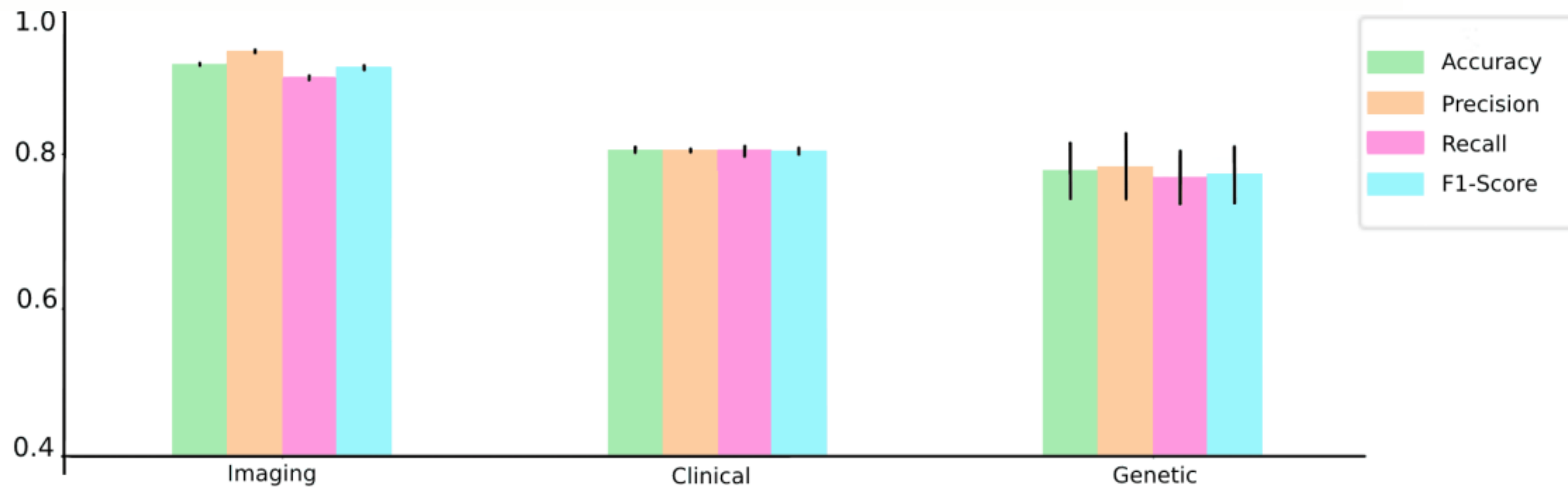
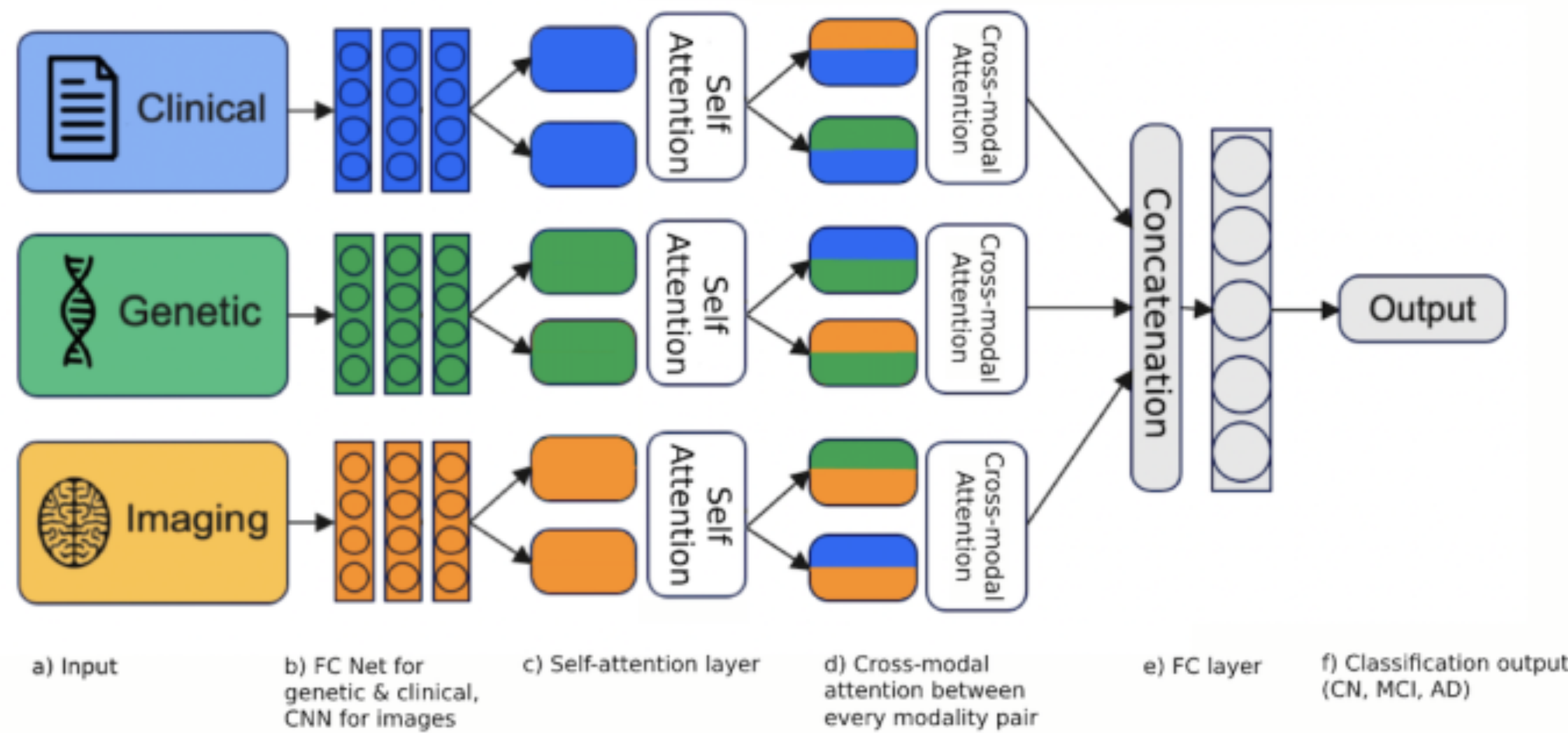
WILLIAM RUDMAN
PHD STUDENT

Cardiac Arrhythmia Localization via ECG Attractor Imaging

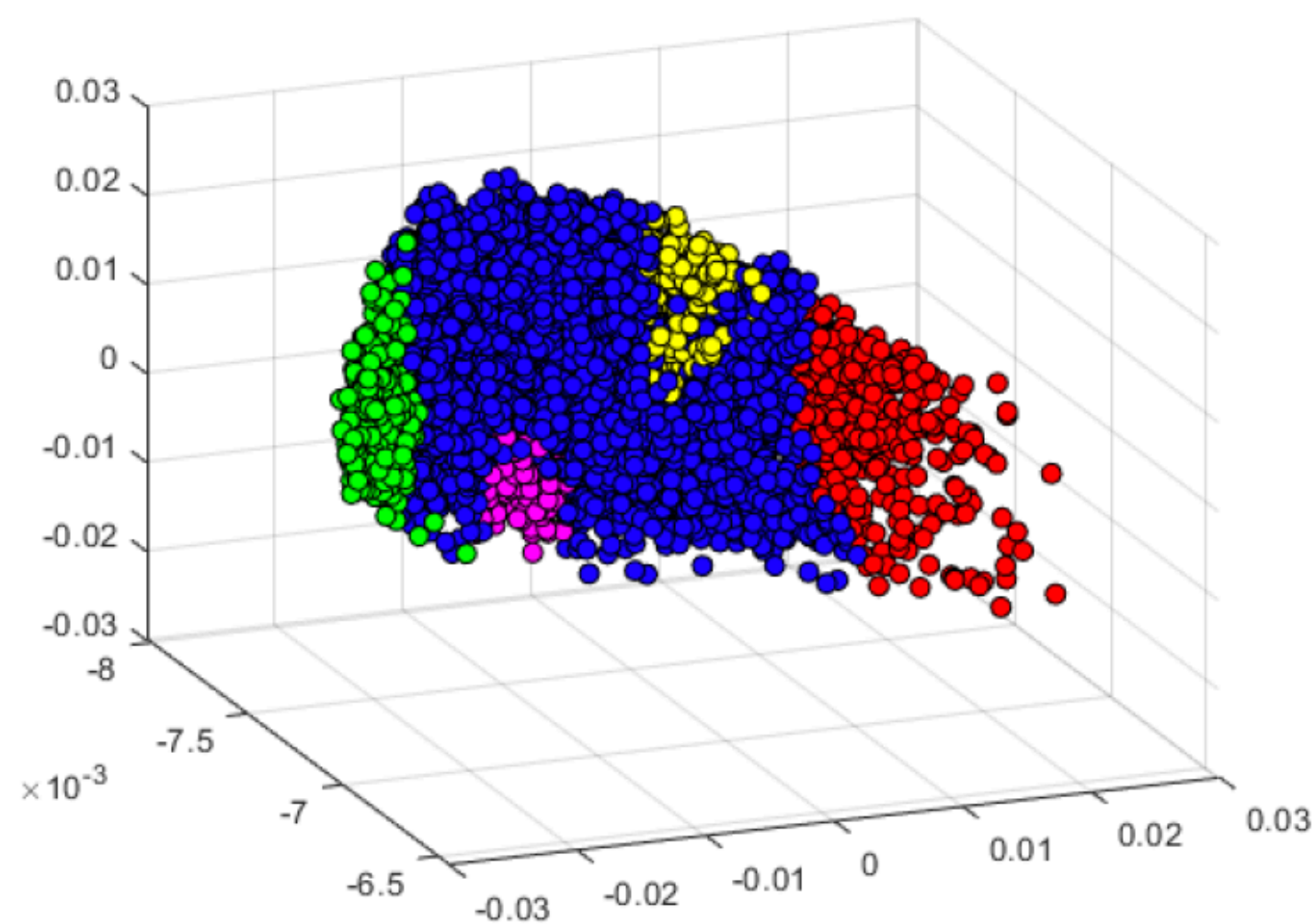
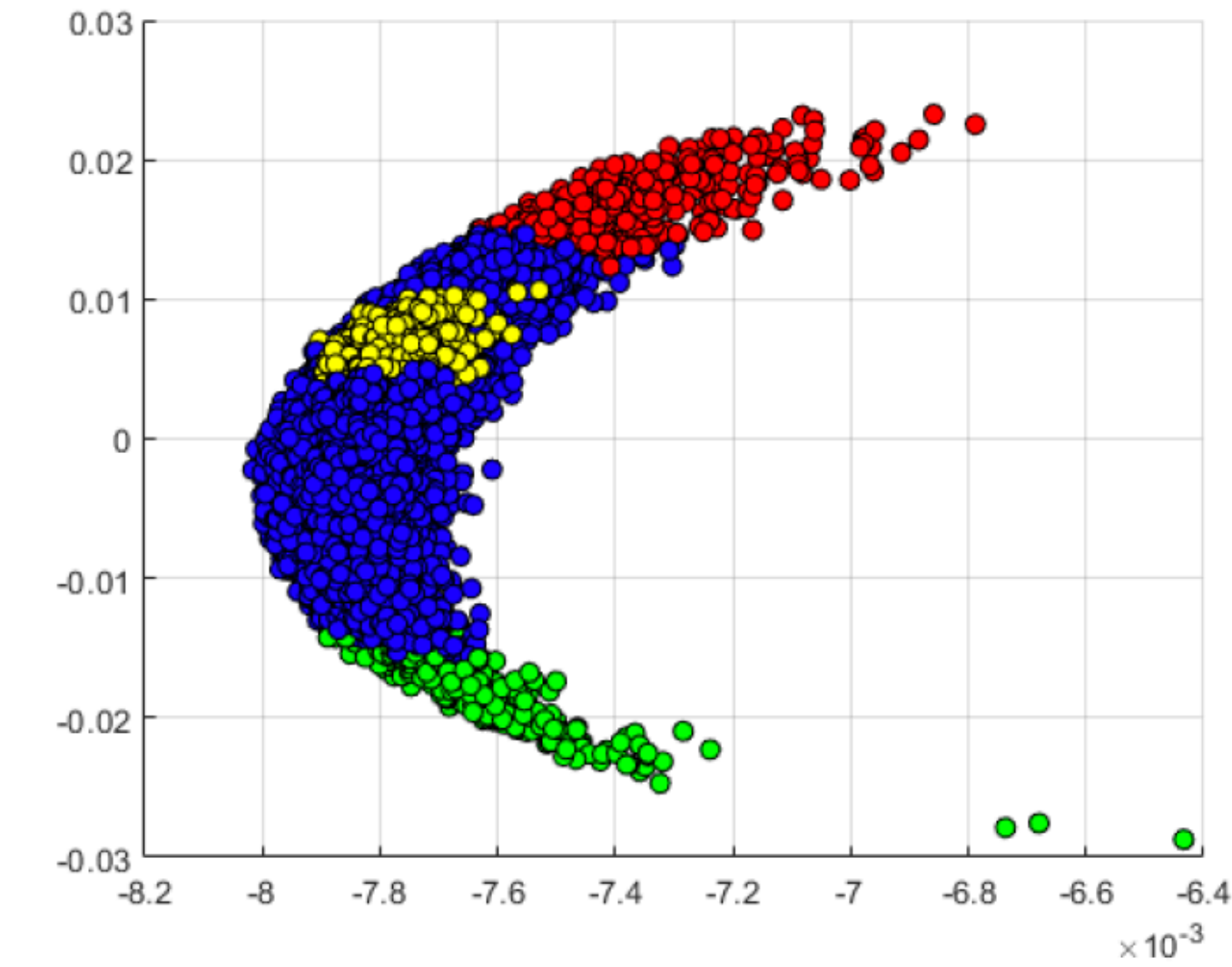


