



**WPI**

# Towards Fairer Health Recommendations:

## Finding informative unbiased samples via Word Sense Disambiguation

Gavin Butts<sup>1\*</sup>, Pegah Emdad<sup>2\*</sup>, Jethro Lee<sup>3\*</sup>

Shannon Song<sup>2</sup>, Chiman Salavati<sup>4</sup>, Wilmar Sosa Diaz<sup>4</sup>, Roberto  
Montenegro<sup>5</sup>, Scott Hale<sup>6</sup>, Shiri Dori-Hacohen<sup>4</sup>, Fabricio Murai<sup>2</sup>

\* Equal Contribution

<sup>1</sup>Loyola Marymount University, <sup>2</sup>Worcester Polytechnic Institute, <sup>3</sup>Northeastern University,

<sup>4</sup>University of Connecticut, <sup>5</sup>University of Washington School of Medicine, <sup>6</sup>Meedan & Oxford Institute

# Table of contents

**01**

## **Introduction**

Fairness and Health  
Recommender Systems

**02**

## **Related Work**

Manual and Computational  
Bias Detection

**03**

## **Dataset**

BRICC Dataset for Bias  
Reduction

**04**

## **WSD**

Evaluation of ML for  
Word Sense  
Disambiguation

**05**

## **Bias Classification**

Evaluation of Bias  
Detection Models

**06**

## **Conclusion**

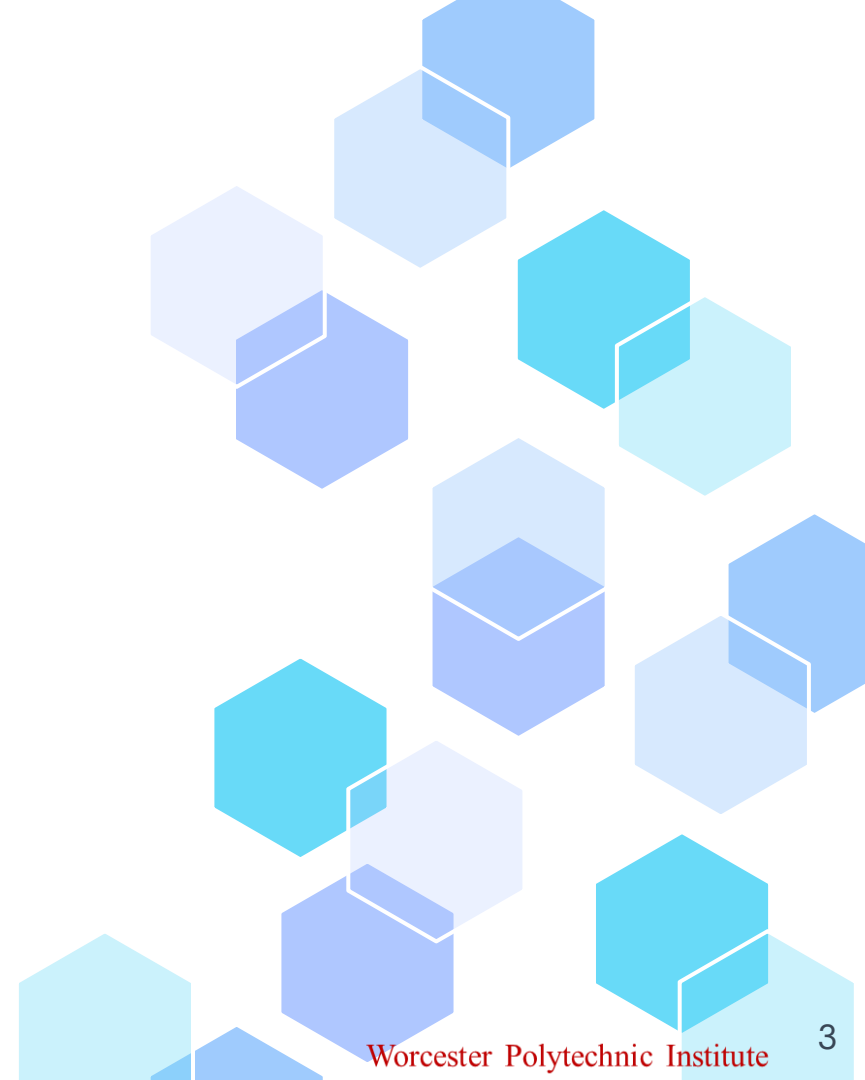
Key Findings and Future  
Implications

---

# 01

# Introduction

Fairness and Health Recommender Systems



# Introduction



Health Recommender  
Systems (HRS)



