

# Language Models in Complex Diagnostics

Carsten Eickhoff

EBERHARD KARLS  
UNIVERSITÄT  
TÜBINGEN

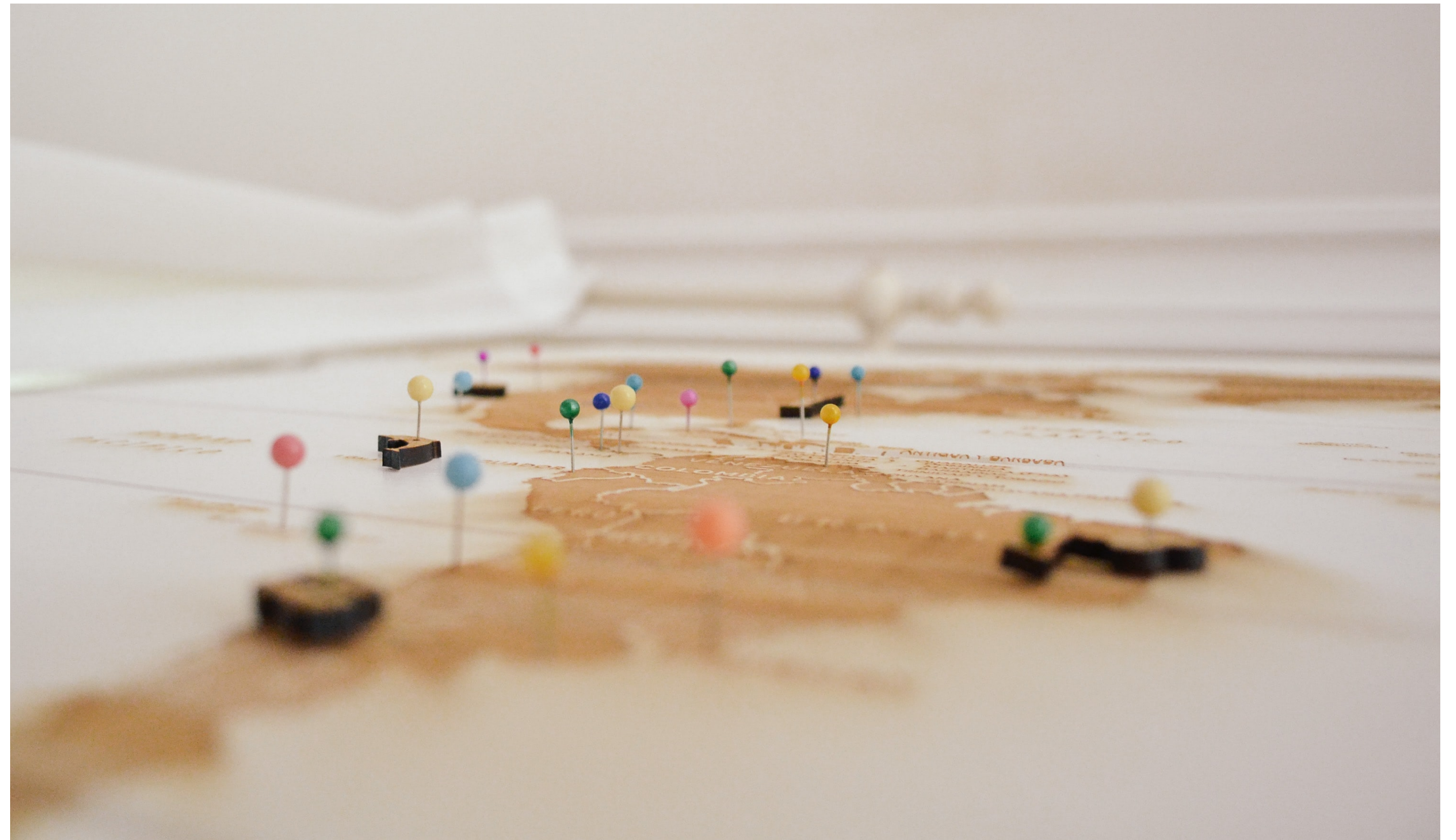


BROWN



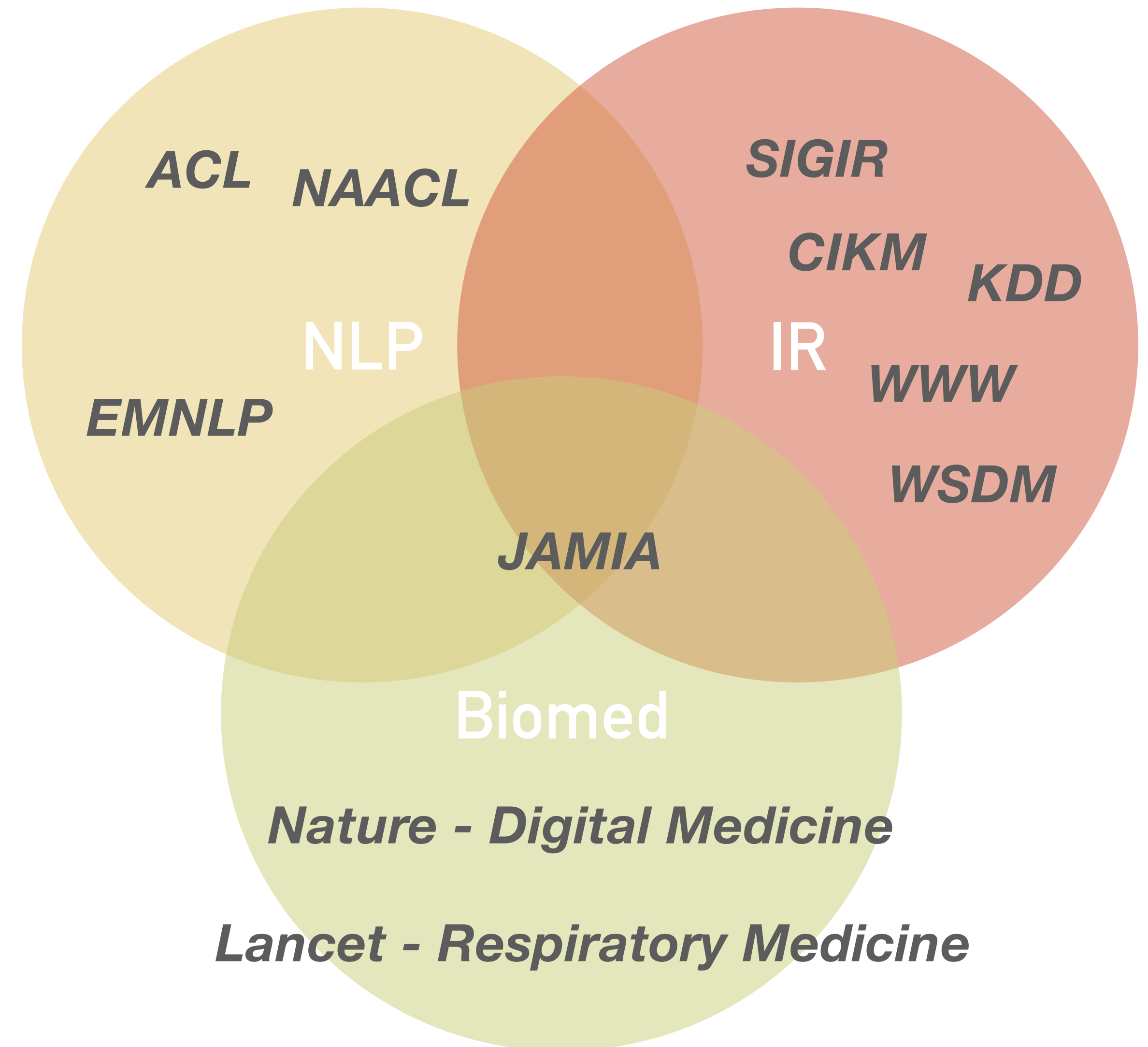
# About

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



# Interests

- **Dense Retrieval**
- **Explainable IR**
- **Uncertainty Aware IR**
- **Grounded Language Modeling**
- **Manifold Learning for Neural LMs**
- **Conditional Text Generation**
- **Clinical Decision Support**