	Normal	AFib	Other	Noise	Avg. Score
Zabihi et. al. [33]	90.49	79.43	75.64	61.11	81.85
Datta et. al. [9]	91.00	79.00	77.00	—	—
Hong et. al. [15]	92.04	86.92	80.68	81.56	85.30
Spectrogram	76.78	43.08	44.71	54.55	54.78
Signal	94.84	96.77	91.93	90.41	93.49
Naïve Attractor	93.93	85.71	94.40	94.87	92.23
01 PCA Attractor	99.66	98.95	98.46	100	99.27

Multimodal Classification



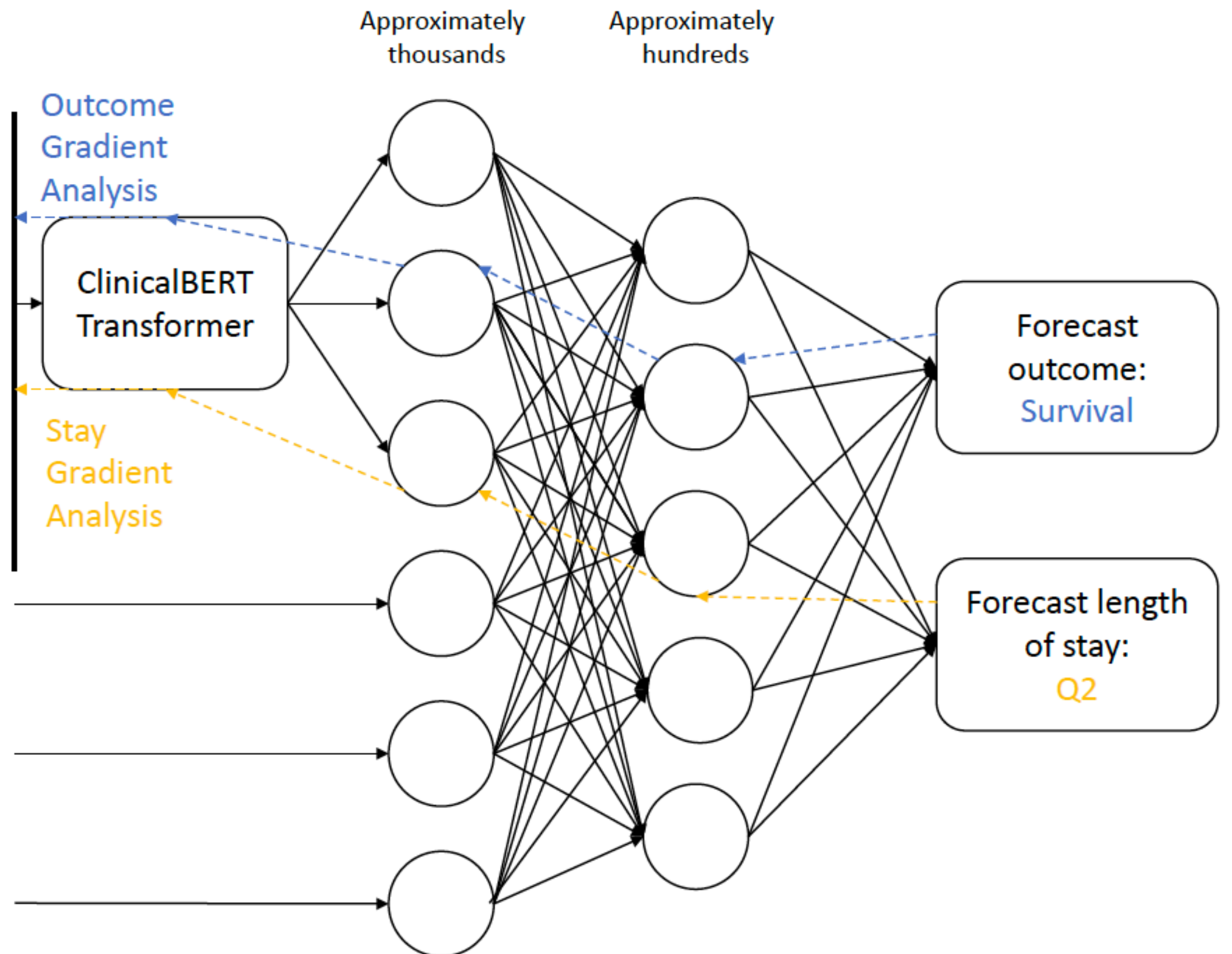
Health Outcome Disparities on the EHR



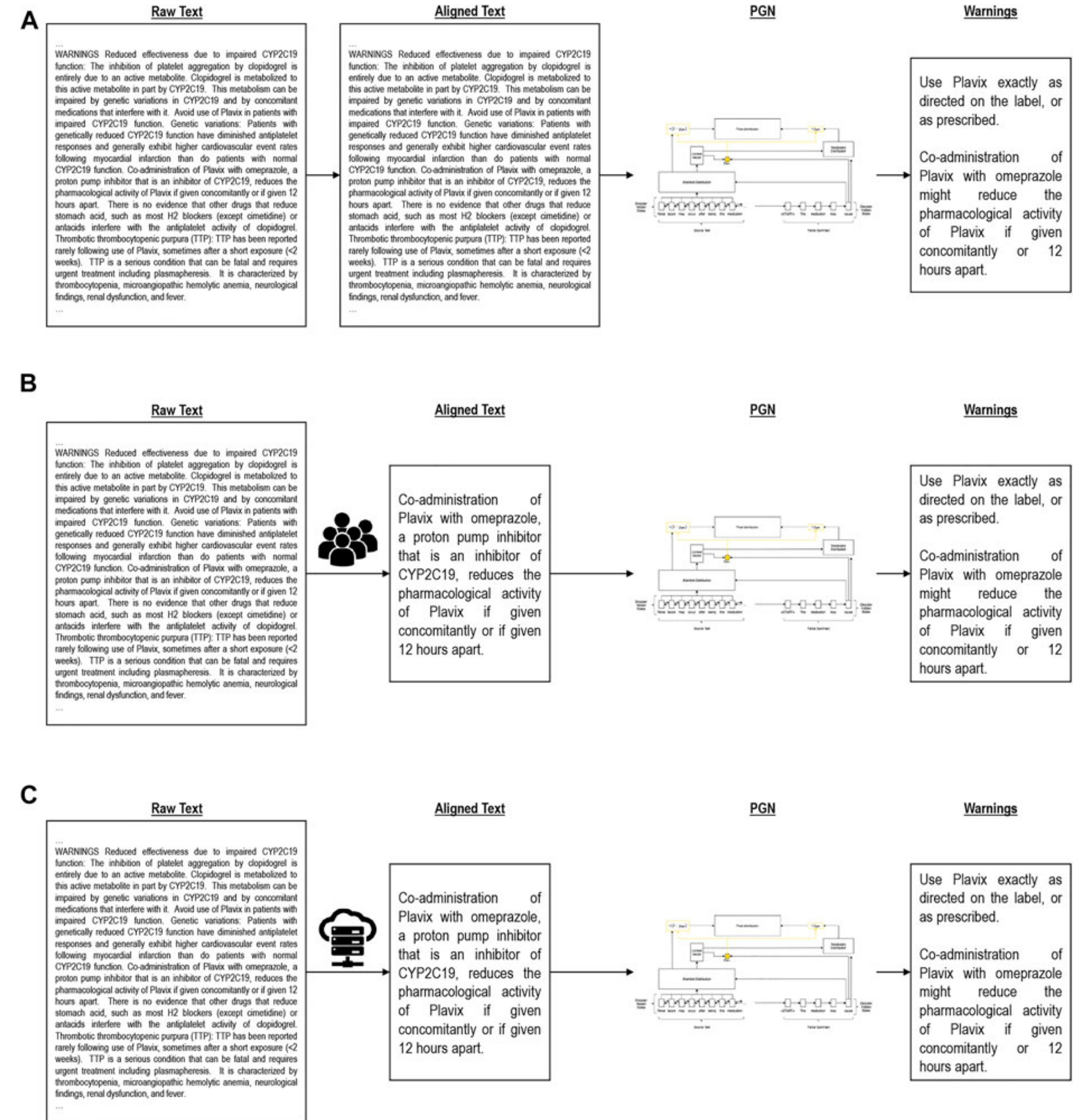
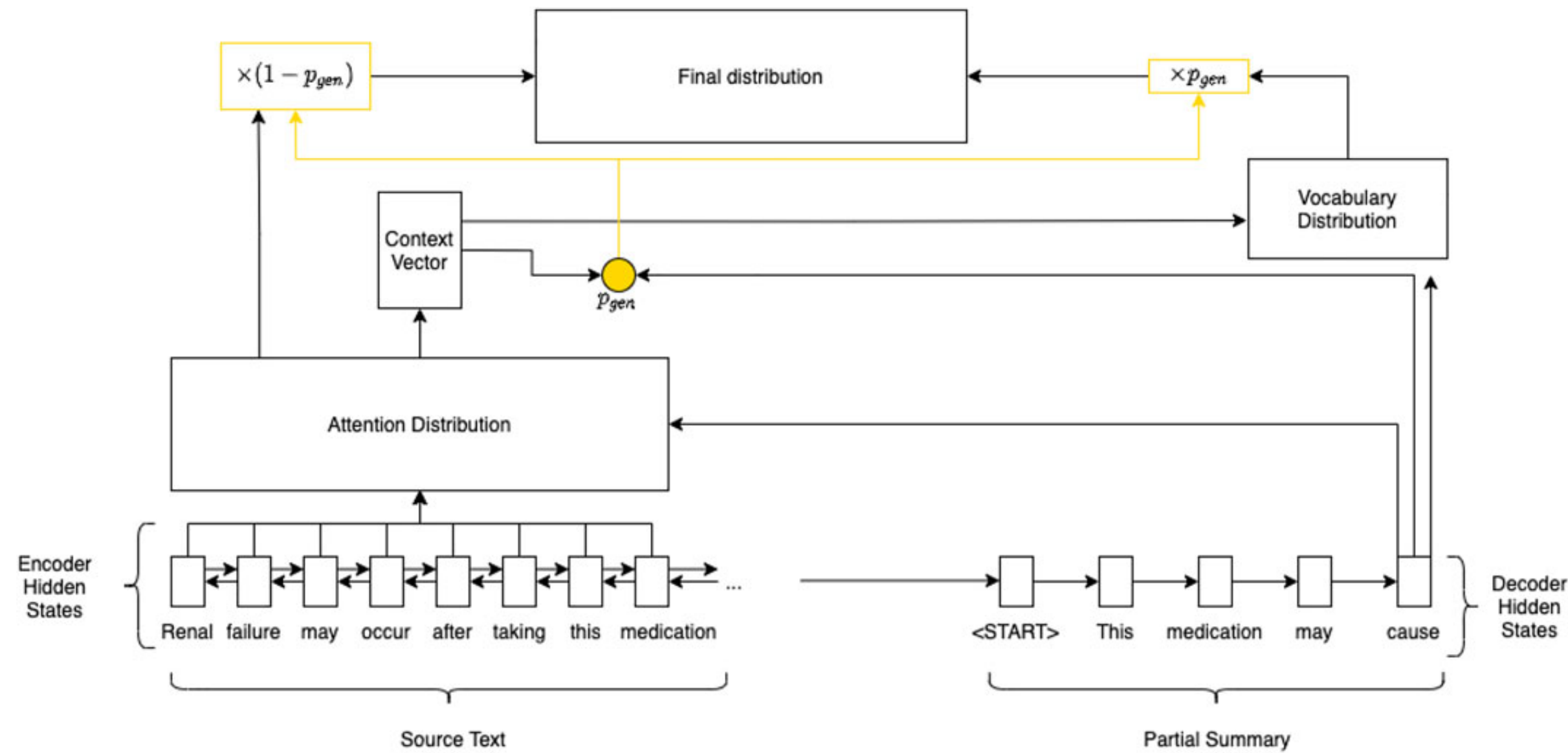
- Self Attention Contextual Embeddings

the patient reported that he had been feeling well without chest pain, shortness of breath, or dyspnea on exertion. The patient underwent a cardiac catheterization on the morning of arrival with pci to the native rca and stents and brachytherapy to the vein graft. the patient tolerated the procedure well and approximately hours later developed a chest pain noted as out of substernal radiating to his throat and back without shortness of breath, diaphoresis, nausea or vomiting. ekg at that time revealed st elevation in ii, iii, and avf