Impacts of Biased Data

# Introduction

Personalized **health recommendations** deploying  
**machine learning** and **information retrieval**

Dependence of **Recommender systems' reliability** on  
the **quality** of their **training data**

**Negative impact** of **biased predictions** on **patient care**  
widening **health disparities**

# Disambiguous Sample

**Table 1: The term “white” in a racial vs. non-racial context**

Race-Related	Not Race-Related
“5 Year Relative Survival: overall 84% for <b>white</b> women, 62% for black women, 95% for local disease, 69% regional disease (spread to lymph node), 17% for distant disease.”	“ <b>White</b> matter within the spinal cord contains the axons of neurons that are ascending and descending to transmit signals to and from the brain, respectively.”

---

# 02

## Related Work

Manual and Computational Bias Detection

# Related Work



Manual Bias  
Detection



Computational Bias  
Detection



LLMs and Prompt  
Engineering



# Related Work



## Manual Bias Detection

Growing concerns over **biased AI models** in **healthcare recommender systems** due to their use in **high-stakes decisions**

### Our Approach:

- Exploring AI models for debiasing medical text.
- Augmenting unbiased samples and evaluating a wider range of models, including LLMs
- Data refinement using WSD

# Related Work



Computational  
Bias Detection

## Our Approach:

- Apply **Large Language Models (LLMs)** for bias detection.
- Use **TinyLlama**, an efficient version of **Llama 2**, for bias classification.
- Implement **Word Sense Disambiguation (WSD)** to improve **data refinement** and enhance the set of **negative samples**.

# Related Work



## LLMs and Prompt Engineering

- LLMs perform **on par** with encoder-only models like **BERT** in NLP tasks **without fine-tuning**.

### Prompting Techniques:

- **Zero-shot**, **Few-shot**, and **Chain of Thought (CoT)** prompting are crucial for improving **model quality** and **output accuracy**.

### Our Approach:

- We are the first to evaluate **zero-** and **few-shot prompting** for detecting **bias in medical curricular content**.