# Team



**CARSTEN EICKHOFF**  
PROFESSOR



**GEORGE ZERVEAS**  
PHD STUDENT



**AMINA ABDULLAHI**  
PHD STUDENT



**LAURA MERCURIO MD**  
ASSISTANT PROFESSOR



**JACK MERULLO**  
PHD STUDENT



**CATHERINE CHEN**  
PHD STUDENT



**ALI BAHRAINIAN**  
POSTDOC



**MICHAL GOLOVANEVSKY**  
PHD STUDENT



**WASIWASI MGONZO**  
PHD STUDENT



**AUGUSTO GARCIA-AGUNDEZ**  
POSTDOC



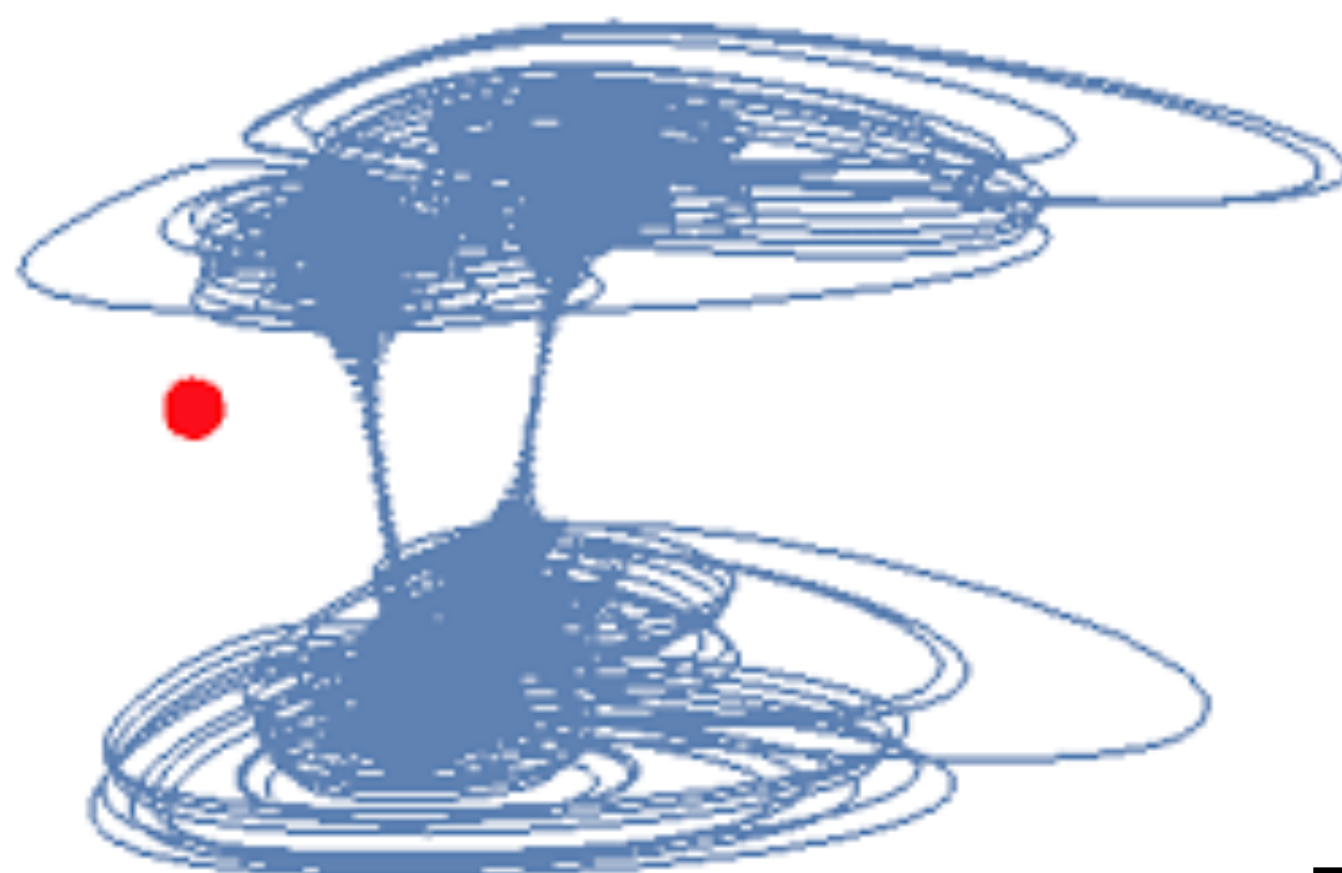
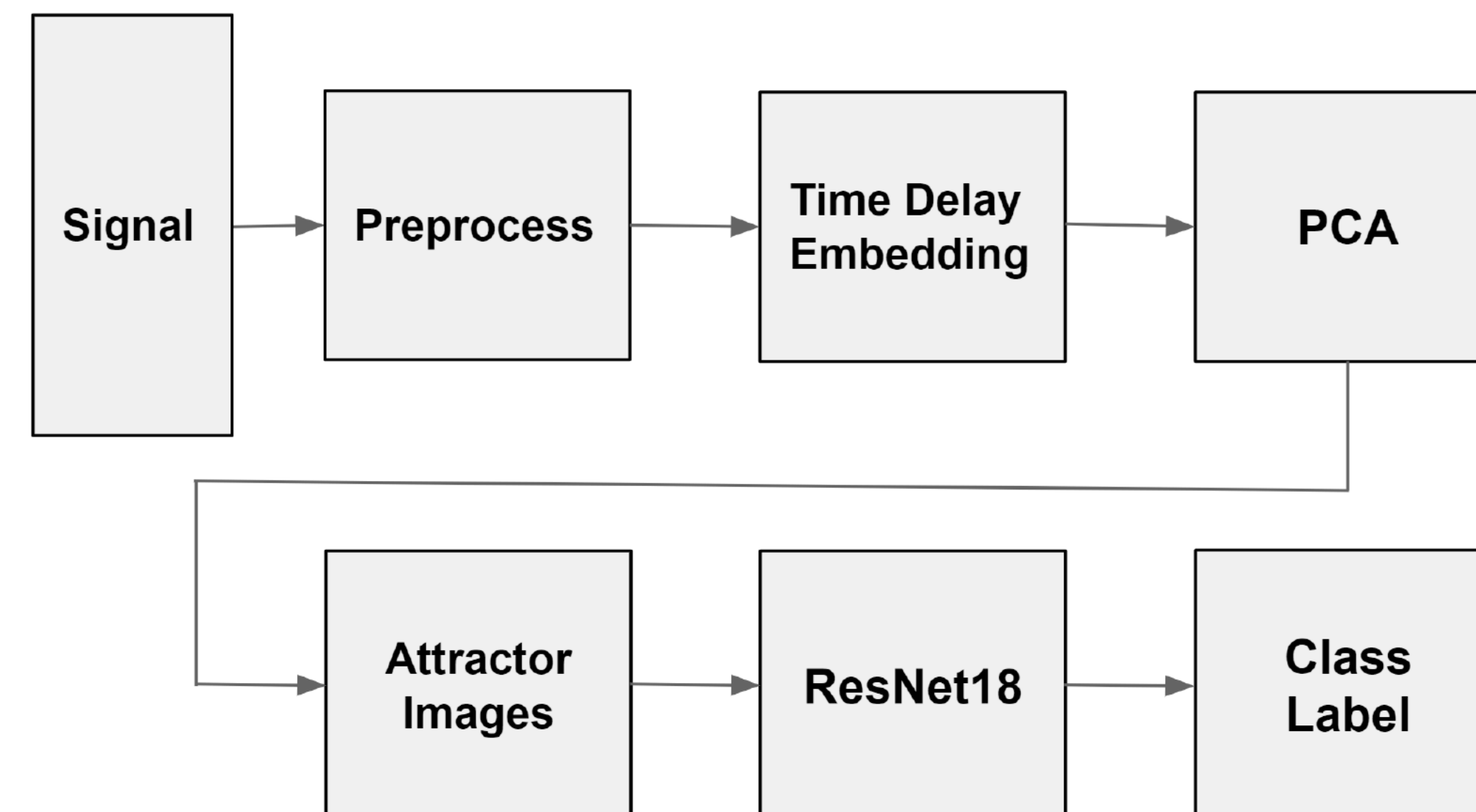
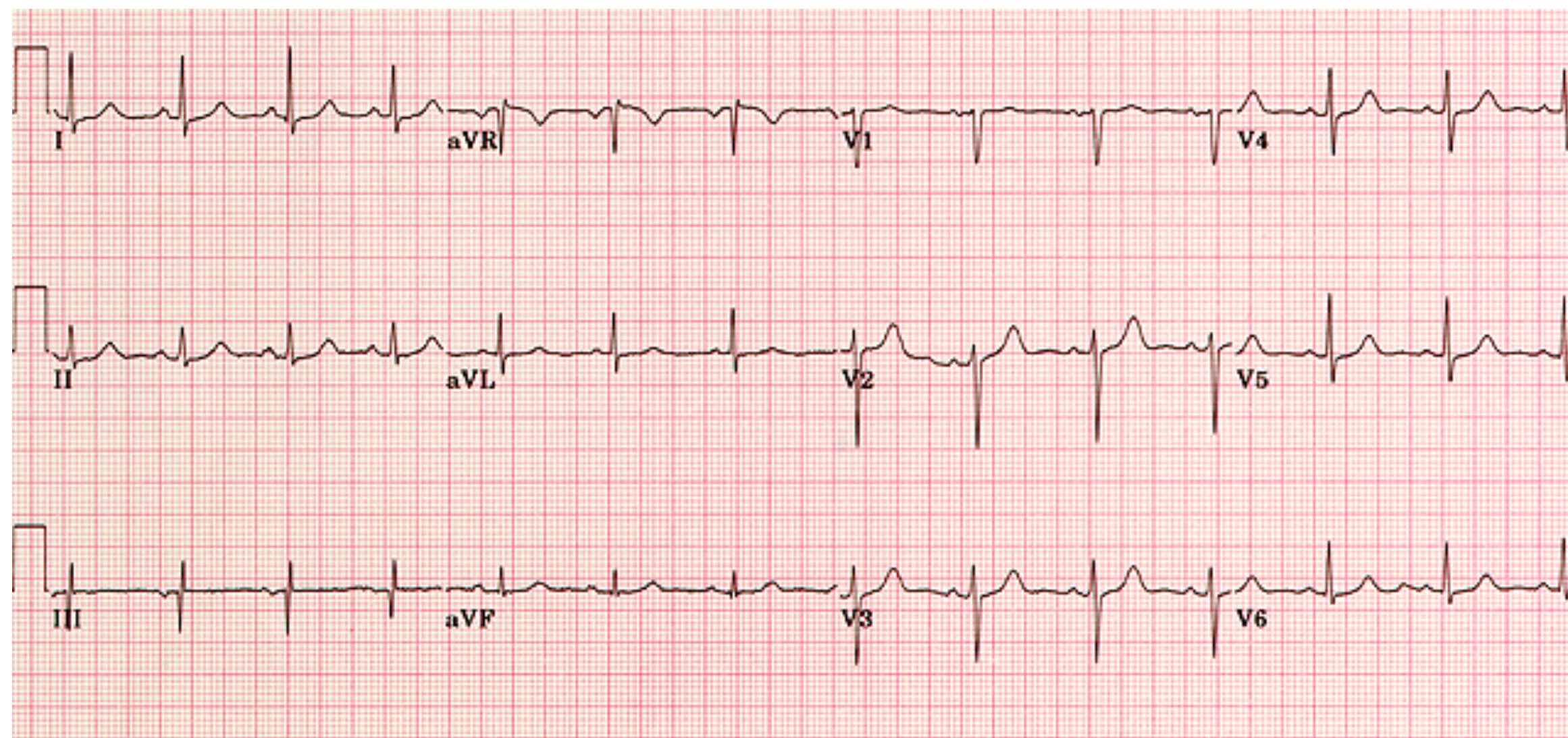
**RUOCHEN ZHANG**  
PHD STUDENT