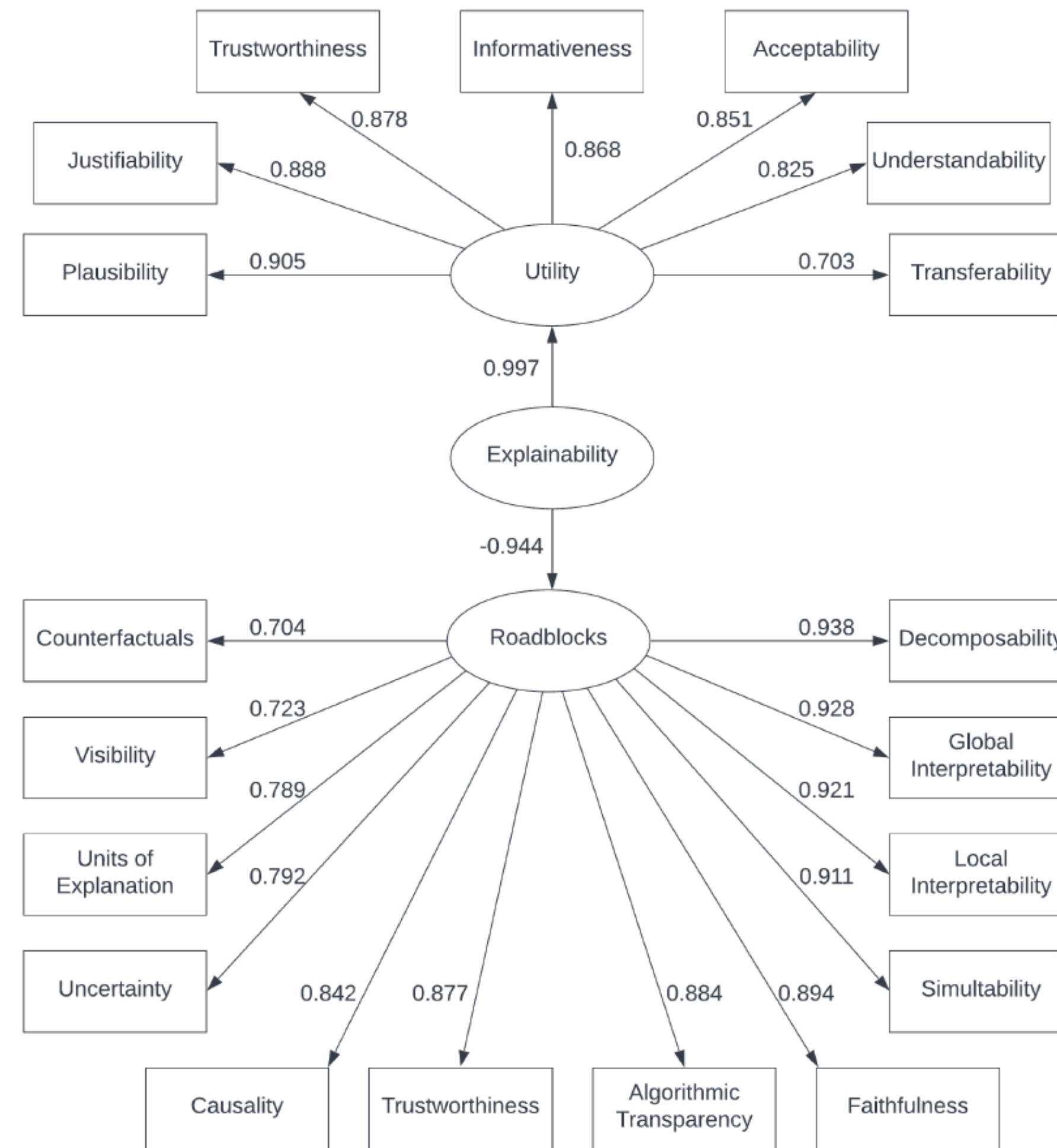
- Age, sex, insurer...
- Comorbidity 1 yes/no
- ...
- Comorbidity n yes/no



Generating Drug Leaflets



Measuring System Explainability



Factors & Items	Original Factor Label	Std. Loading	Std. Error
Factor 1		0.997	0.000
41	Plausibility	0.905	0.096
49	Justifiability	0.888	0.097
47	Trustworthiness	0.878	0.096
21	Informativeness	0.868	0.094
37	Acceptability	0.851	0.110
19	Understandability	0.825	0.087
15	Transferability	0.703	0.000
Factor 2		-0.944	0.065
2	Decomposability	0.938	0.094
22	Global Interpretability	0.928	0.098
24	Local Interpretability	0.921	0.098
0	Simultability	0.911	0.099
38	Faithfulness	0.894	0.097
4	Algorithmic Transparency	0.884	0.097
46	Trustworthiness	0.877	0.091
6	Causality	0.842	0.095
8	Uncertainty	0.792	0.092
34	Units of Explanation	0.789	0.094
12	Visibility	0.723	0.092
26	Counterfactuals	0.704	0.000

Language Generation in a Code-switched World

Hi, 提醒我明天早上10点的
final project presentation. ZH
EN
(Hi, can you remind me of my final project
presentation tomorrow morning at 10?)