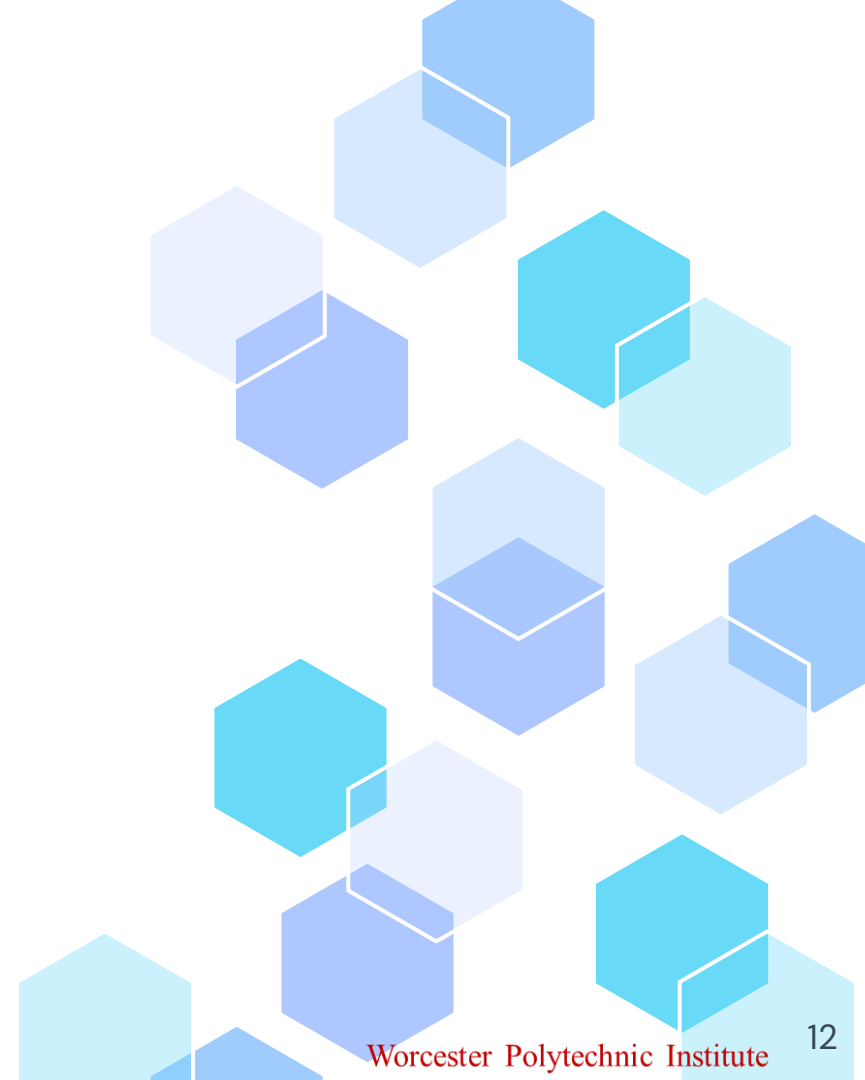
---

# 03

# Dataset

BRICC\* Dataset for Bias Reduction in  
Curricular Content

\*Salavati et al. (2024)



# Data Labels

## First-level: Identify Social Demographic

*'Sex,' 'Gender,' 'Race,' 'Ethnicity,' 'Age,' and 'Geography'*

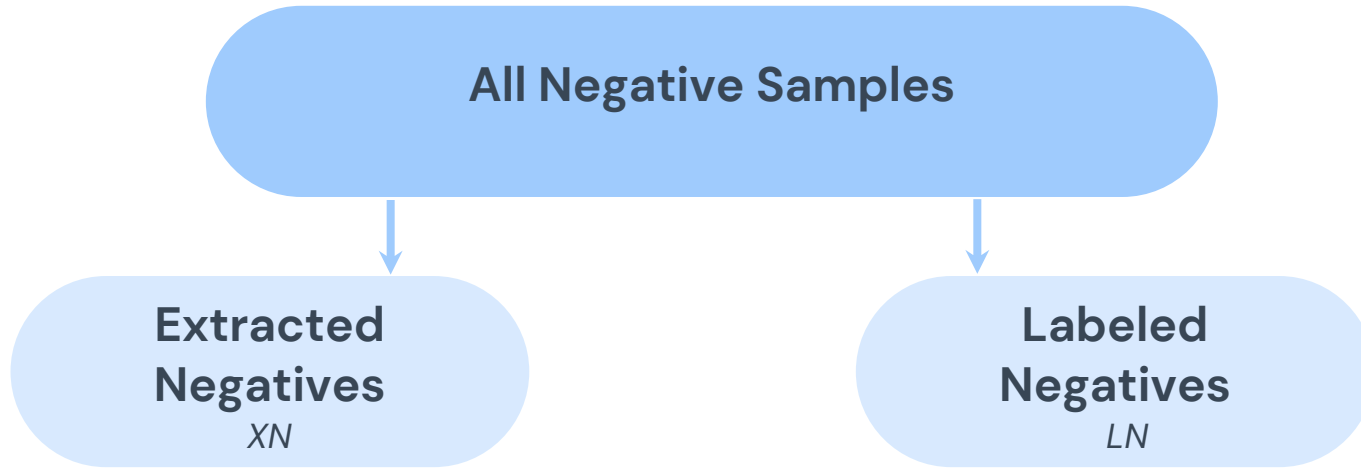
## Second-level: Identify Bias

*'Biased,' 'Potentially Biased,' 'Non-Biased,' and 'Review'*

## Third-level: Identify Link Between Social Demographic and Medical Condition

*Ex. 'Race-Disease'*

# Negative Samples



- Samples marked as biased without any other label

- Samples containing all labels

# Extracted Negatives

**Table 1: The term “white” in a racial vs. non-racial context**

Race-Related	Not Race-Related
“5 Year Relative Survival: overall 84% for <u>white</u> women, 62% for black women, 95% for local disease, 69% regional disease (spread to lymph node), 17% for distant disease.”	“ <u>White</u> matter within the spinal cord contains the axons of neurons that are ascending and descending to transmit signals to and from the brain, respectively.”



**Uninformative Extracted Negative**

# Extracted Negatives

**Table 1: The term “white” in a racial vs. non-racial context**

Race-Related	Not Race-Related
“5 Year Relative Survival: overall 84% for <u>white</u> women, 62% for black women, 95% for local disease, 69% regional disease (spread to lymph node), 17% for distant disease.”	“ <u>White</u> matter within the spinal cord contains the axons of neurons that are ascending and descending to transmit signals to and from the brain, respectively.”