**WILLIAM RUDMAN**  
PHD STUDENT



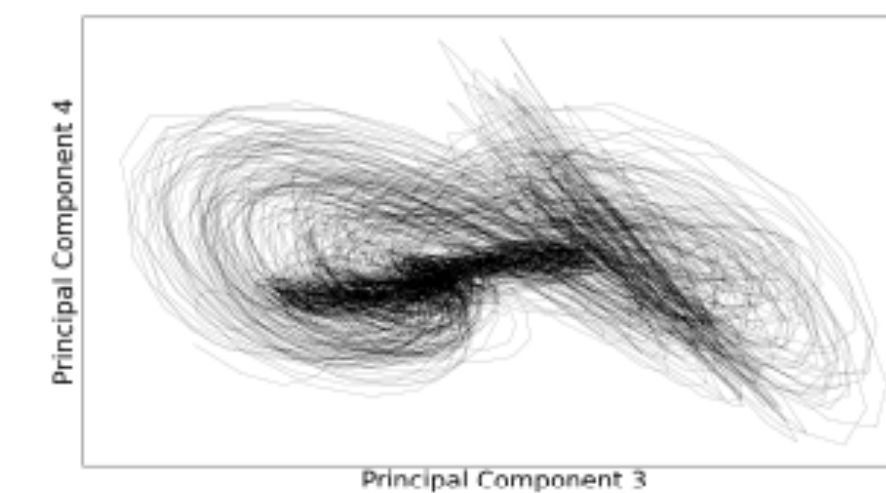
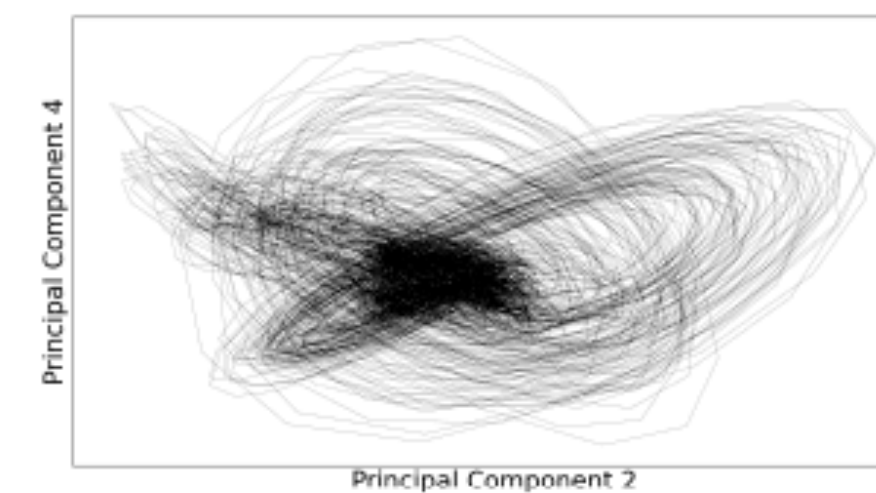
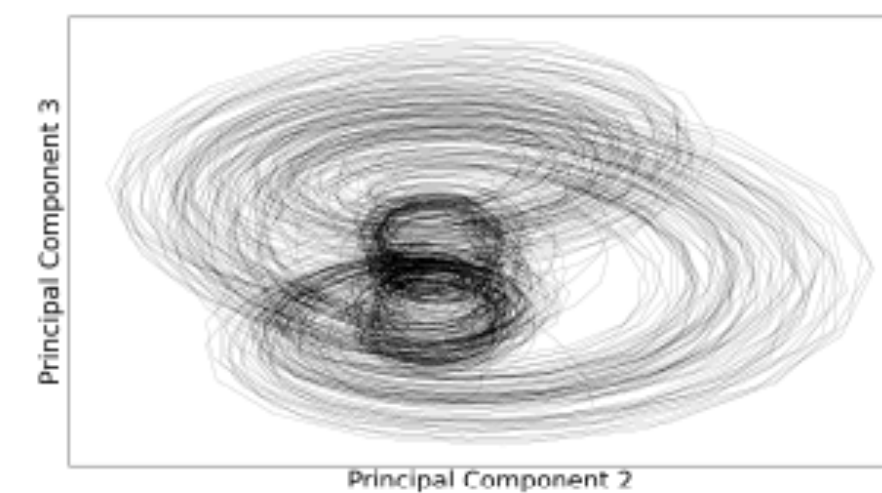
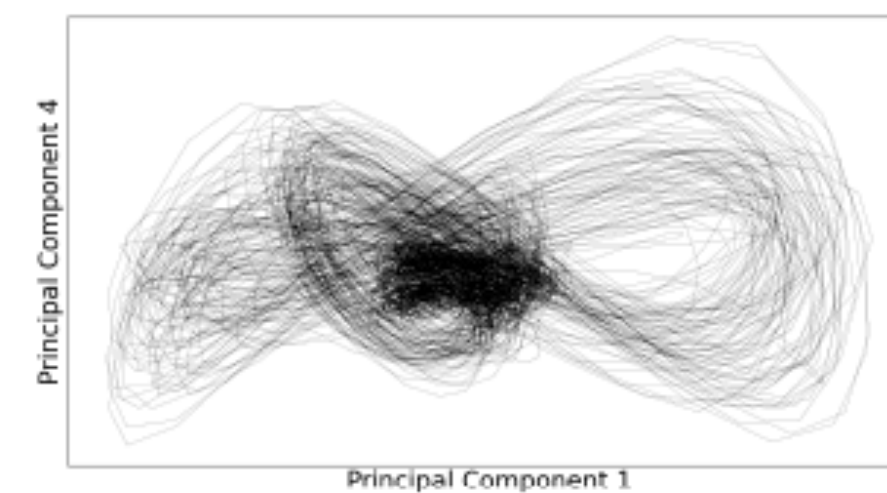
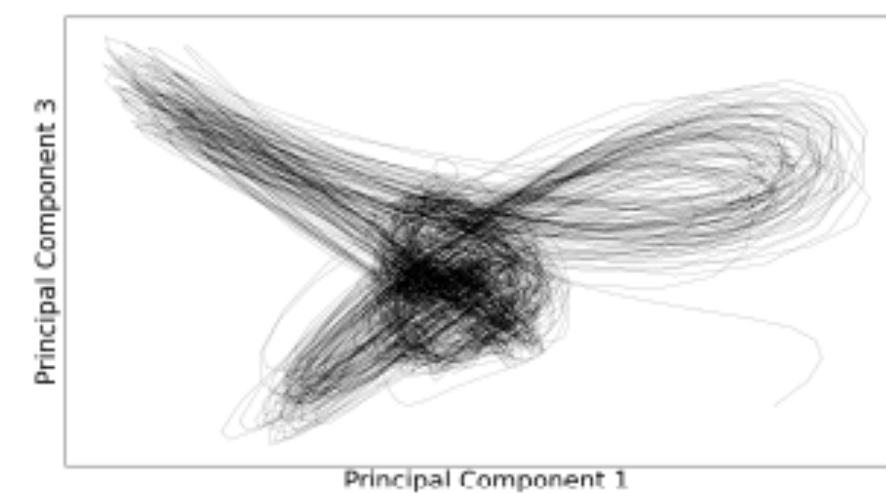
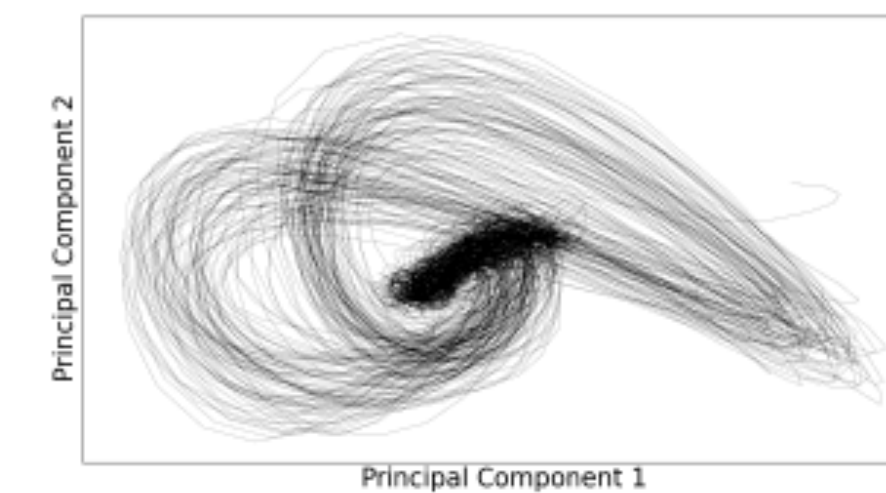
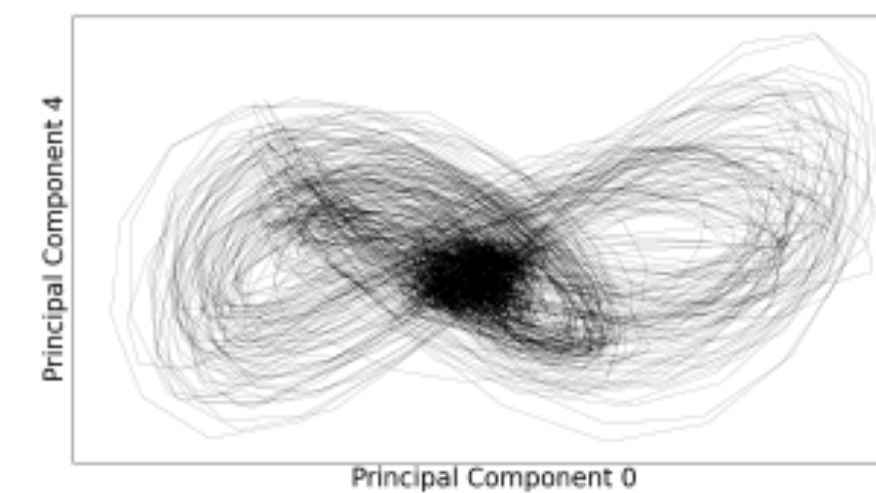
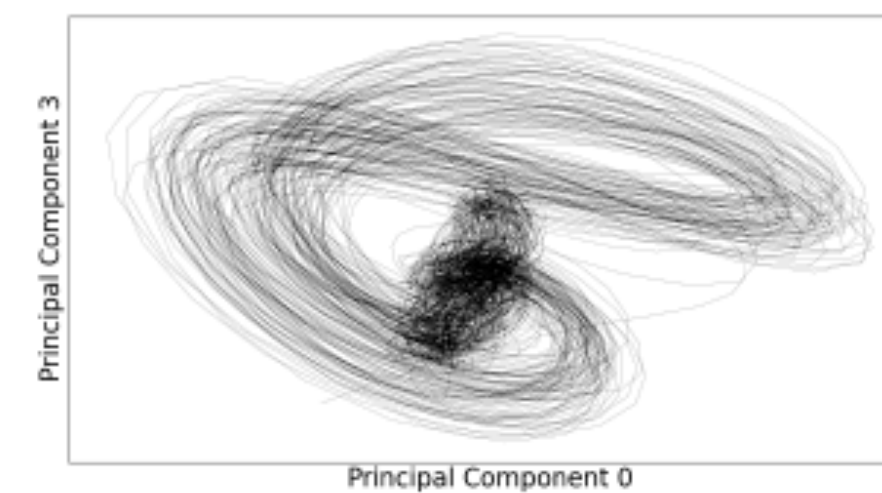
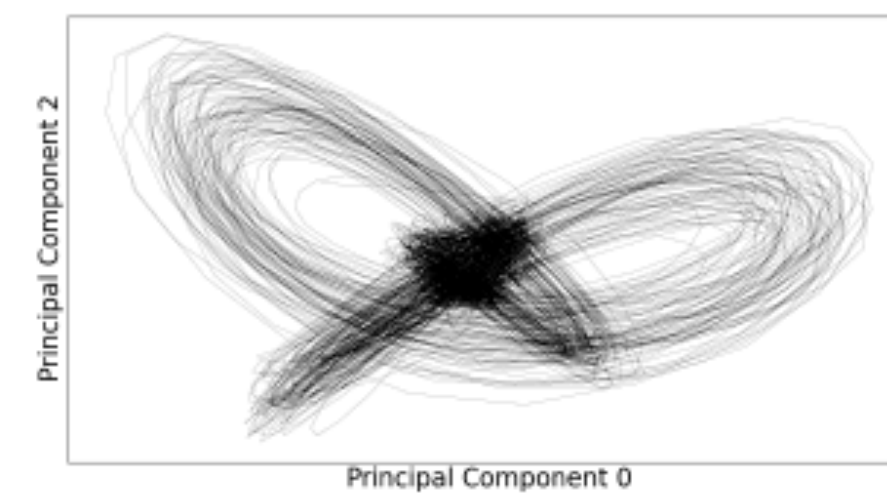
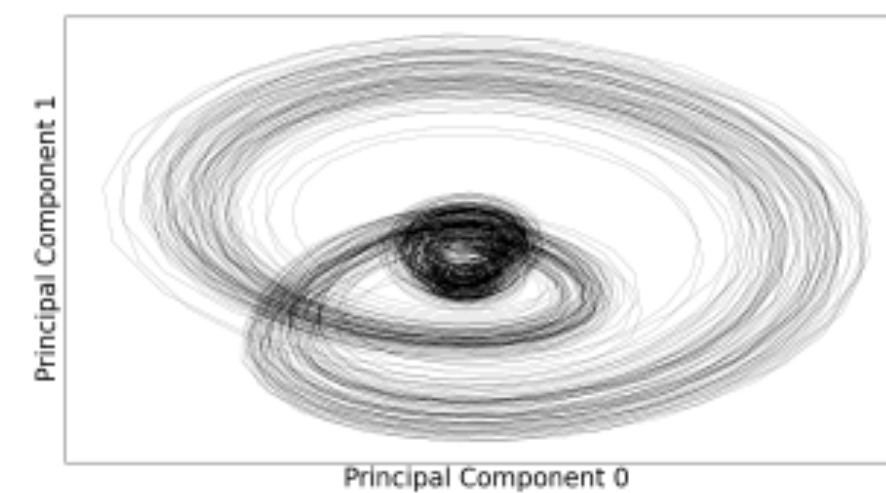