Okay done!

- Final Project Presentation
@ 10 AM TOMORROW
- Feed cat
- Do Laundry
- ● ●
-

Tell mom that "Mom you need
to think of one thing at a time,
פרה פרה. (hebrew slang for
one thing at a time)" HE
EN

Message sent!

Mom you need to think of one
thing at a time, פרה פרה.

Read Liam's new message.

Liam said, "After this I'm just
gonna go home drink summ hot
chocolate con bolillo and sleep." ES
EN

Are Language Models World Models?

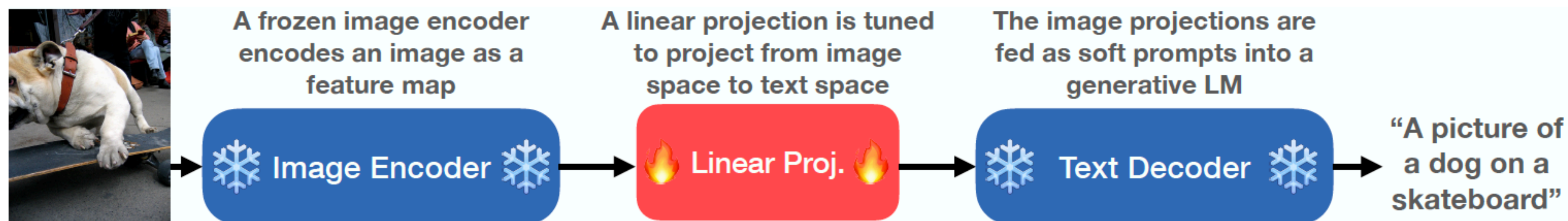
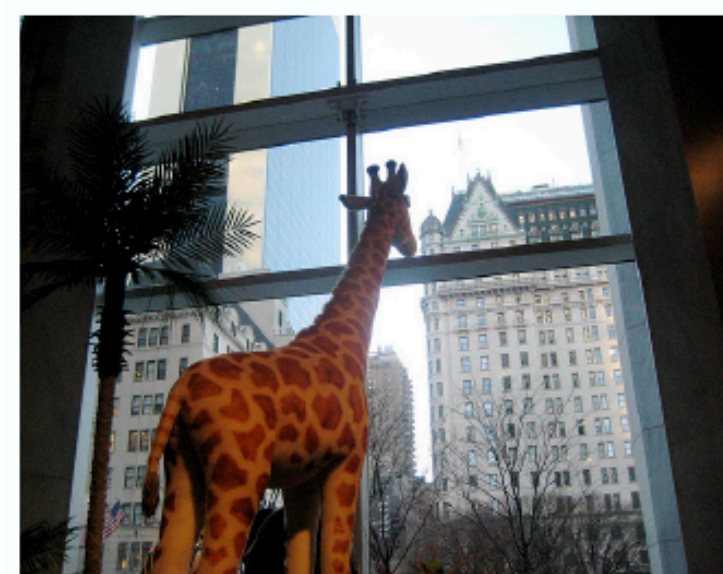


Image Captioning



CLIP	a giraffe in the lobby of the building
NFRN50	the giraffe in the zoo.
BEIT	a peacock in the garden
NFRN50 Random	a man and a woman in a field of flowers



CLIP	tennis player in action
NFRN50	tennis player at the tennis tournament.
BEIT	tennis player during a tennis match.
NFRN50 Random	the new logo for the team

Visual Question Answering



CLIP	He is surfing a wave.
NFRN50	He is surfing the waves.
BEIT	He is jumping into the water.
NFRN50 Random	He is swimming in the pool.

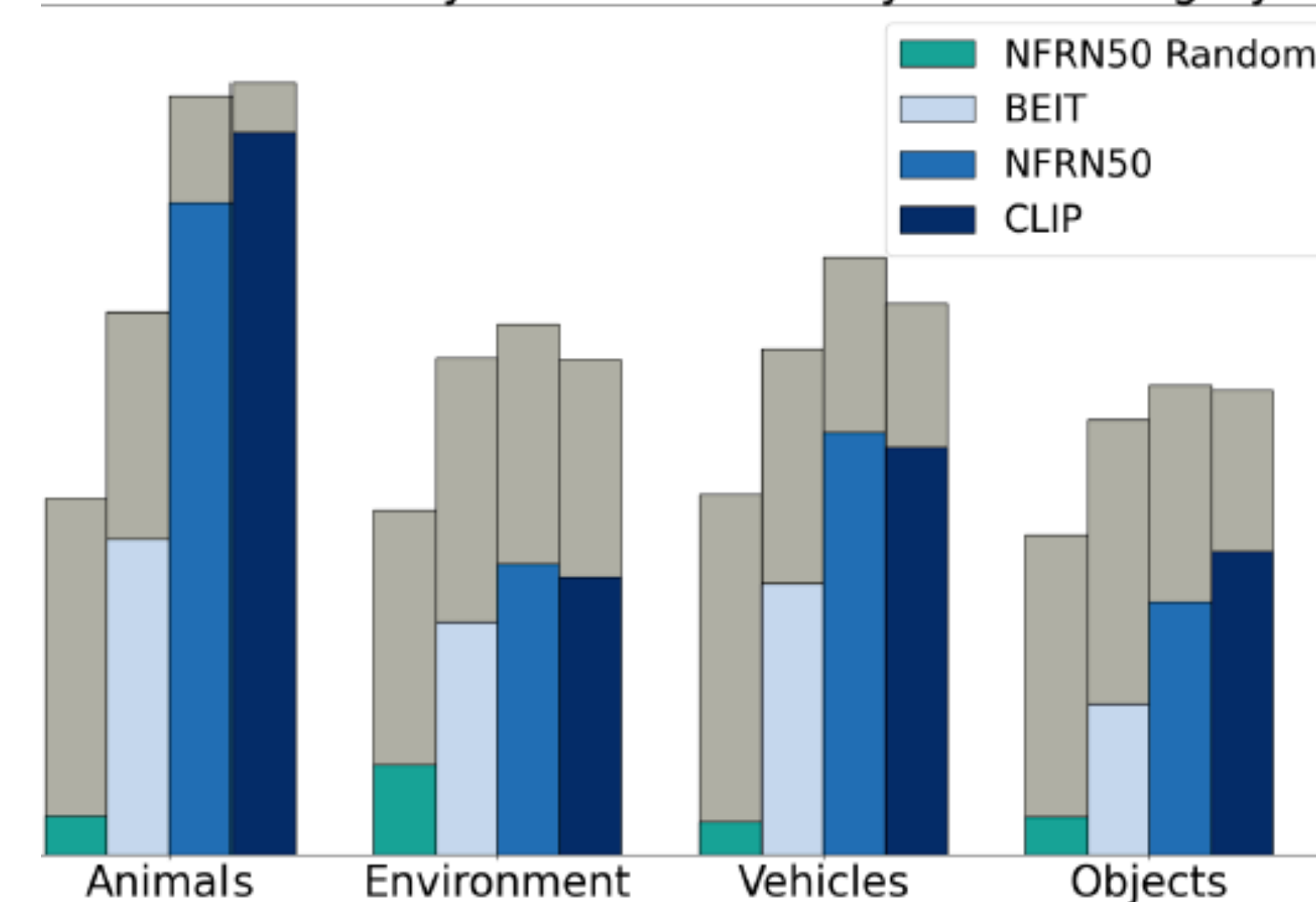
Q: What is the person doing?
A: surfing



CLIP	A tennis racket
NFRN50	A tennis racket
BEIT	A baseball bat.
NFRN50 Random	A tree

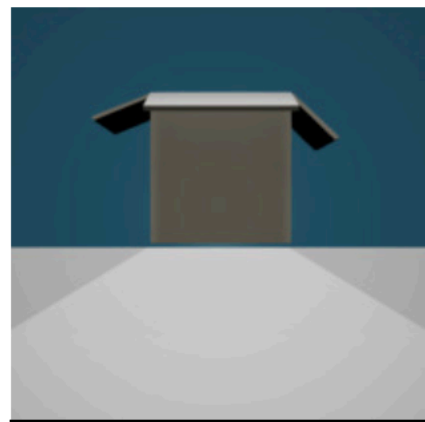
Q: What is the person holding?
A: tennis racket

ℓ_1 -Palmer Similarity for each Model by Noun Category

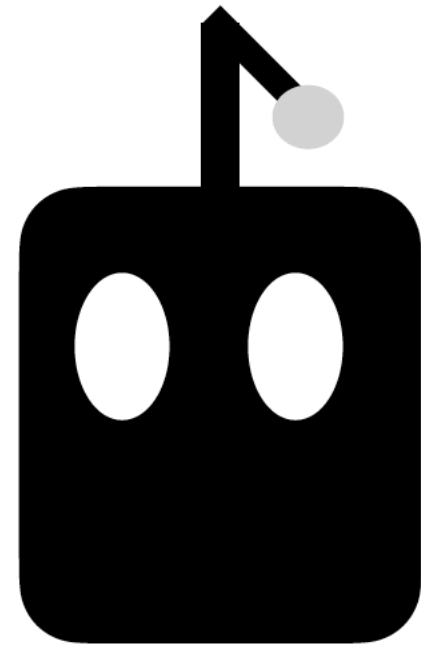


Grounding Language Models in Physics

Play with these objects!