1

0



# 4

## Word Sense Disambiguation

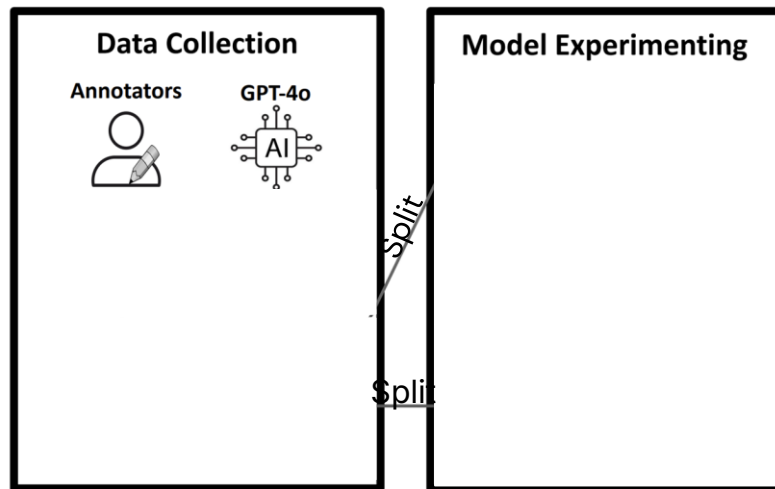
Machine Learning for Word Sense  
Disambiguation and Classification

# What is the most effective model for word sense disambiguation in a medical context? Can we produce accurate results?

**Table 1: The term “white” in a racial vs. non-racial context**

Race-Related	Not Race-Related
“5 Year Relative Survival: overall 84% for <b>white</b> women, 62% for black women, 95% for local disease, 69% regional disease (spread to lymph node), 17% for distant disease.”	“ <b>White</b> matter within the spinal cord contains the axons of neurons that are ascending and descending to transmit signals to and from the brain, respectively.”

# Word Sense Disambiguation (WSD) Experiments



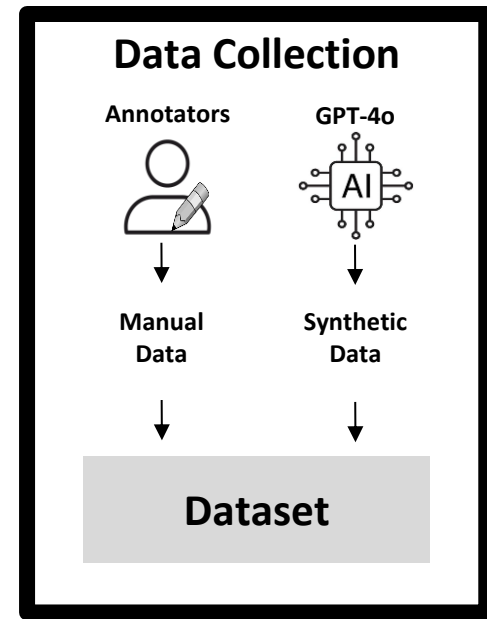
**Figure 2: WSD training and evaluation. Excerpts manually labeled as race-related or not plus GPT-generated sentences are used to train and evaluate the WSD models.**

# Data for WSD

Extracted negatives randomly sampled

- Human expert labels data
  - Label: 1 if sample relates to a social demographic
  - Label: 0 otherwise
- Additional samples synthetically generated using GPT-4o
  - More accurate results?

Demographics of interest include Race and Ethnicity



# WSD Problem Statement

Given a set of words  $\mathcal{W}$  and a set of senses  $\mathcal{S}_w = \{s_w^{(1)}, \dots, s_w^{(k)}\}$  for each  $w \in \mathcal{W}$

and a context (i.e. an ordered sequence of words)

$x = (x_1, \dots, x_{i-1}, w, x_{i+1}, \dots, x_n) \in \mathcal{X}$

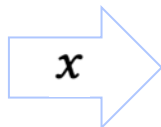
We are interested in determining if a term  $w$  is related to a sense  $s$  in an excerpt  $x$

$\text{IsRelated}(w, x, s_w) \in \{\text{TRUE}, \text{FALSE}\}$

# WSD Example

$S$   
= race/ethnicity

$\mathcal{W}$   
= {'white', 'Black'...}

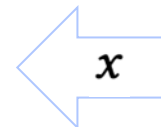


**Black** youth less likely to be diagnosed with MDD, Bipolar, or substance use disorder than **white** youth



$\text{IsRELATED}(w, x, S_w) = \text{TRUE}$

**White** matter and Grey matter anatomy of the spinal cord.



$\text{IsRELATED}(w, x, S_w) = \text{FALSE}$

# Evaluation of WSD models

**Table 2: Performance metrics for WSD on manually-annotated+GPT excerpts. Best result for each metric shown in bold. GlossBERT and GPT-4o are tied as the best models.**