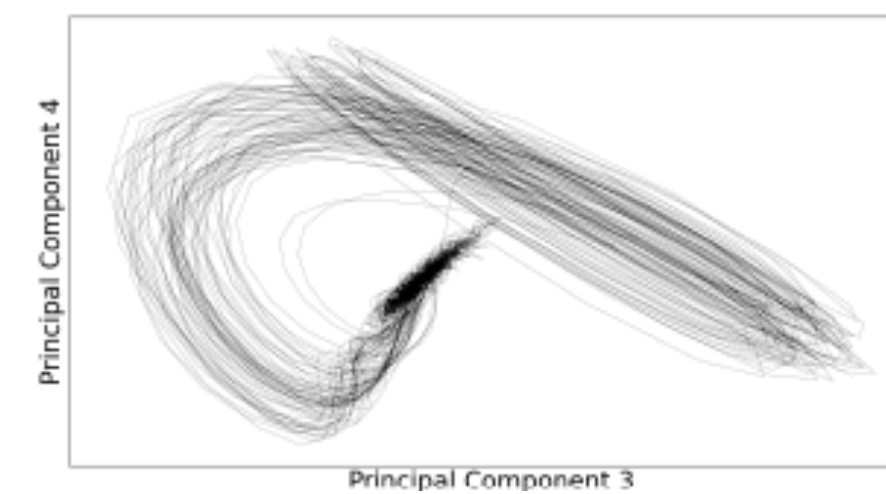
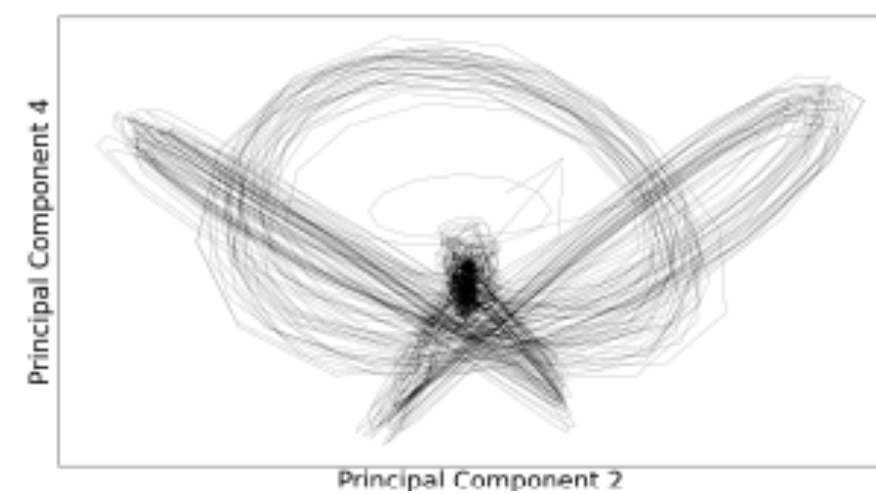
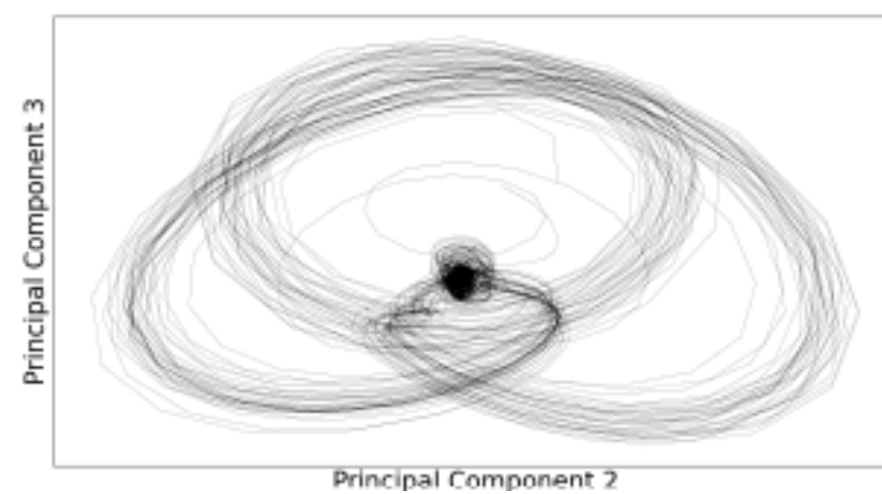
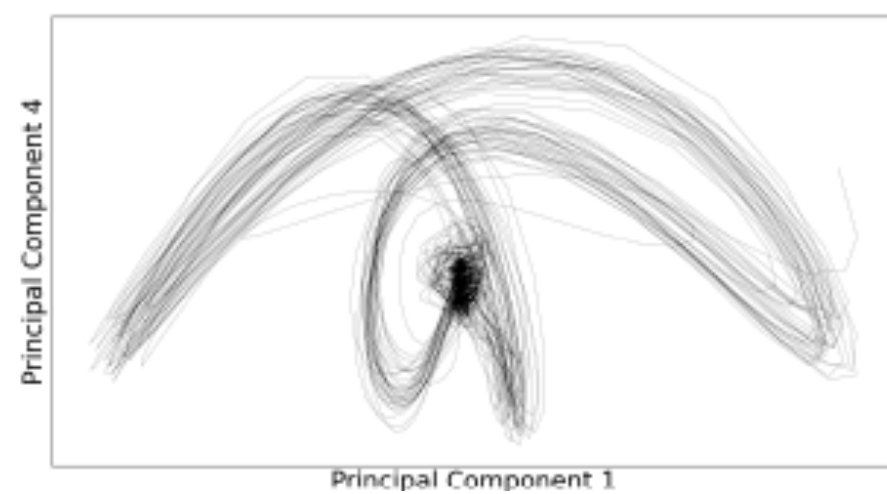
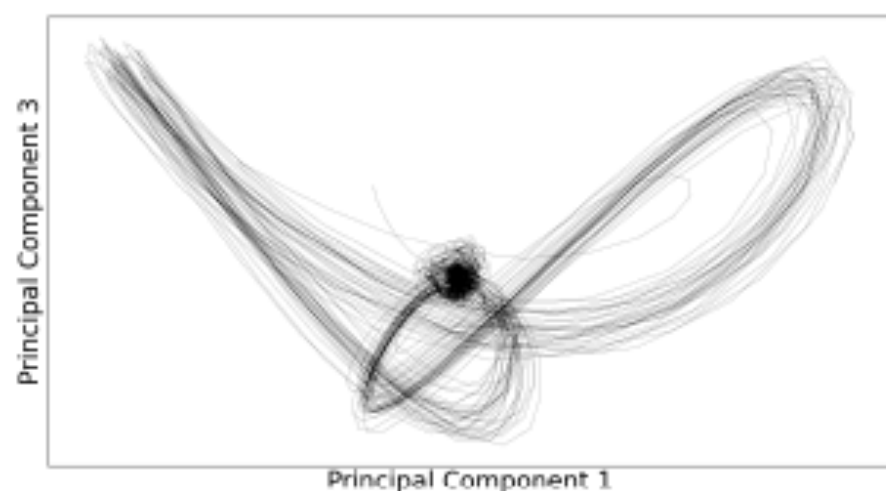
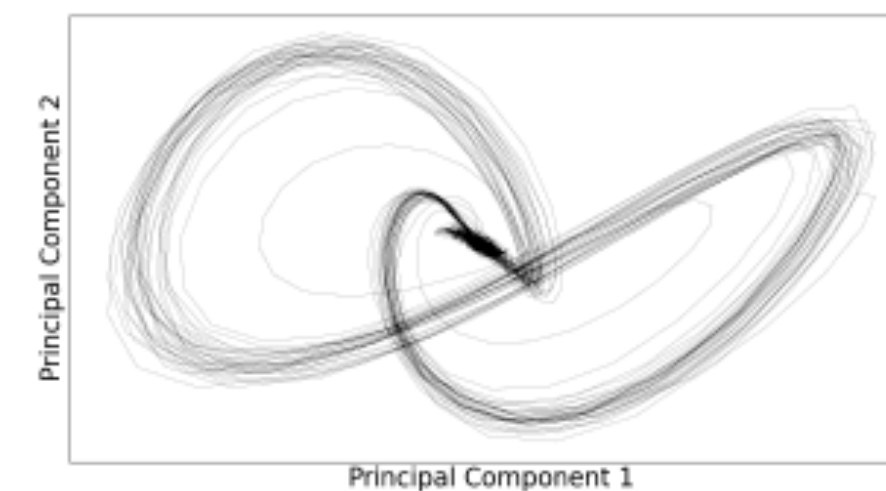
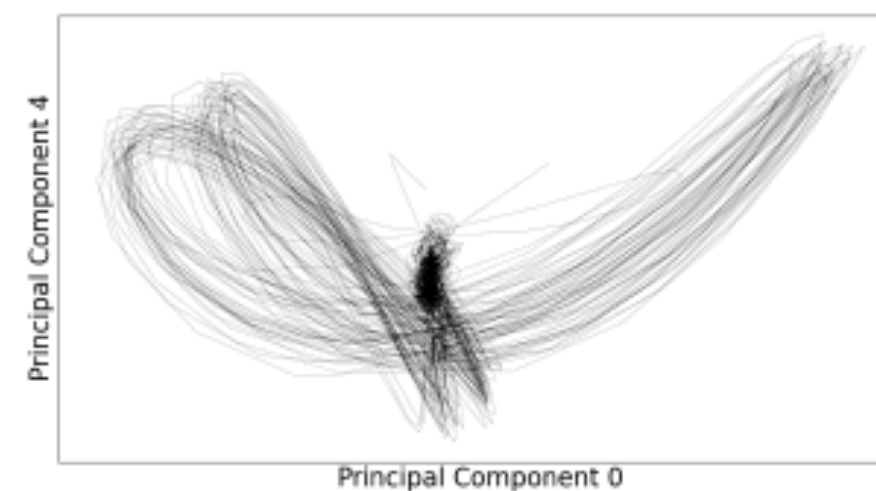
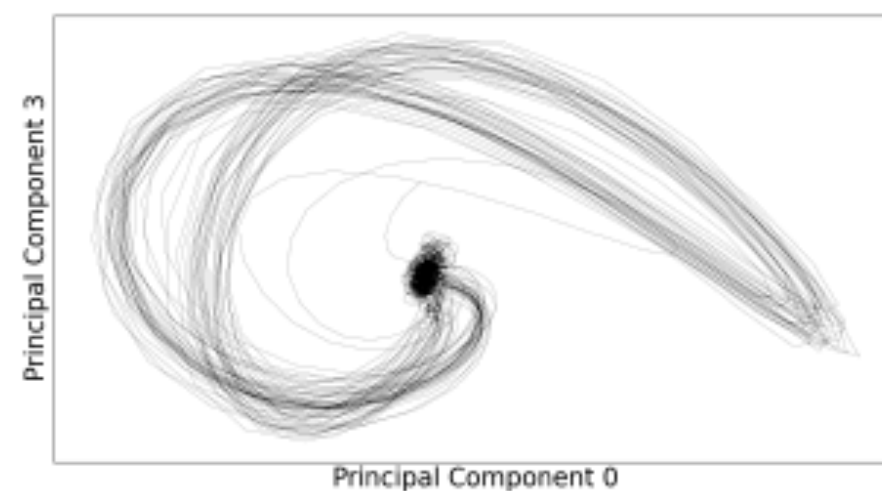
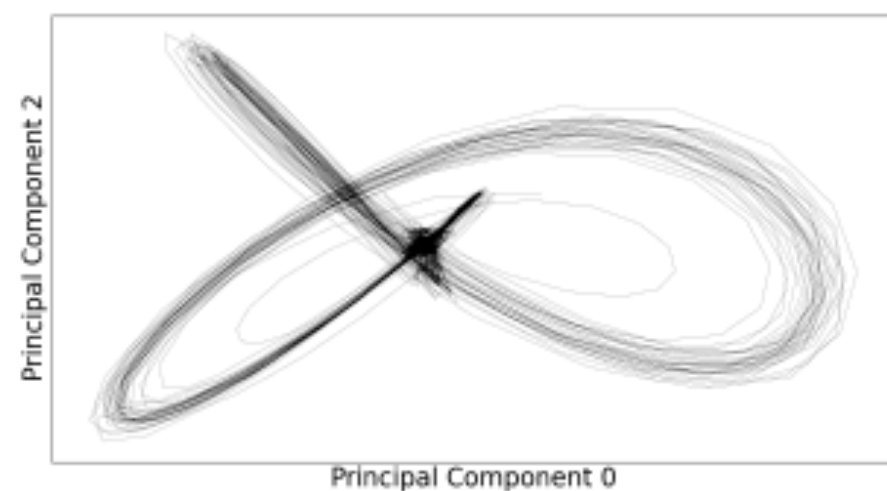
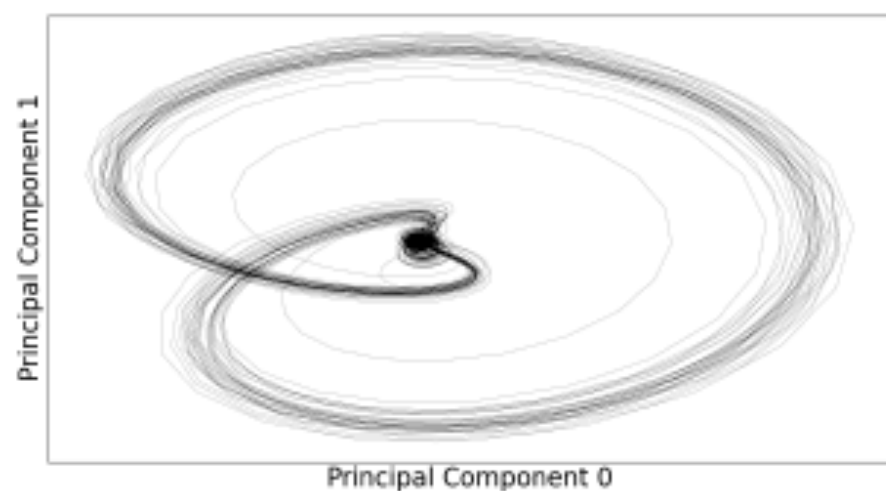
# Cardiac Arrhythmia Localization via ECG Attractor Imaging