Seen



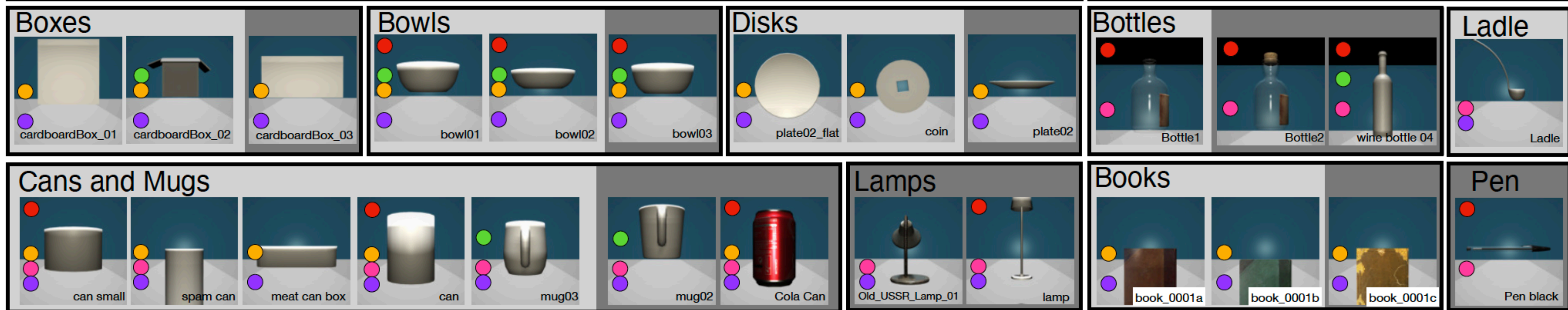
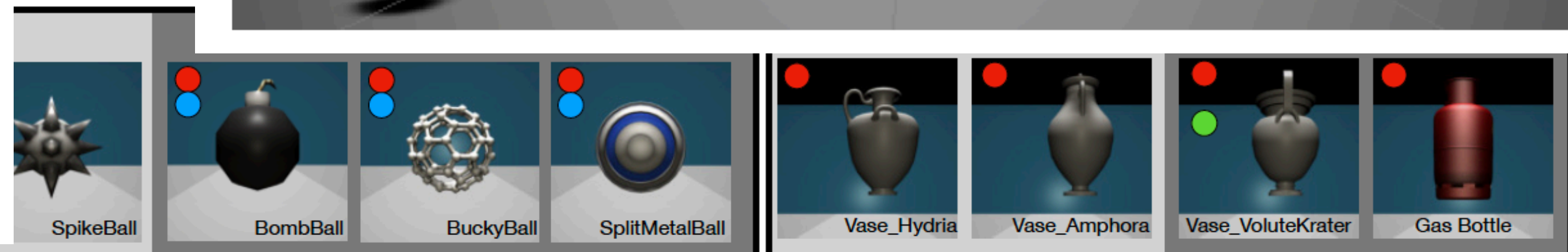
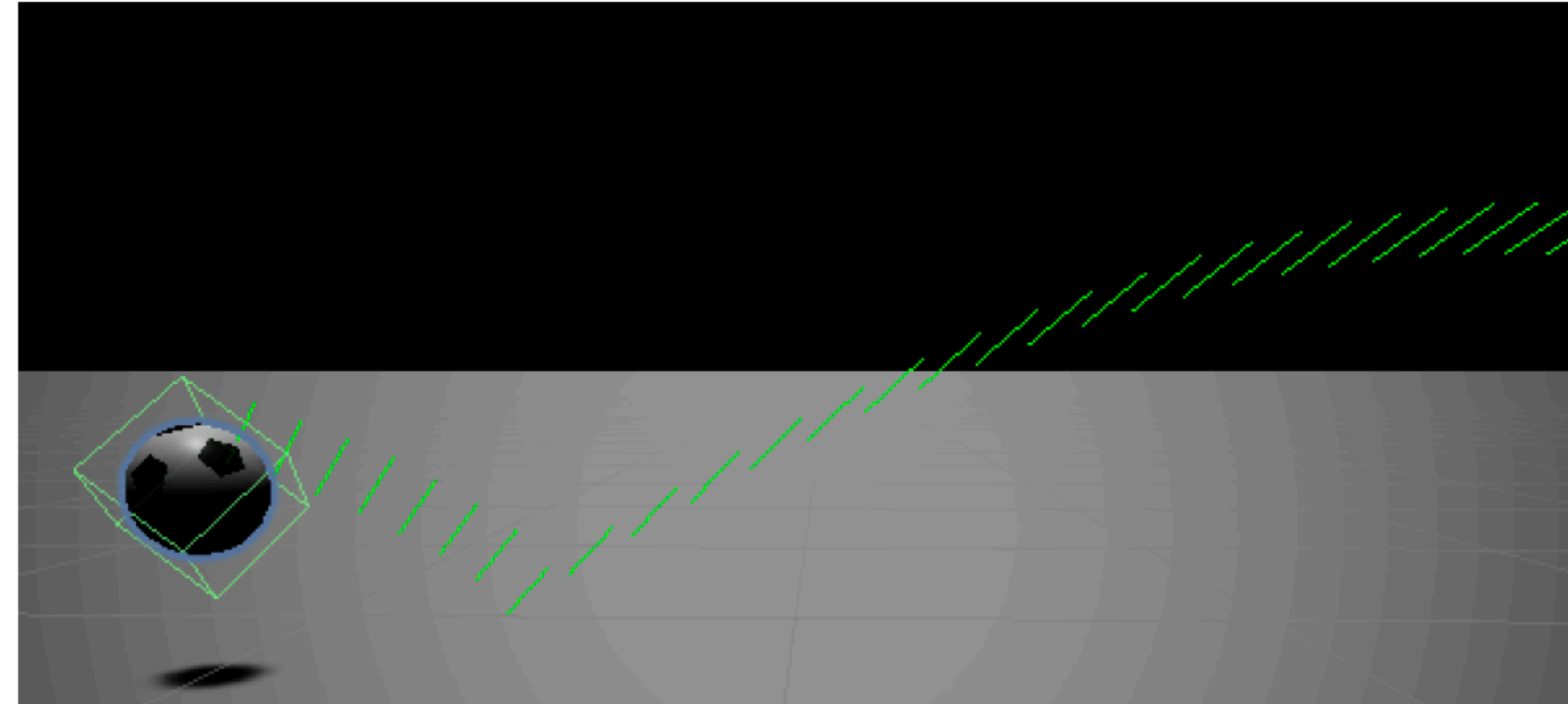
Now play with this one!



Can it roll?

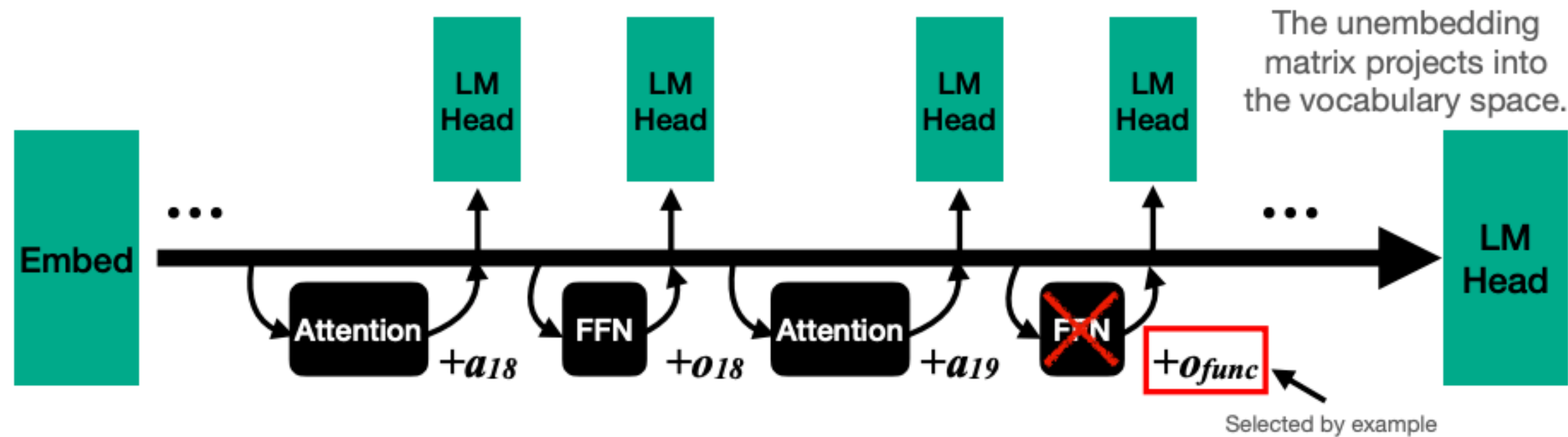
Can it contain other objects?

Unseen



● Roll ● Bounce ● Contain ● Stack ● W-Grasp ● Slide

LLM Vector Arithmetics



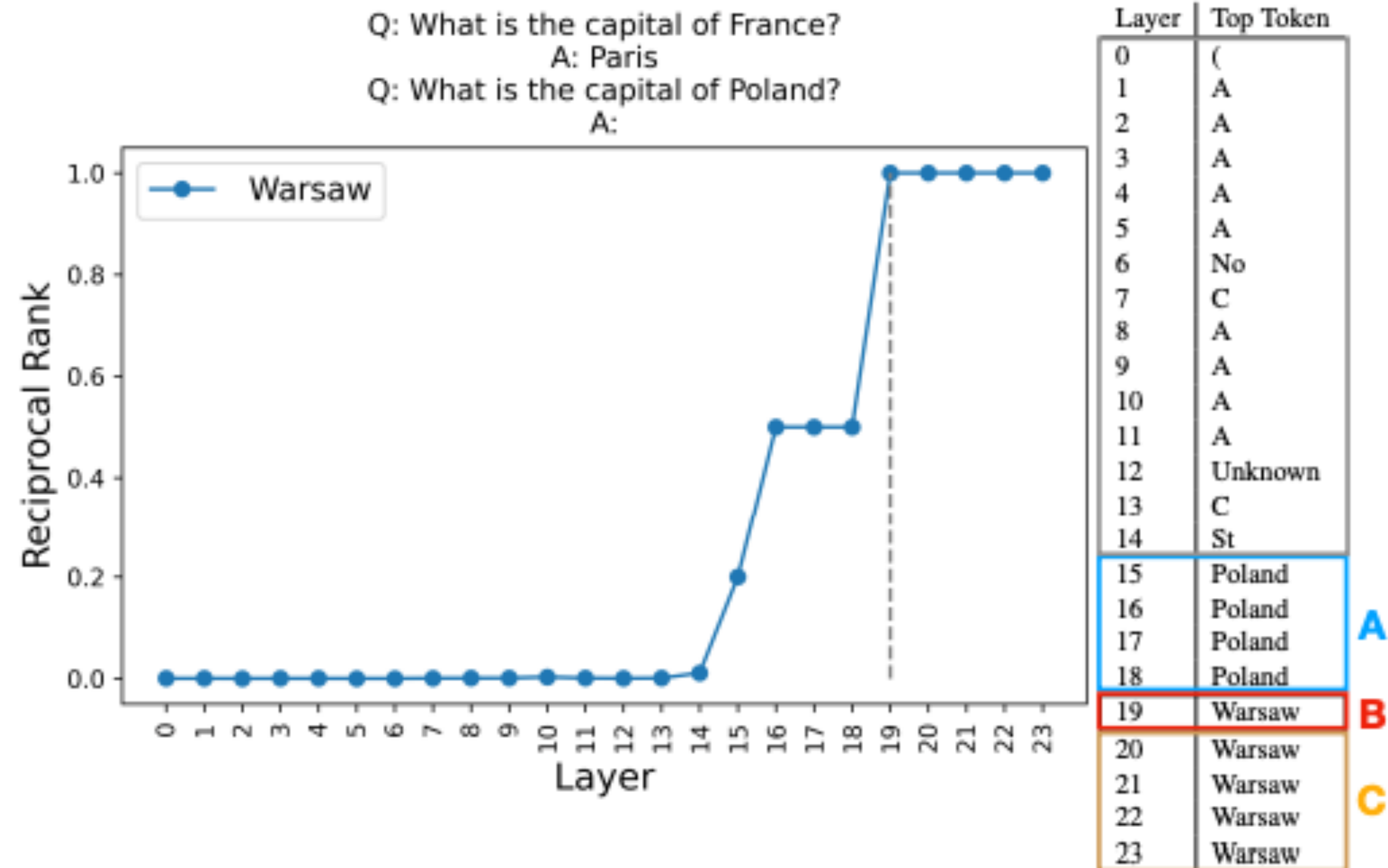
- *We find simple vector arithmetic in Transformers*