|                         | 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        |
| <b>01 PCA Attractor</b> | <b>99.66</b> | <b>98.95</b> | <b>98.46</b> | <b>100</b> | <b>99.27</b> |

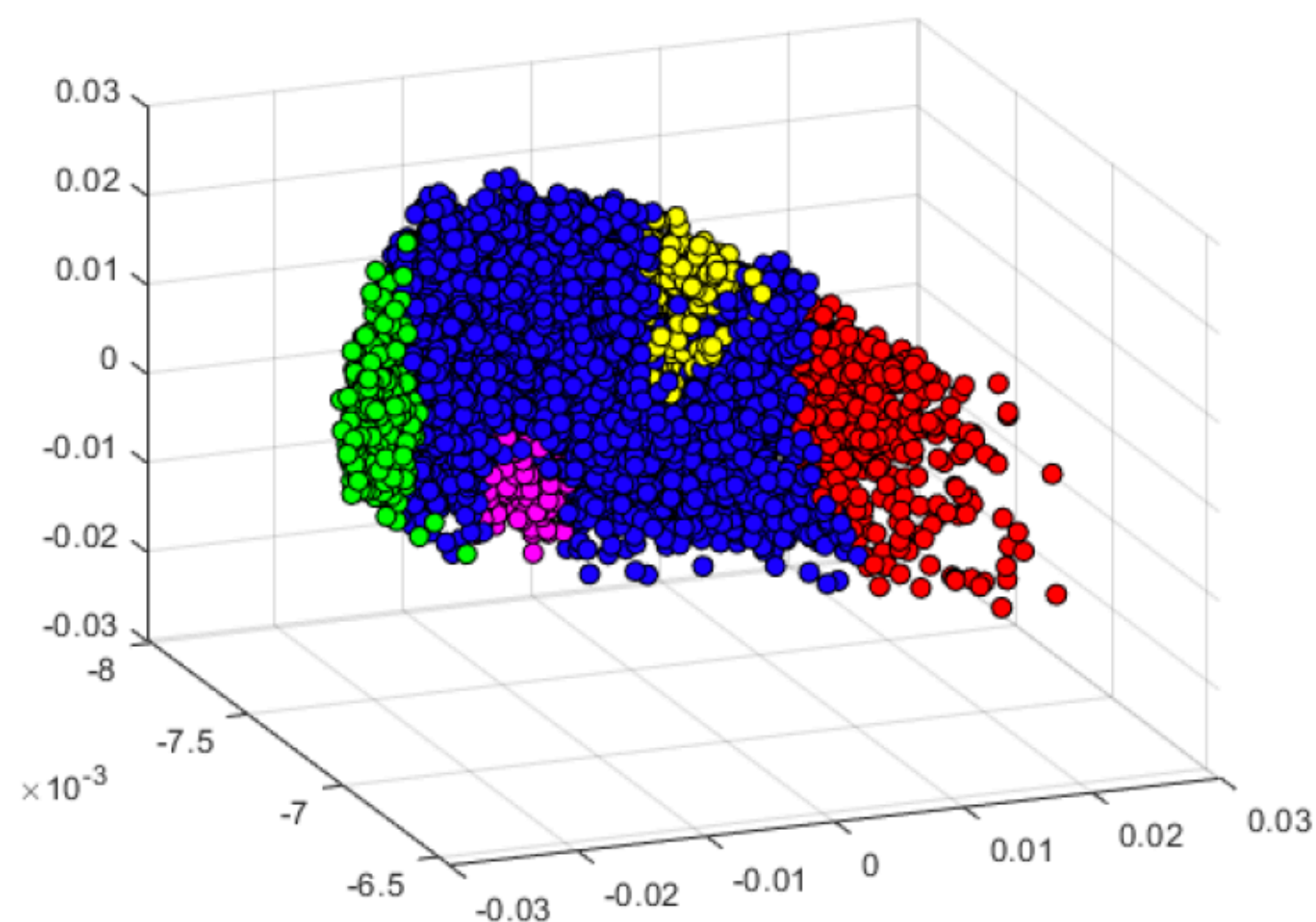
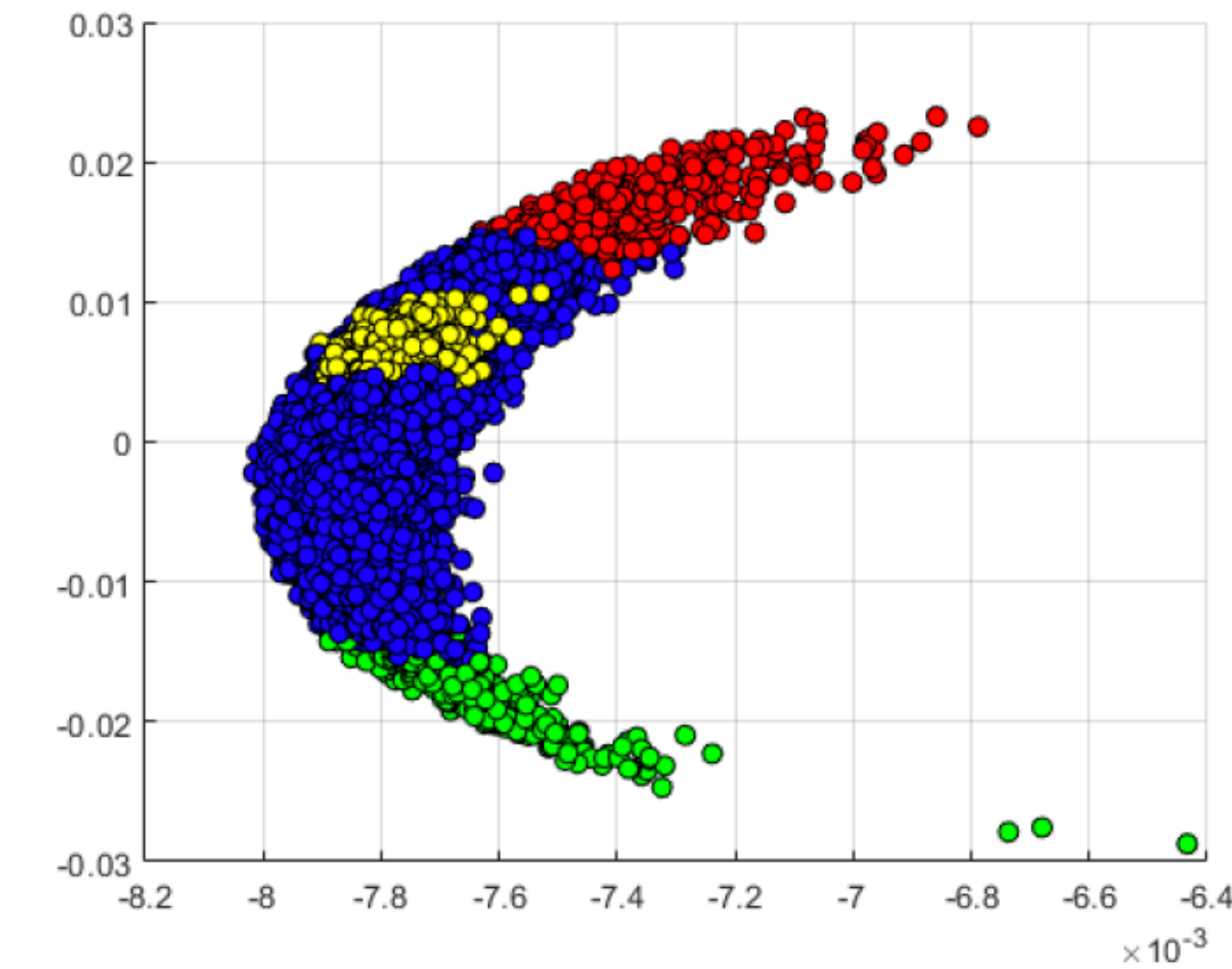


# Some Examples





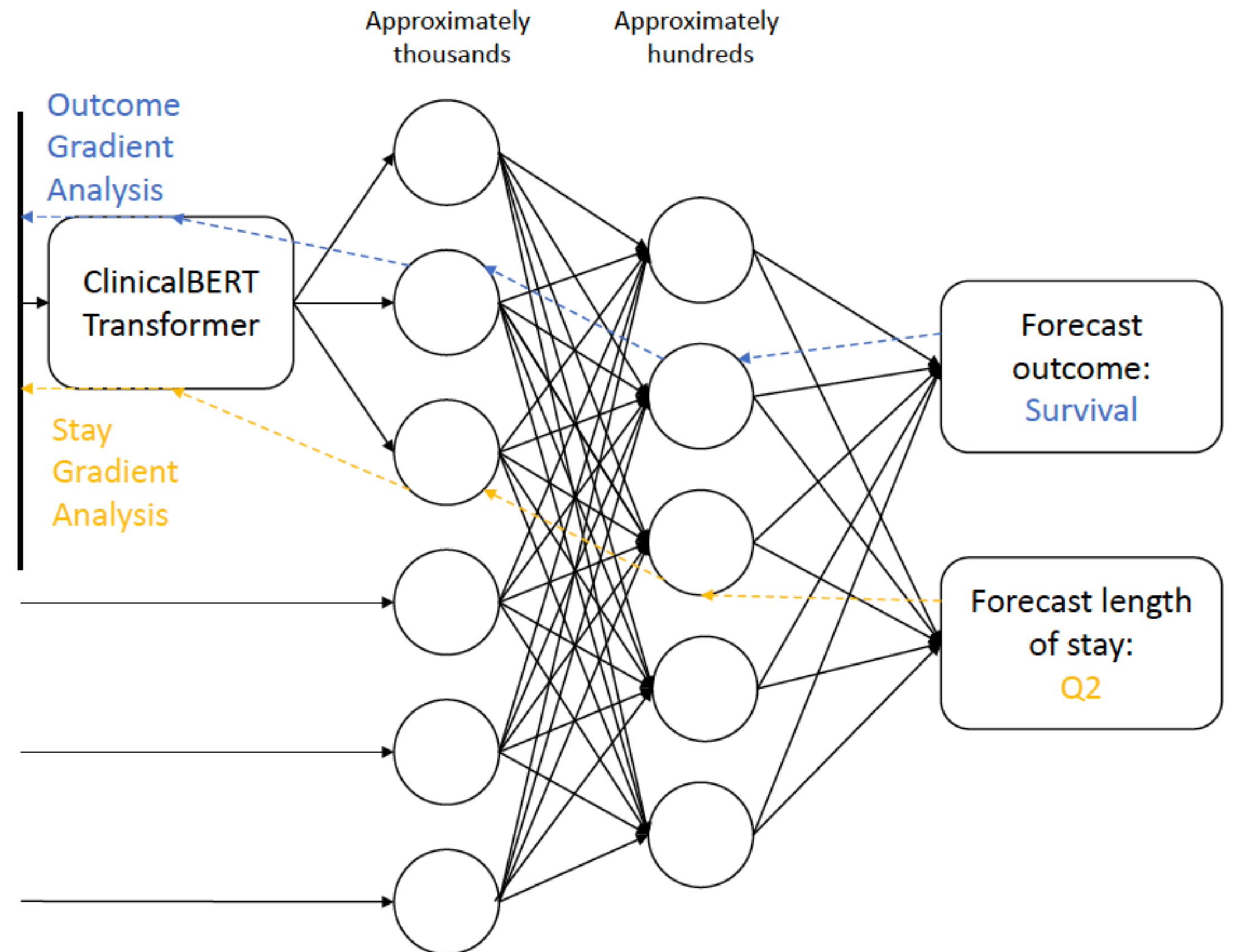
# Health Outcome Disparities on the EHR



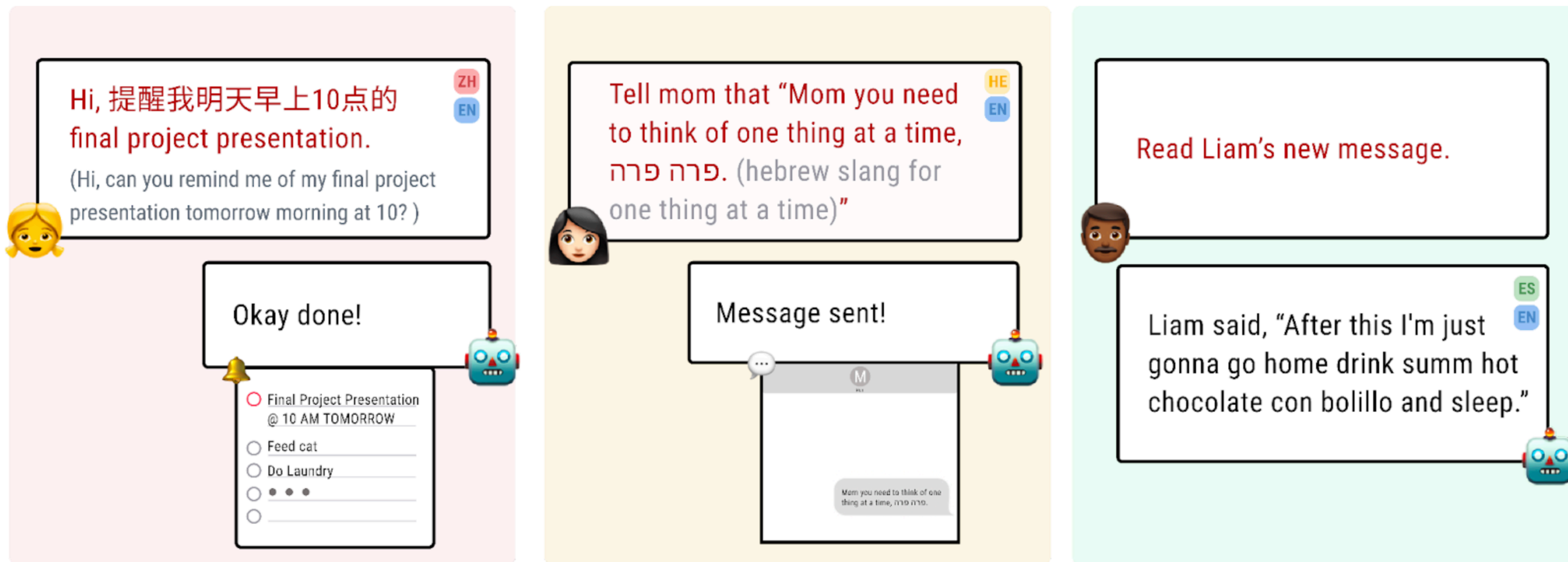
- Self Attention Contextual Embeddings

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

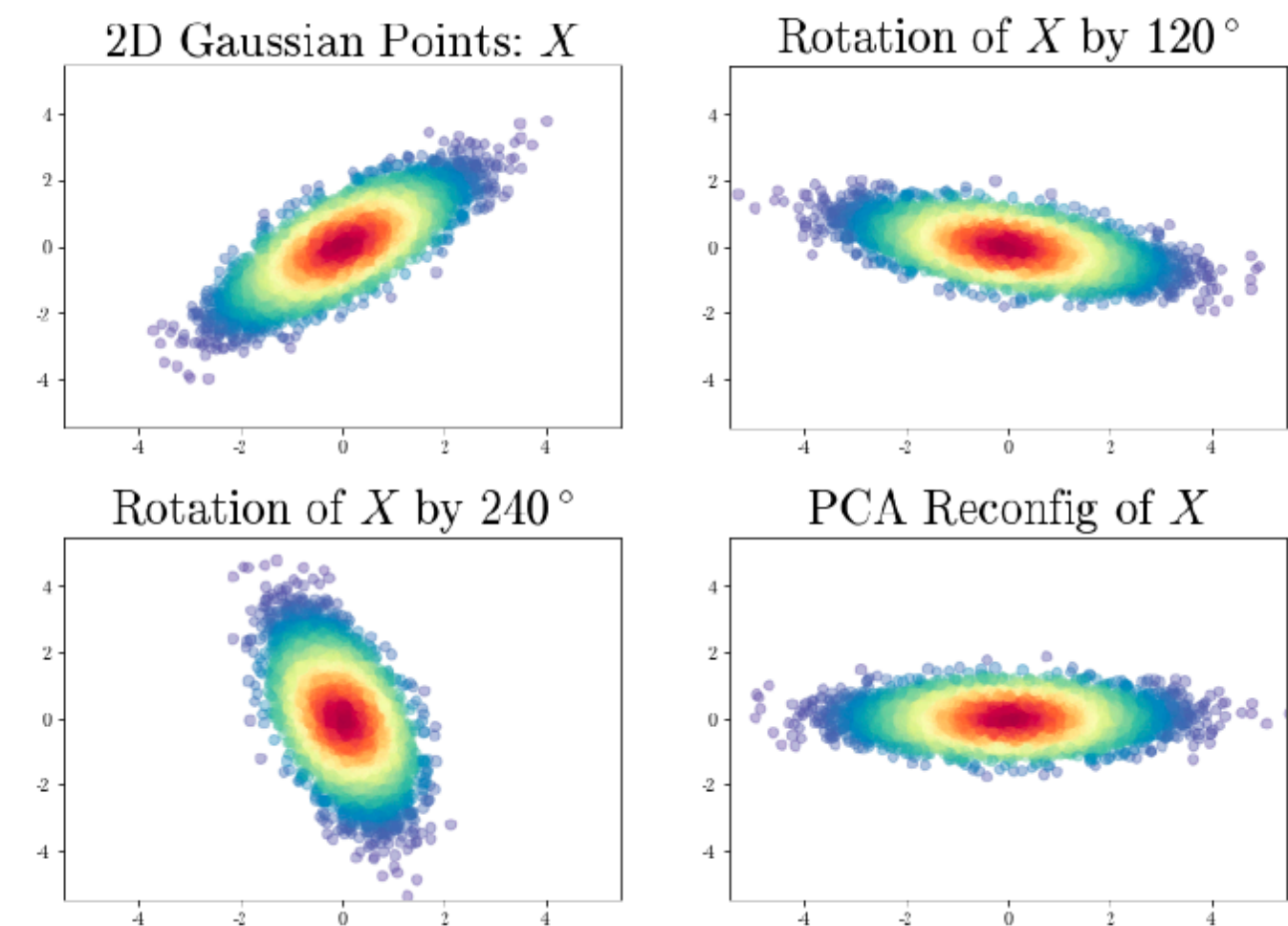
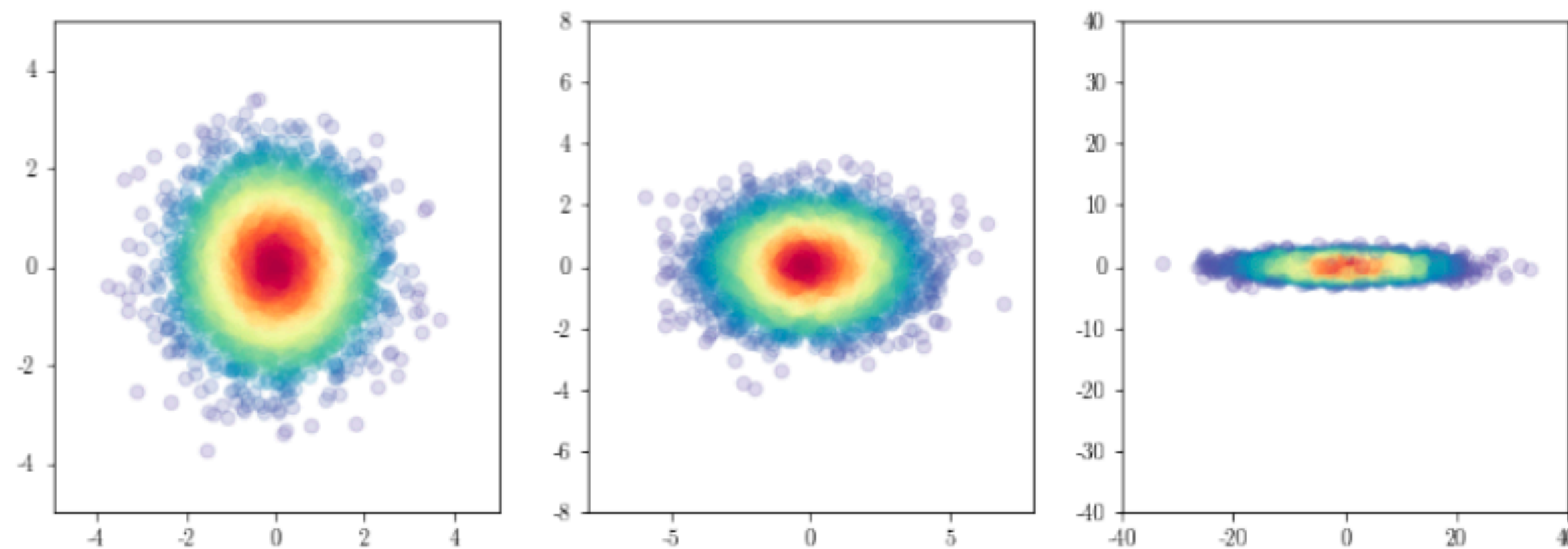


# Language Generation in a Code-switched World

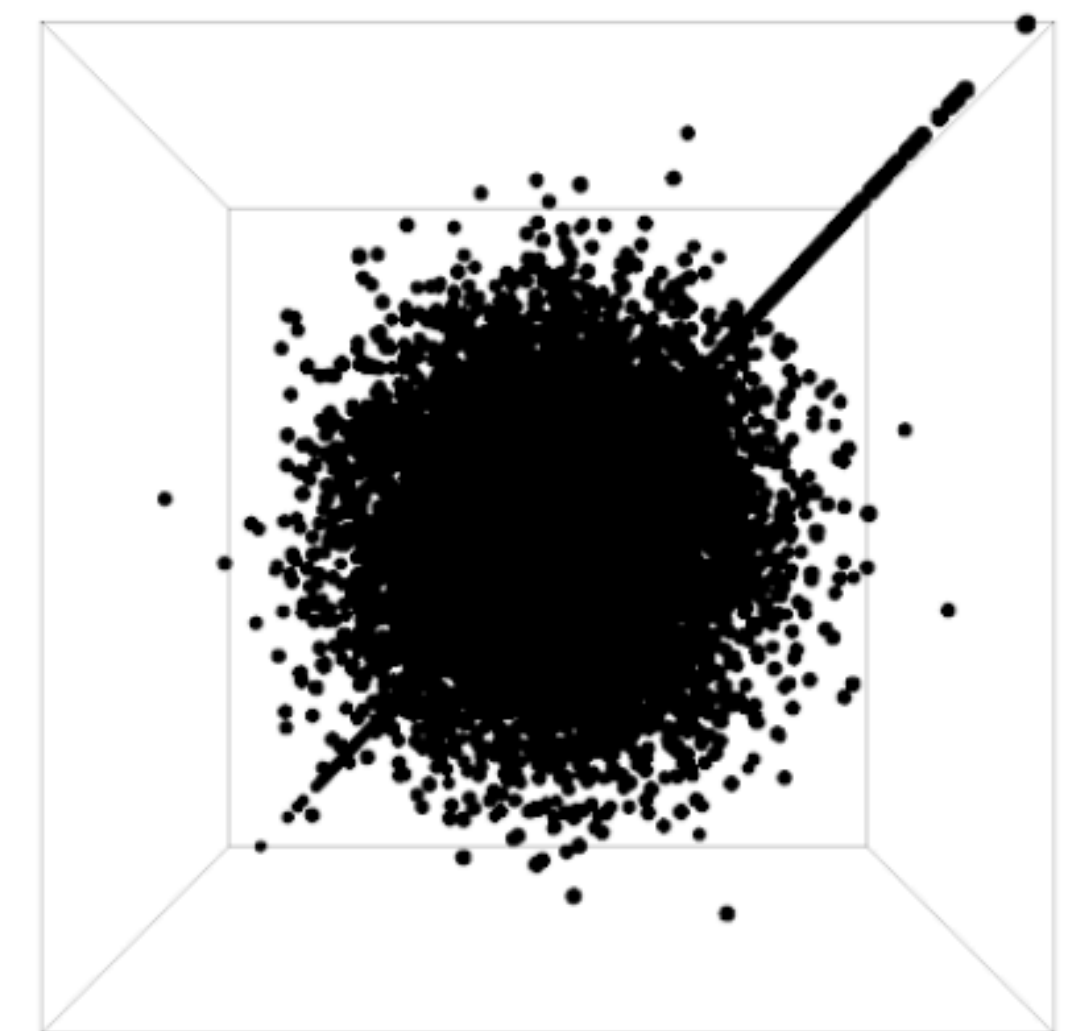




# Topological Embedding Analysis

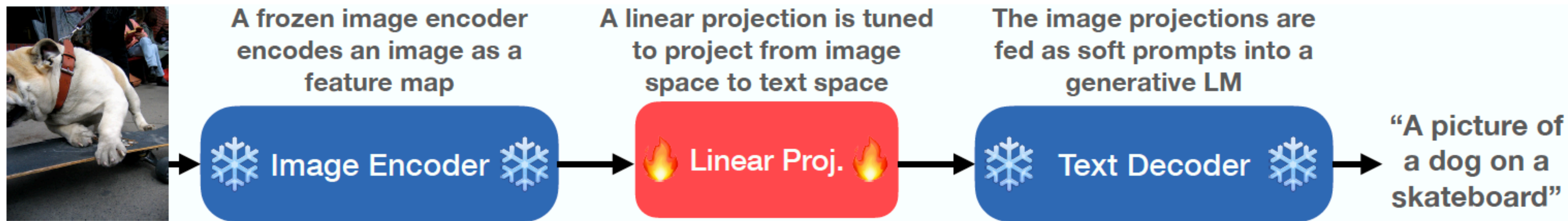


| Test                   | IsoScore | AvgRandCosSim | Partition | ID Score | VarEx |
|------------------------|----------|---------------|-----------|----------|-------|
| 1. Mean Agnostic       | ✓        | ✗             | ✗         | ✓        | ✓     |
| 2. Scalar Covariance   | ✓        | ✗             | ✗         | ✓        | ✓     |
| 3. Maximum Variance    | ✓        | ✗             | ✓         | ✗        | ✗     |
| 4. Rotation Invariance | ✓        | ✓             | ✗         | ✓        | ✓     |
| 5. Dimensions Used     | ✓        | ✗             | ✗         | ✗        | ✗     |
| 6. Global Stability    | ✓        | ✗             | ✓         | ✓        | ✗     |





# Are Language Models World Models?



## Image Captioning



|                      |   |
|----------------------|---|
| <b>CLIP</b>          | a giraffe in the lobby of the building  |
| <b>NFRN50</b>        | the giraffe in the zoo.                 |
| <b>BEIT</b>          | a peacock in the garden                 |
| <b>NFRN50 Random</b> | a man and a woman in a field of flowers |



|                      |   |
|----------------------|---|
| <b>CLIP</b>          | tennis player in action                 |
| <b>NFRN50</b>        | tennis player at the tennis tournament. |
| <b>BEIT</b>          | tennis player during a tennis match.    |
| <b>NFRN50 Random</b> | the new logo for the team               |

## Visual Question Answering



|                      |                               |
|----------------------|-------------------------------|
| <b>CLIP</b>          | He is surfing a wave.         |
| <b>NFRN50</b>        | He is surfing the waves.      |
| <b>BEIT</b>          | He is jumping into the water. |
| <b>NFRN50 Random</b> | He is swimming in the pool.   |

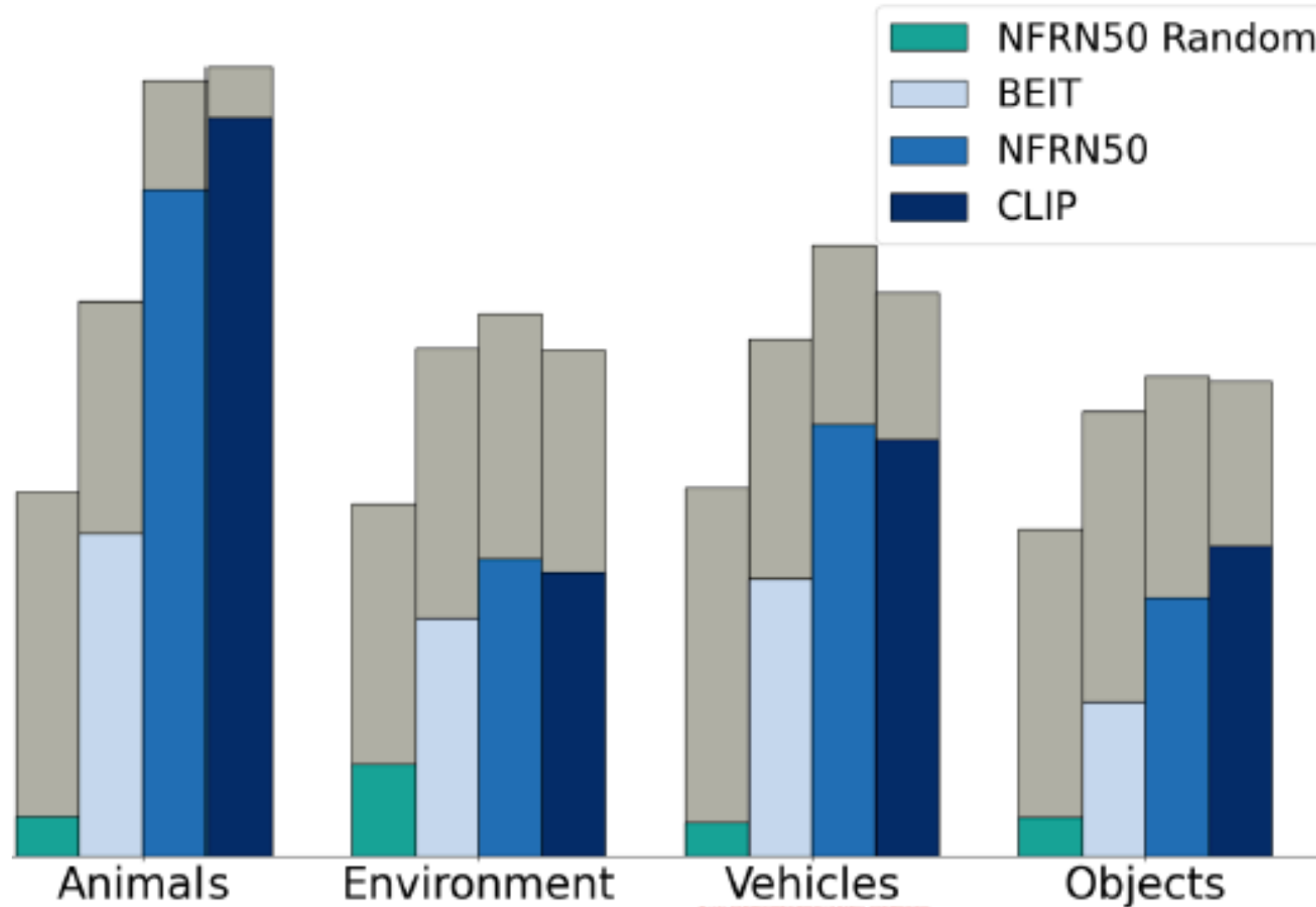
**Q: What is the person doing?**  
**A: surfing**



|                      |                 |
|----------------------|-----------------|
| <b>CLIP</b>          | A tennis racket |
| <b>NFRN50</b>        | A tennis racket |
| <b>BEIT</b>          | A baseball bat. |
| <b>NFRN50 Random</b> | A tree          |