- **3 Models**

- **GPT2 (124M)**

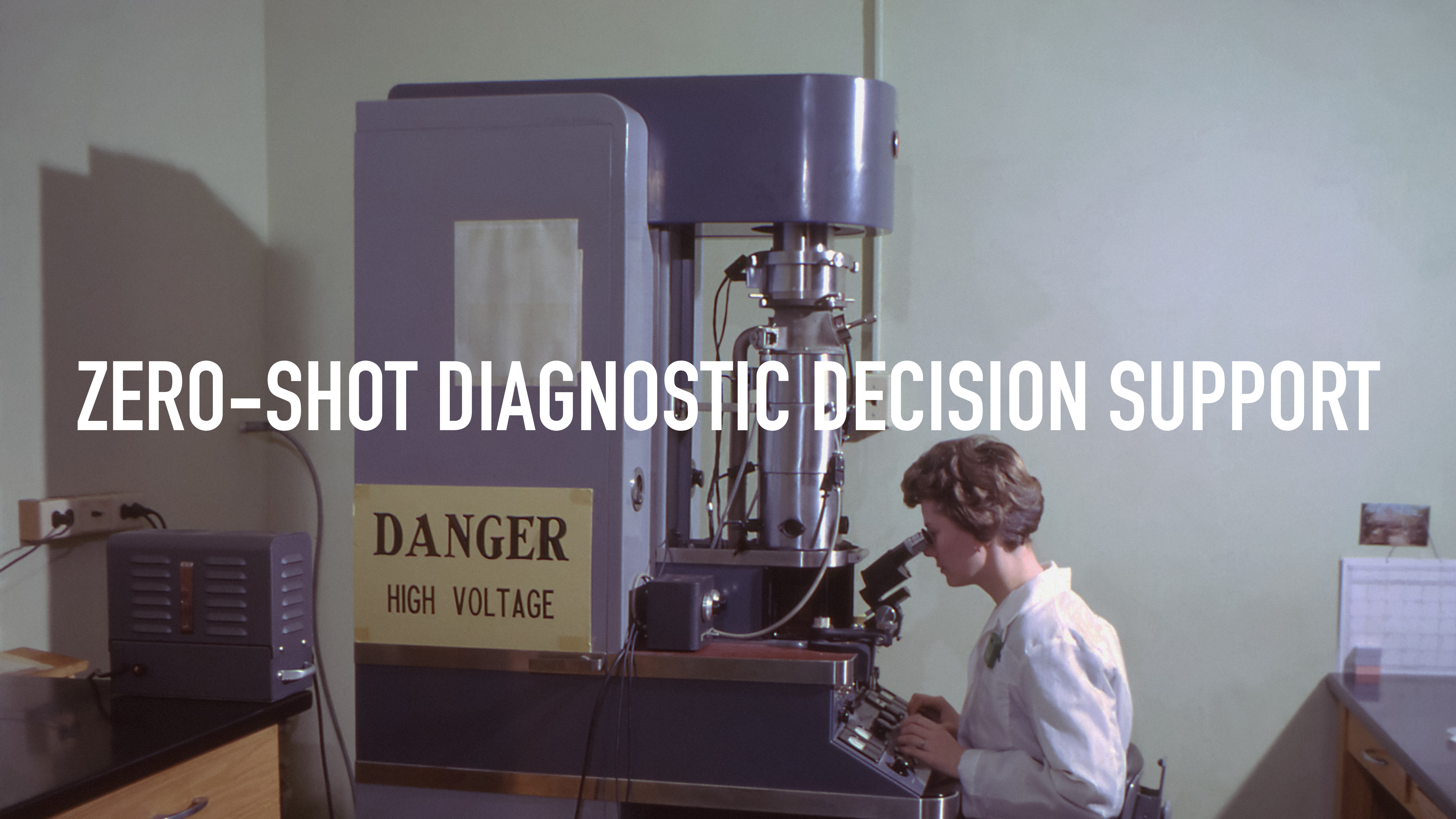
- **GPT-J (6B)**

- **Bloom (176B)**



ZERO-SHOT DIAGNOSTIC DECISION SUPPORT

DANGER
HIGH VOLTAGE



Surgical &
Medication
Errors

5%
of outpatient
office visits

10%
of hospital
inpatient deaths

Diagnostic Errors

12%
of hospital
adverse events

18 MILLION
diagnostic **ERRORS** each year

74,000
deaths each year

“ Nearly every person will experience
a **diagnostic error** in their lifetime ”

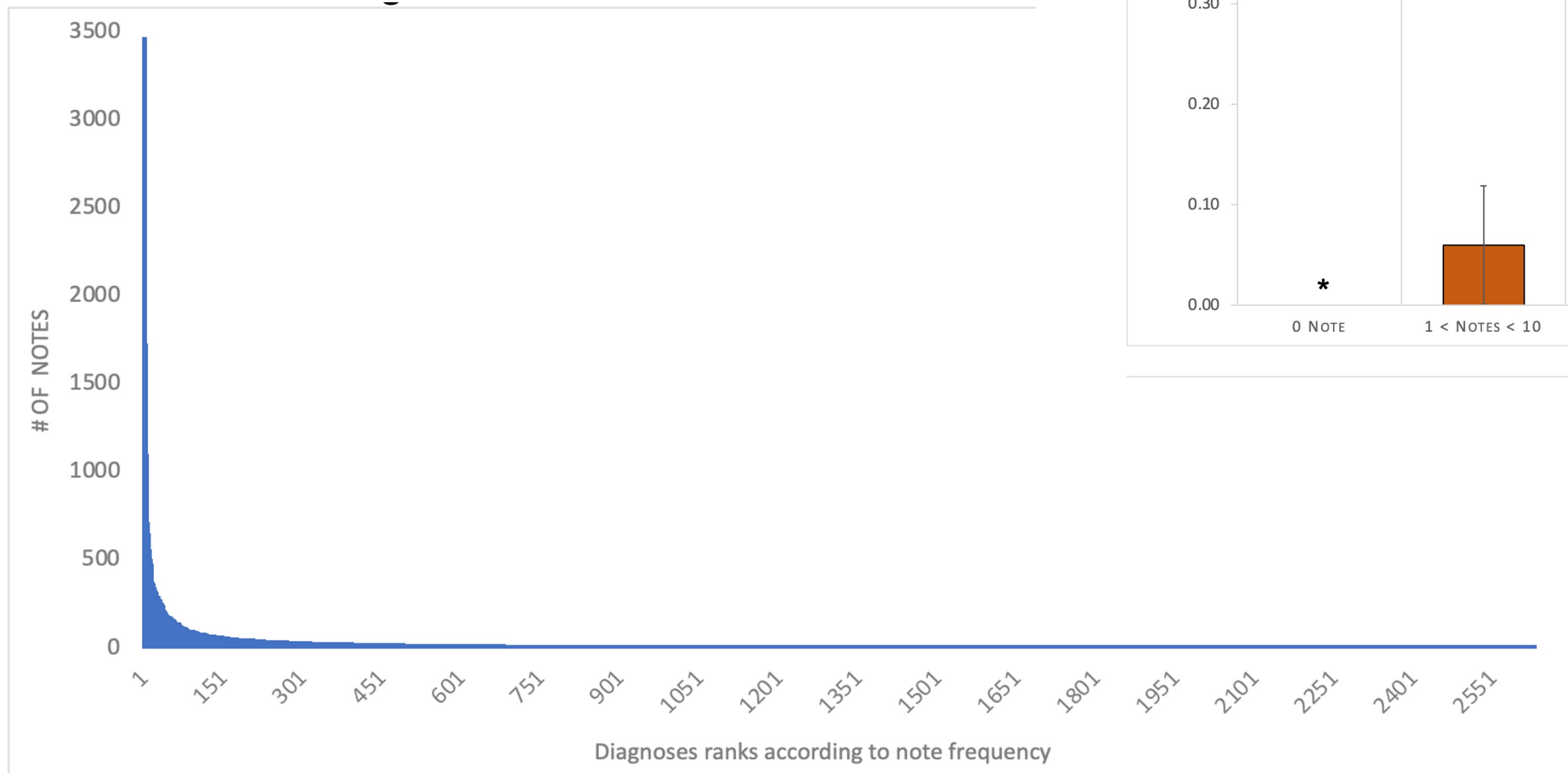
ML to the Rescue

- *Golovanevsky et al. 2022: Alzheimers (92.28%)*
- Delahanty et al. 2018: Sepsis (97%)
- *Gulshan et al. 2016: Diabetic Retinopathy (99.1%)*
- *Rudman et al. 2022: Cardiac Arrhythmias (99.27%)*
- ...