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

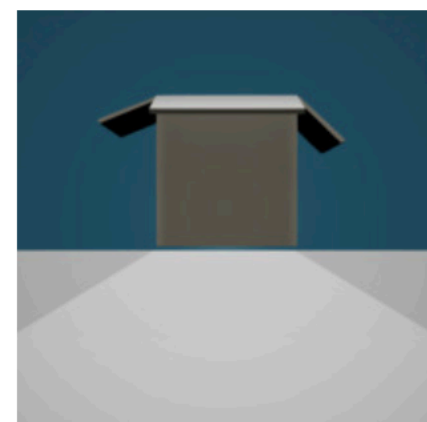
l-Palmer Similarity for each Model by Noun Category





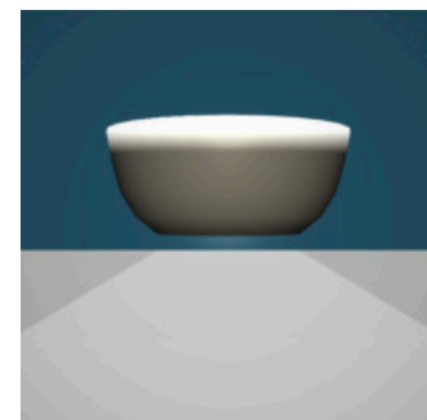
# Grounding Language Models in Physics

Play with  
these objects!



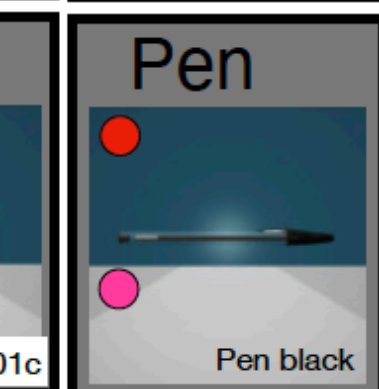
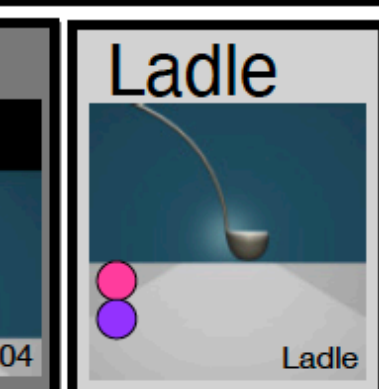
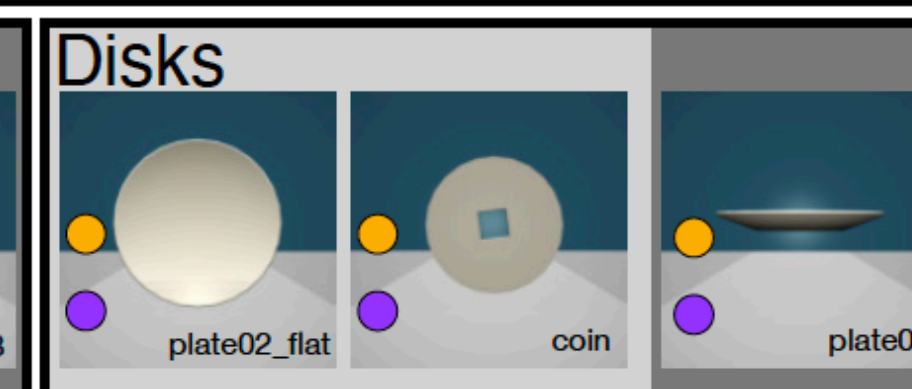
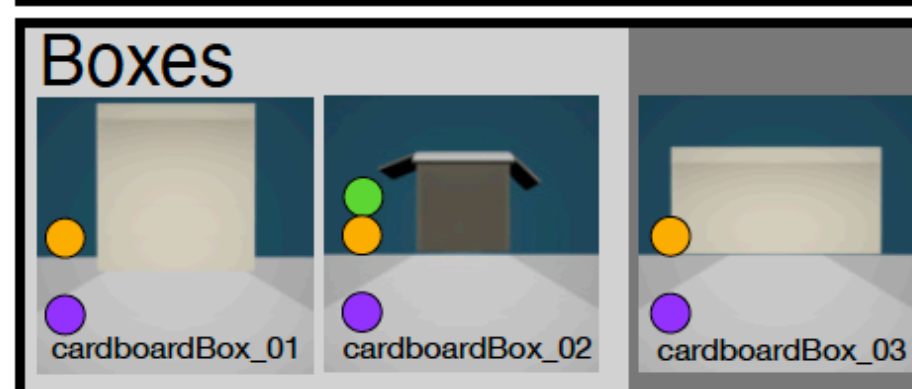
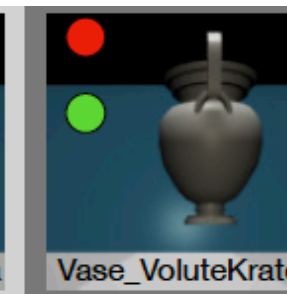
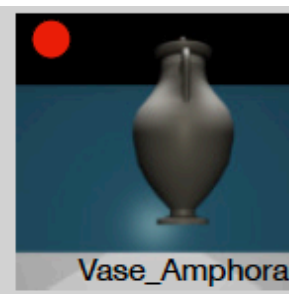
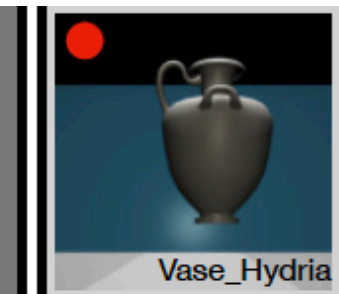
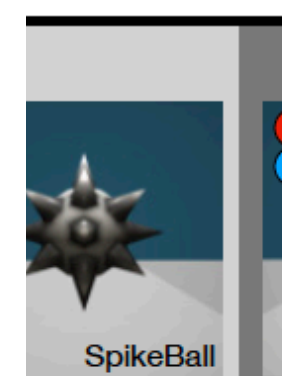
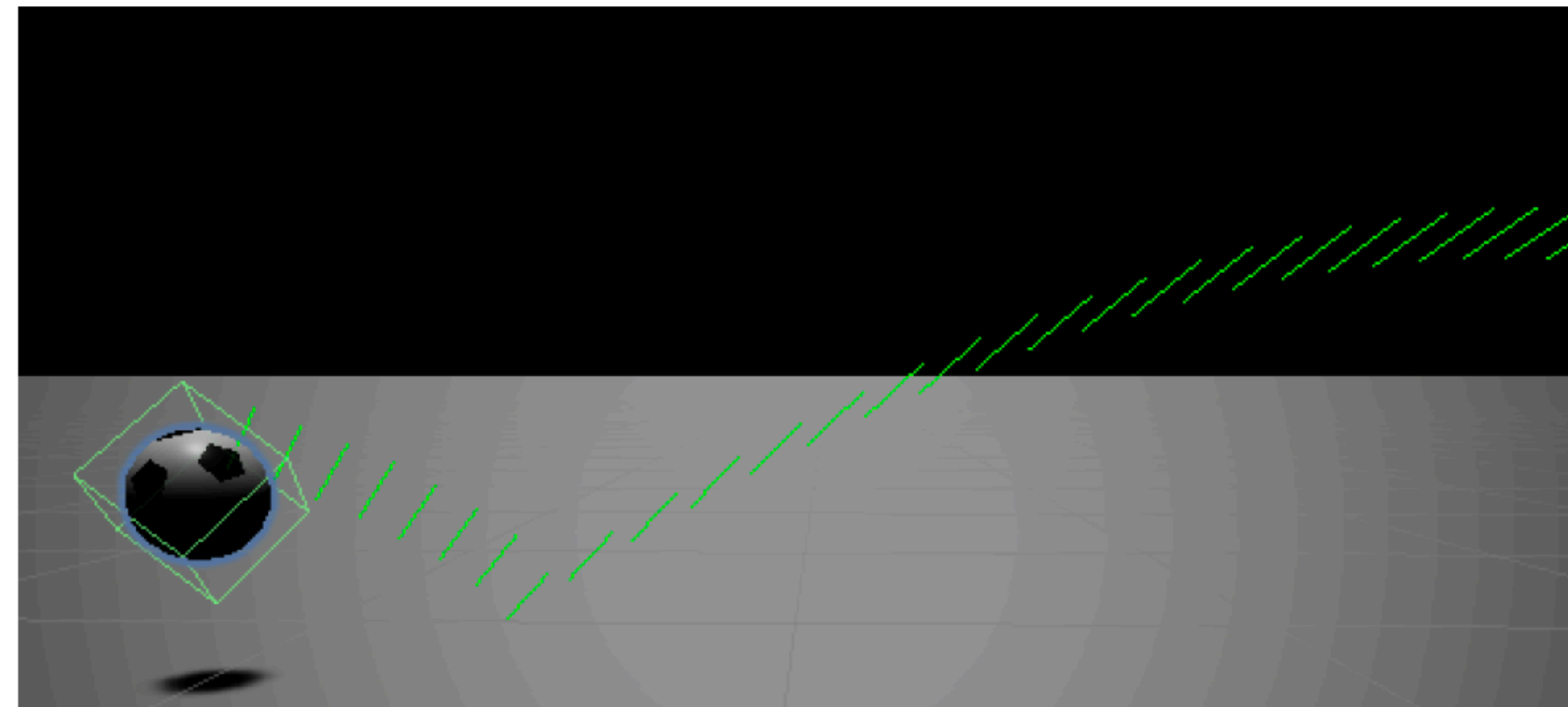
Seen

Now play  
with this one!



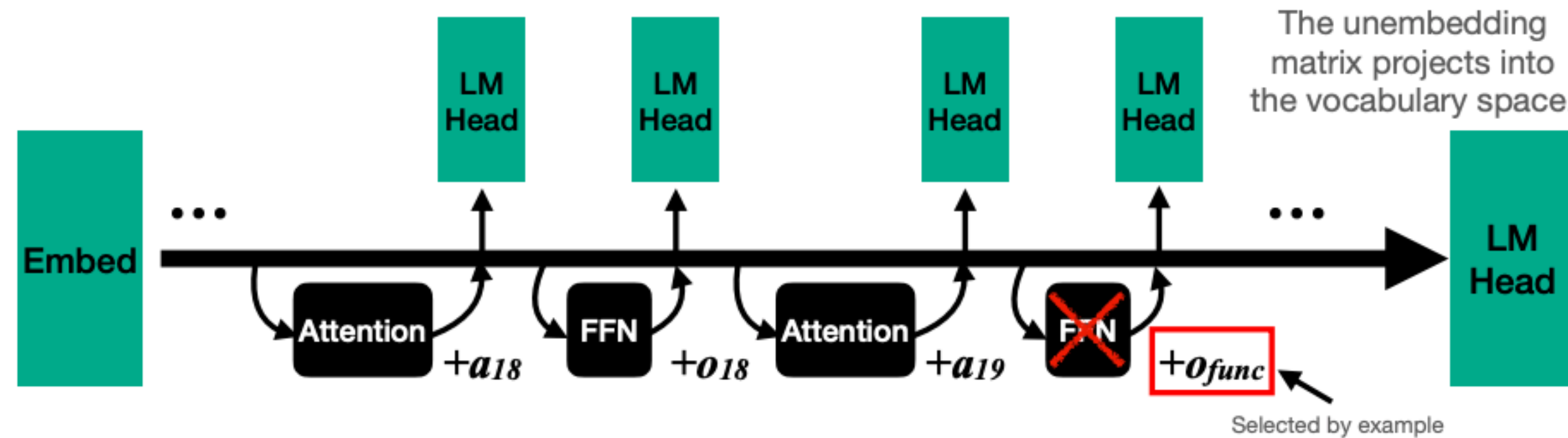
Can it roll?  
Can it contain  
other objects?

Unseen



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

# LLM Vector Arithmetics



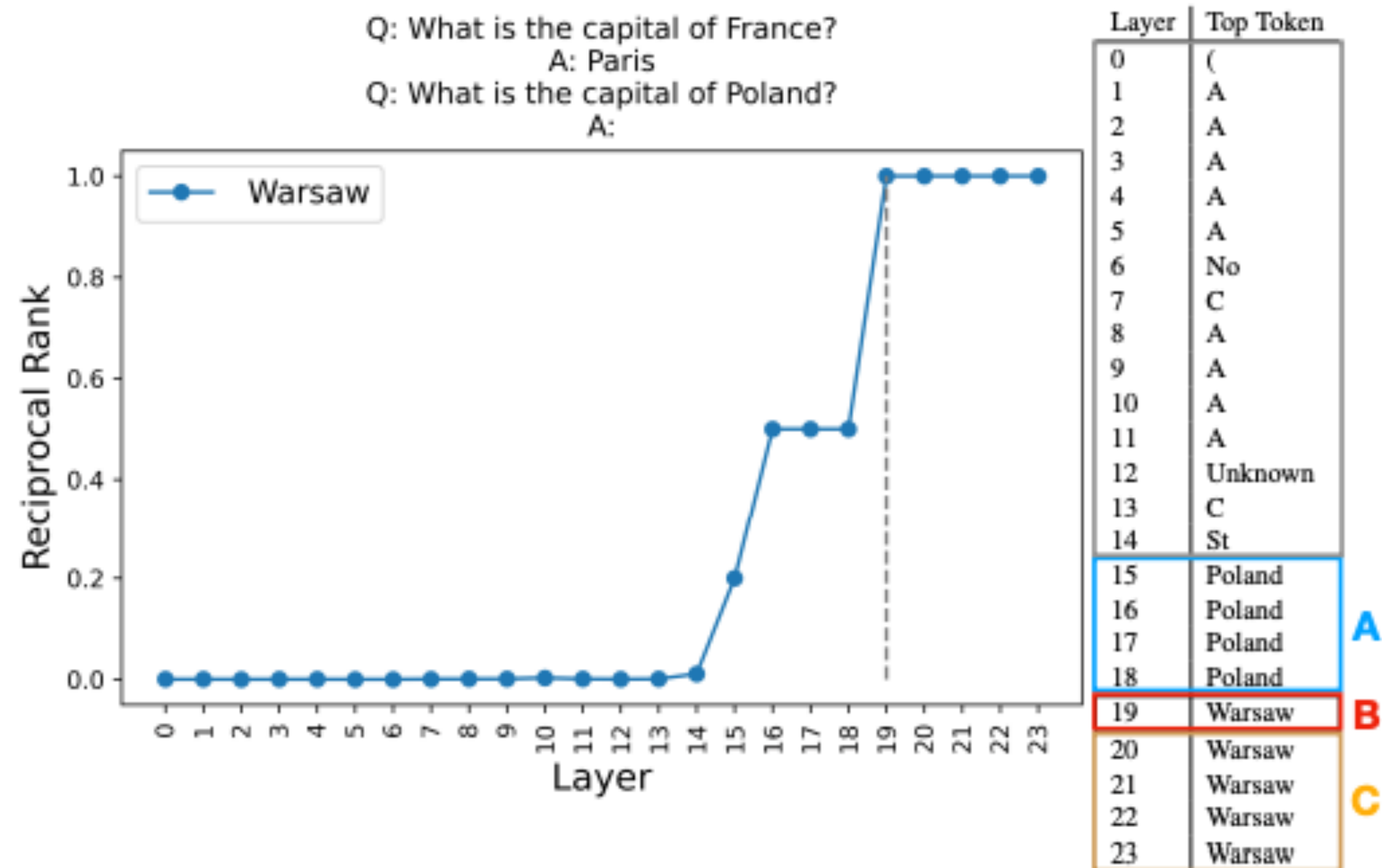
- *We find simple vector arithmetic in Transformers*

- *3 Models*

- *GPT2 (124M)*

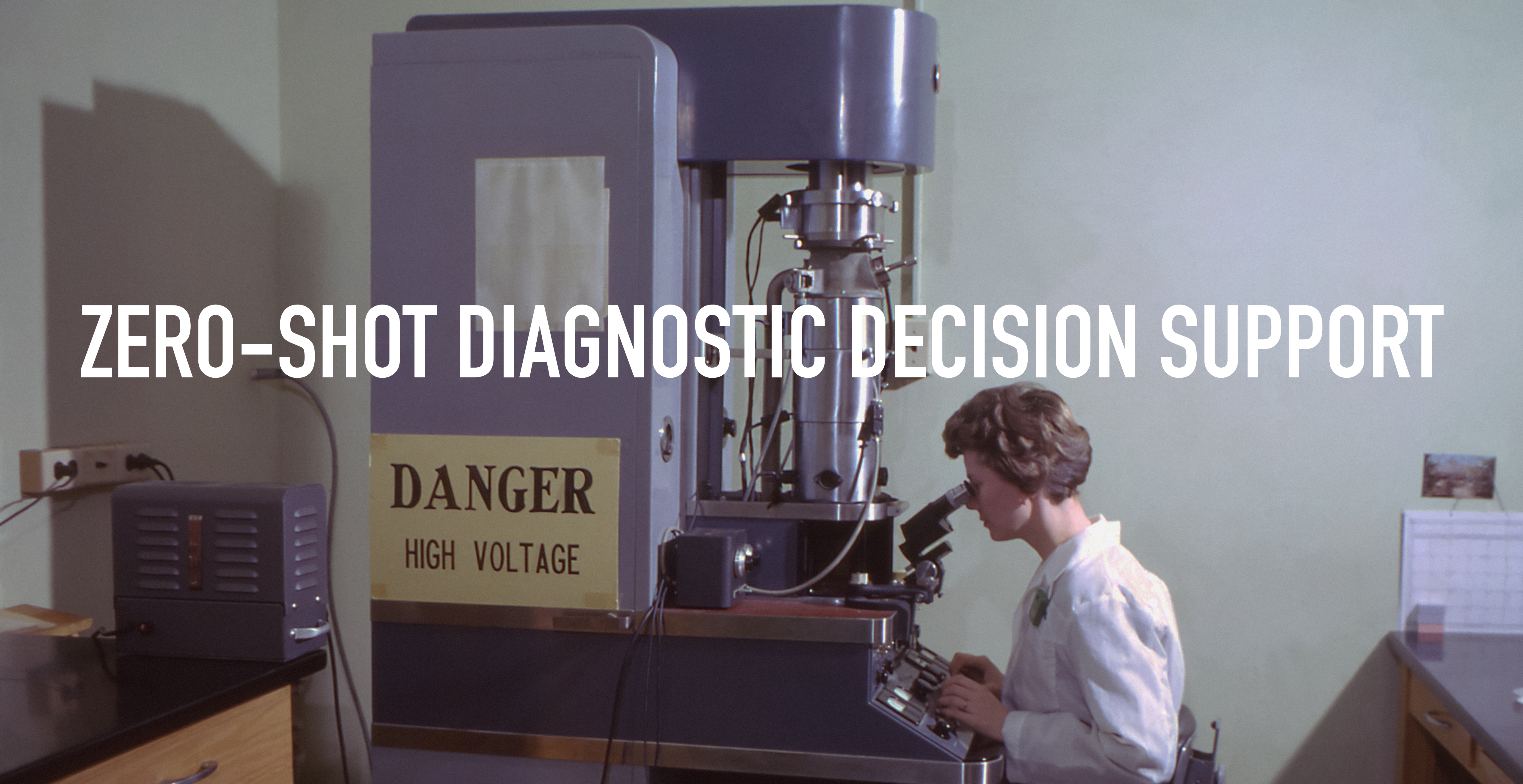
- *GPT-J (6B)*

- *Bloom (176B)*





# ZERO-SHOT DIAGNOSTIC DECISION SUPPORT





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%)
- ...



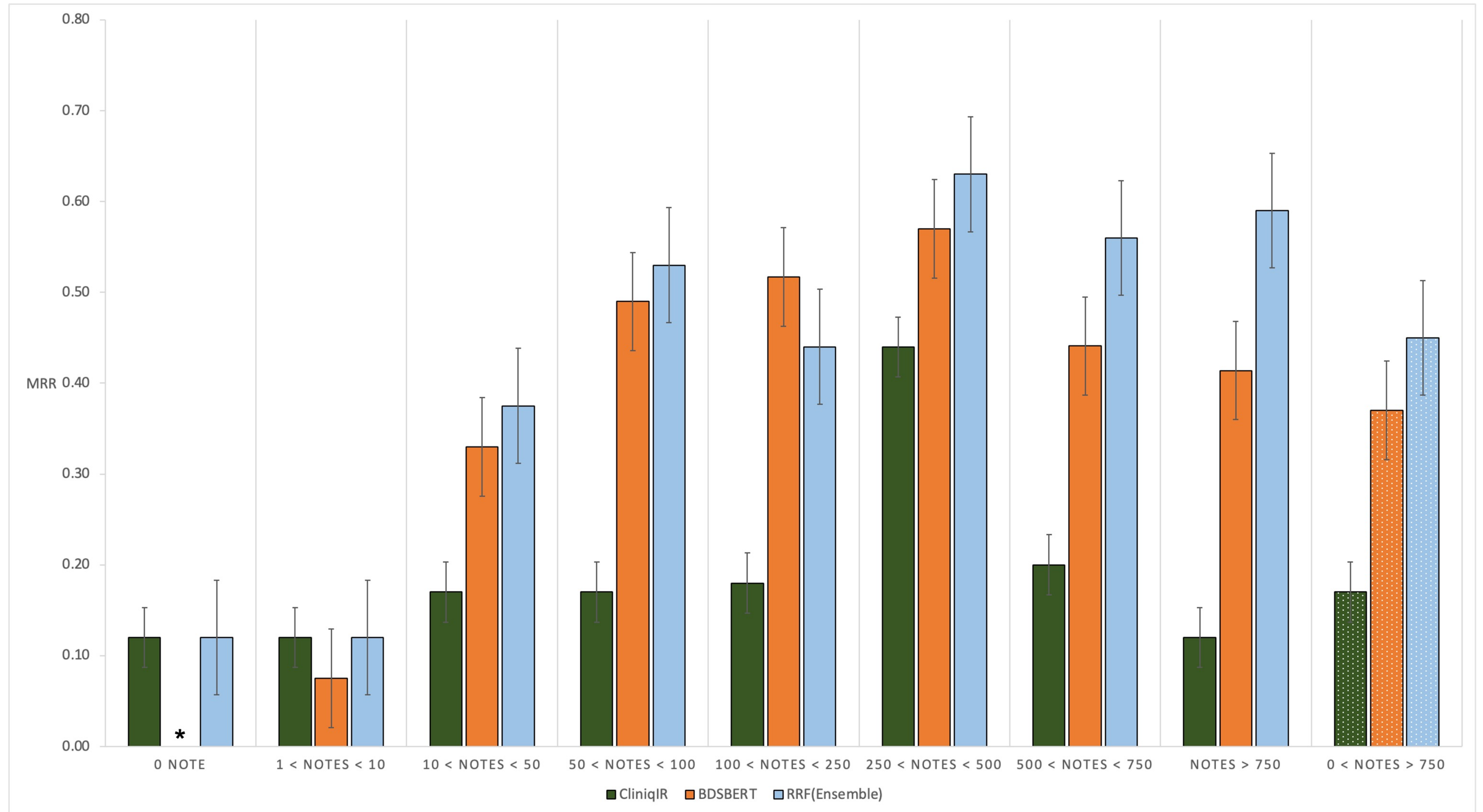
# ML to the Rescue

- *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]
- ...



# Results II

B.



# Clinical LLMs

## Week (05/15/23)- Can Large Language Models Perform Complex Diagnoses?

- We tested these models on the DC3 dataset
- It contains 30 difficult to diagnose case challenges

| Models  | Task I<br>(Open-ended) | Task II<br>(Multiple-choice) | Task III<br>(Multiple-choice) | Total     |
|---------|------------------------|------------------------------|-------------------------------|-----------|
| ChatGPT | 8                      | 13                           | 11                            | 16        |
| GPT4    | 8                      | <b>14</b>                    | <b>13</b>                     | 17        |
| BARD    | <b>10</b>              | <b>14</b>                    | <b>13</b>                     | <b>18</b> |





HEALTH-NLP.ORG



# Why is ChatGPT so good?

- Attention
- Lots of (diverse) training data
- Masked Language Modeling
- Tight Guardrails