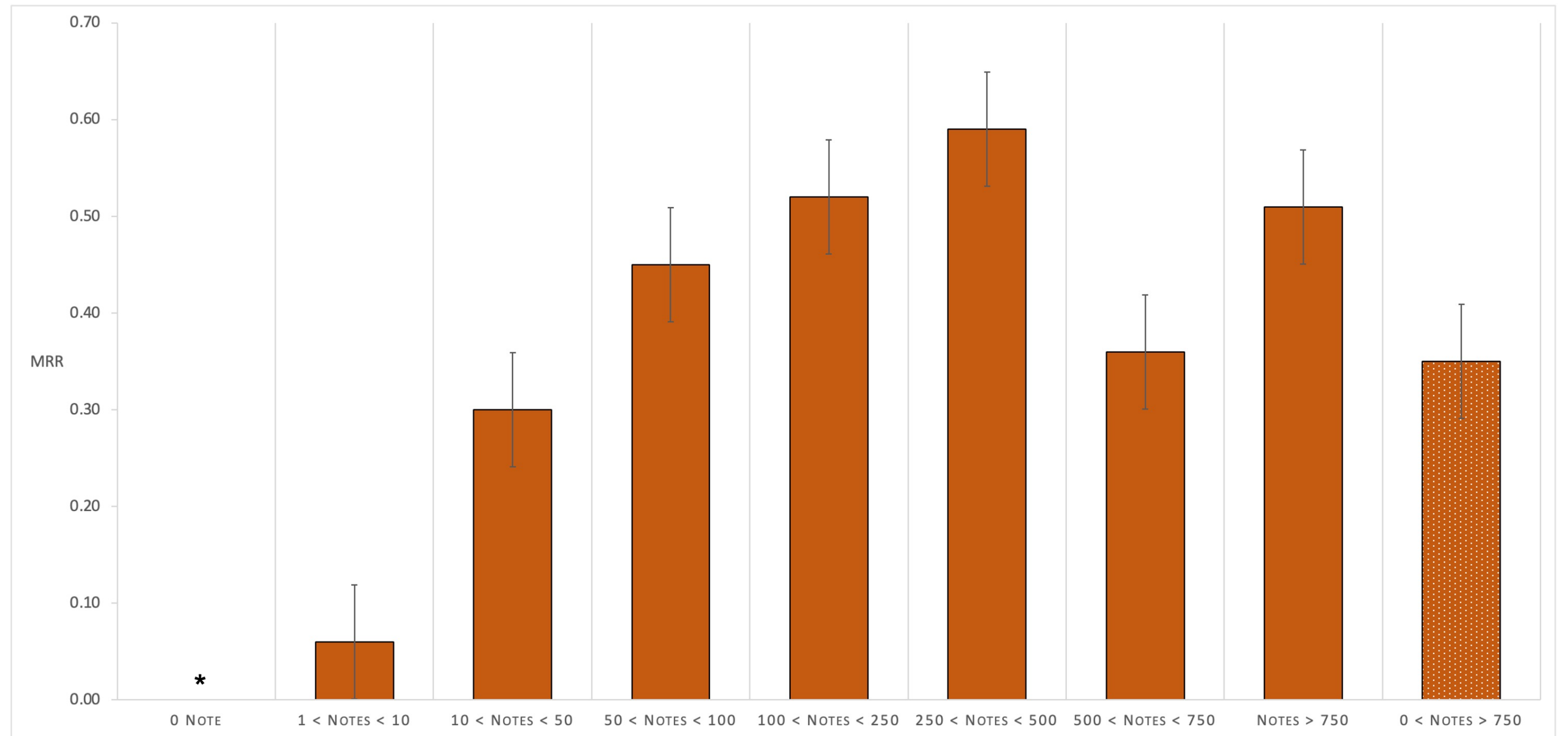
ML to the Rescue

- *Golovanevsky et al. 2022*: Alzheimers (92.28%) [n= 2,384]
- Delahanty et al. 2018: Sepsis (97%) [n= 2,759,529]
- *Gulshan et al. 2016*: Diabetic Retinopathy (99.1%) [n= 128,175]
- *Rudman et al. 2022*: Cardiac Arrhythmias (99.27%) [n= 8,528]
- ...

But only if you have the Data!

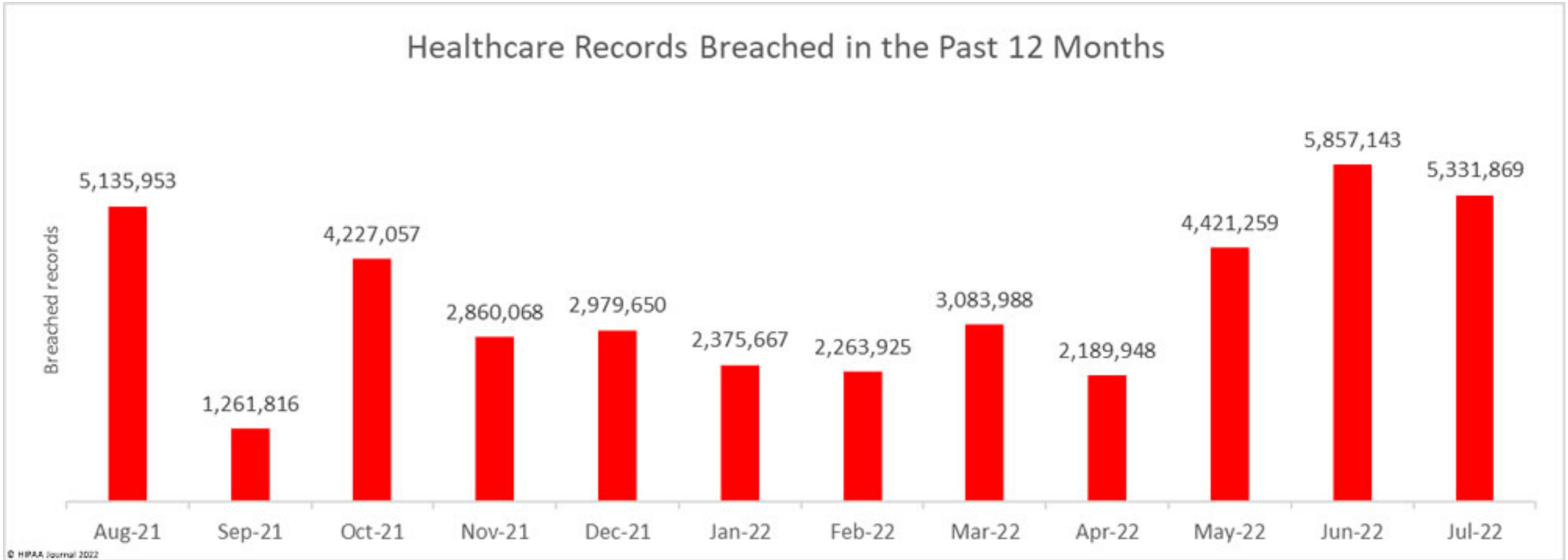


A.

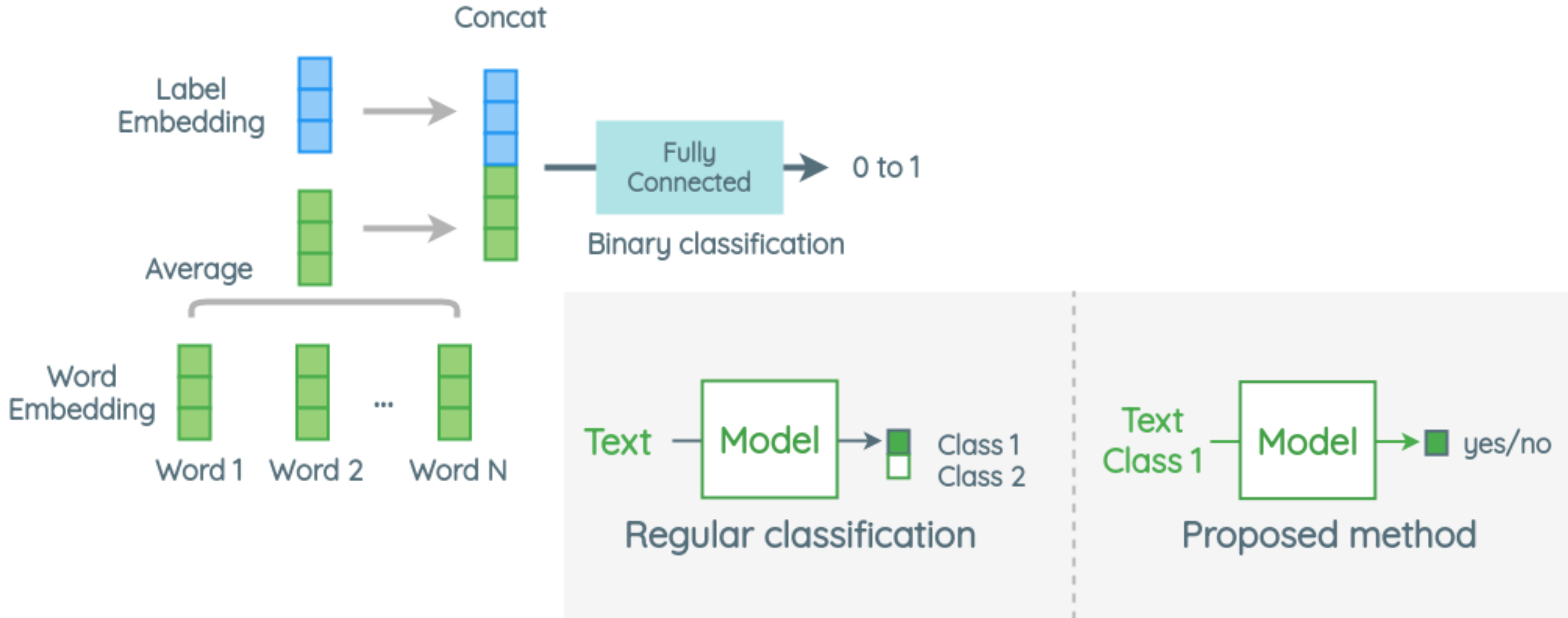


Sharing is ...

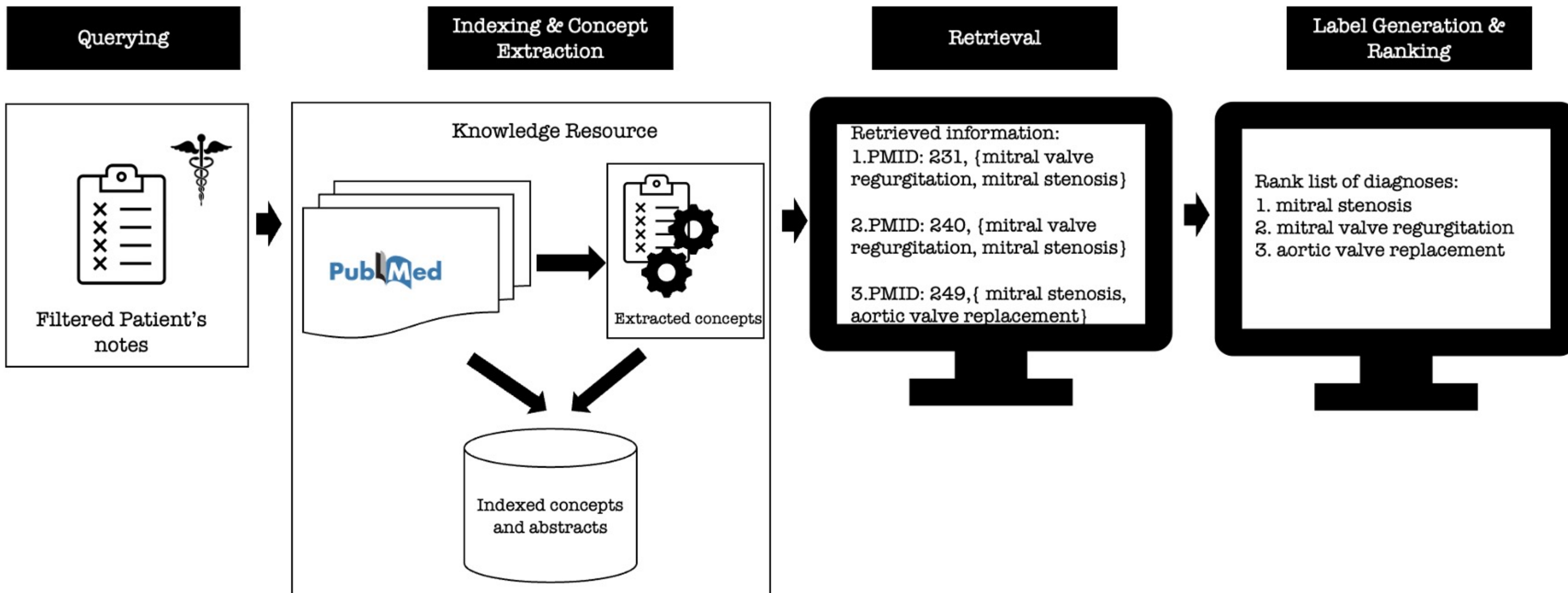
Healthcare Records Breached in the Past 12 Months



Zero-shot Text Classification



Via Retrieval

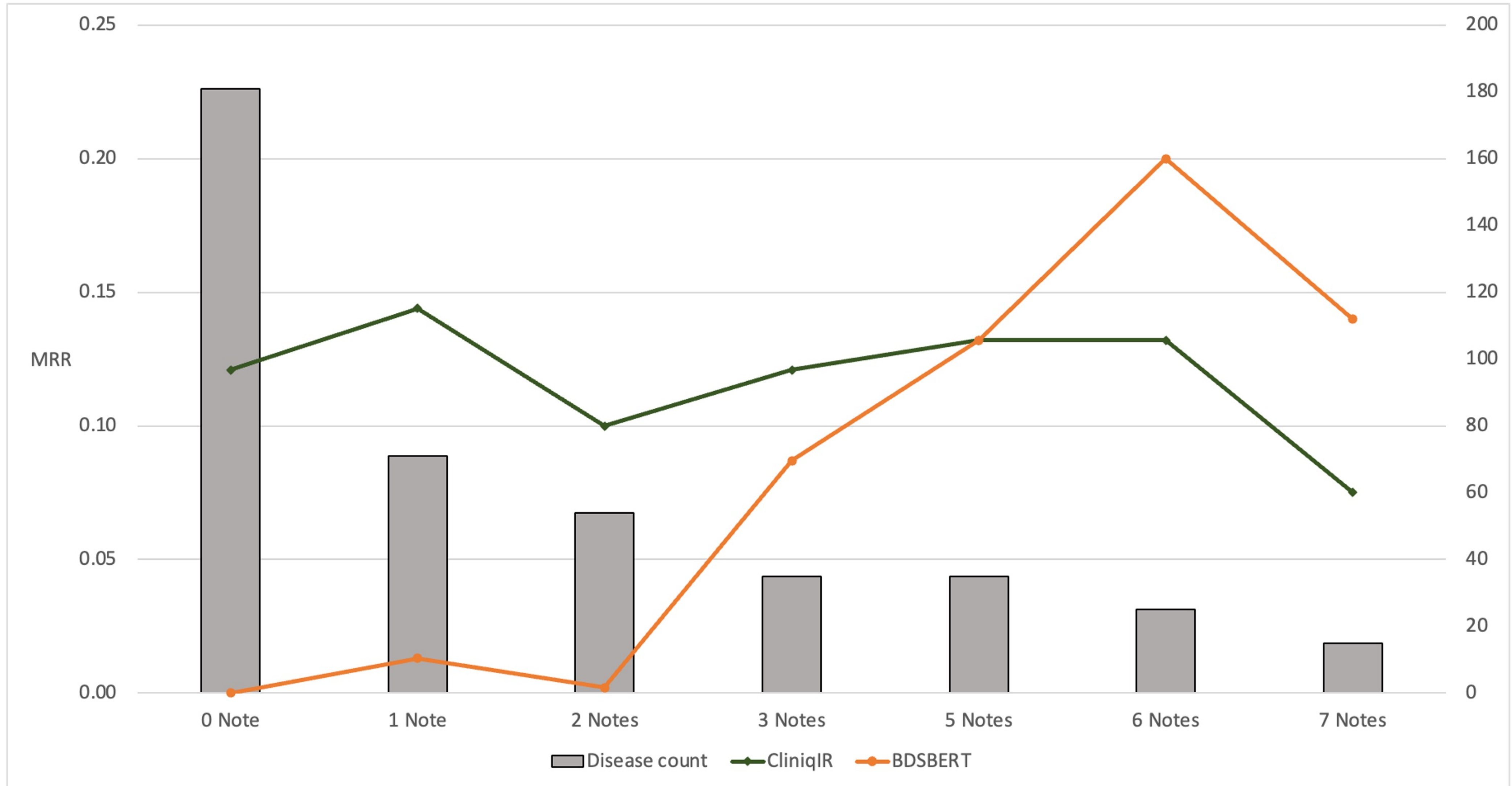


Datasets

- PubMed
 - 33M abstracts
- MIMIC III
 - 50,000 ICU encounters
 - 2643 unique diagnoses
 - 902 of them occur exactly once
- DC3
 - 31 (difficult) case reports

Results I

A.



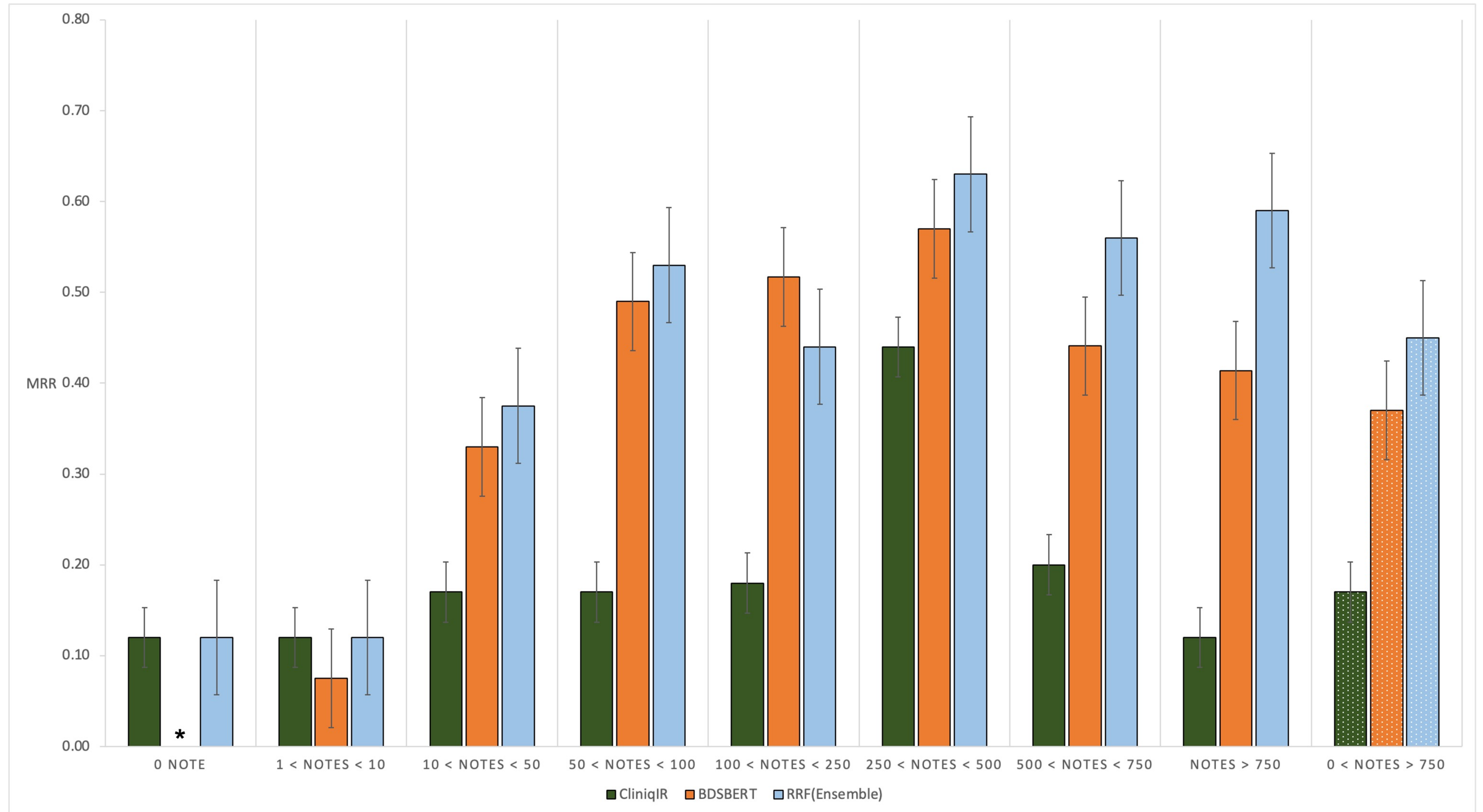
Ensembling

$$P'(d, e) = \frac{P_s(d, e) + \mu P_u(d, e)}{|d| + \mu}$$

- d : diagnosis
- e : encounter
- μ : Dirichlet prior, median # of training examples per diagnosis
- $|d|$: # of training examples for diagnosis d

Results II

B.



Reasons to Care

- The world is Zipfian
- Effective zero-shot approaches rely on structured information (KB triples)
- Unstructured data is abundantly available and growing fast
- Using unstructured collections for unsupervised learning unlocks considerable resources
- Next stop: Going beyond diagnostics

DISCUSSION