<b>Metric</b>	<b>TF-IDF+ Logistic Reg.</b>	<b>ALBERT</b>	<b>Gloss BERT</b>	<b>GPT-3.5 Turbo</b>	<b>GPT-4o mini</b>
Accuracy	0.839	0.926	<b>0.944</b>	0.925	<b>0.944</b>
Precision	0.816	0.935	<b>0.936</b>	0.916	<b>0.936</b>
Recall	0.839	0.977	<b>1.000</b>	<b>1.000</b>	<b>1.000</b>
F1 Score	0.817	0.956	<b>0.967</b>	0.956	<b>0.967</b>

# Evaluation of WSD models

**Table 3: Examples of WSD test cases and GlossBERT predicted probabilities for  $y = 1$ . Each excerpt has a term (bolded) listed among race/ethnicity keywords.**

Input $x$ (label $y$ )	Prediction
Melanoma: increasing in incidence in the <b>white</b> population (CDC). ( $y = 1$ )	0.9998
2015 <b>American</b> Heart Association guidelines suggest treating patients presenting with systolic BP above 150-220 mmHg, but they do not offer a specific BP target. ( $y = 0$ )	0.9998
Calcific plaques are chalky <b>white</b> and arise from cardiac (aortic and mitral) valves. ( $y = 0$ )	0.0001



---

# 05

## Bias Classification

Approach and Evaluation of Bias Detection Models

# Bias Example

"They promote hair growth in the groin, axilla, chest and face, yet they also promote hair loss in the scalp in **men** who are genetically susceptible to androgenetic alopecia."

Label: Biased  
Category: Gender Bias

Reasoning: "Use sex terms when speaking of populations, should be male instead of men. Also, include citation to support this assertion."

# Bias Classification Problem Statement

- Formally, Salavati et al. define type-specific bias as a binary label  $\text{BIAS}(x, t) \in \{\text{TRUE}, \text{FALSE}\}$  indicating whether excerpt  $x$  is biased with respect to a social identifier category  $t$
- In the present work, we consider only the *general* definition of bias, regardless of which category  $t$  in a set  $\mathcal{T}$  it belongs to:

$$\text{BIAS}(x, \mathcal{T}) = \text{TRUE} \iff \exists t \in \mathcal{T} \text{ s.t. } \text{BIAS}(x, t) = \text{TRUE}$$

# Bias Classification Data

LN

Negatives  
labeled by  
human  
annotators

XN

Negatives  
extracted by  
use of social  
identifiers

XN\*

Extracted  
Negatives  
filtered using  
WSD

# Bias Classification Data Sets

LN

Human  
annotated  
dataset

LN+XN

Dataset  
used by  
Salavati et al.

LN+XN\*

Refined  
dataset by  
Salavati et al  
using WSD

# Evaluation of Bias Detection models

Table 4: Performance metrics and 95%-CIs for RoBERTa, TinyLlama trained on dataset variants (LN+XN\*, LN+XN, LN). Best results among each model variants (resp. across all models) and statistical ties shown are bolded (resp. underlined).

Metric	RoBERTa			TinyLlama		
	LN+XN*	LN+XN	LN	LN+XN*	LN+XN	LN
Precision	<b><math>0.613 \pm 0.015</math></b>	<b><math>0.640 \pm 0.021</math></b>	$0.526 \pm 0.029$	<b><u><math>0.675 \pm 0.008</math></u></b>	<b><u><math>0.693 \pm 0.028</math></u></b>	$0.536 \pm 0.020$
Recall	<b><math>0.692 \pm 0.024</math></b>	$0.667 \pm 0.023$	<b><u><math>0.719 \pm 0.026</math></u></b>	<b><math>0.548 \pm 0.030</math></b>	$0.519 \pm 0.029$	<b><math>0.607 \pm 0.035</math></b>
F1 Score	<b><u><math>0.650 \pm 0.013</math></u></b>	<b><u><math>0.652 \pm 0.017</math></u></b>	$0.606 \pm 0.017$	<b><math>0.604 \pm 0.021</math></b>	<b><math>0.593 \pm 0.017</math></b>	<b><math>0.568 \pm 0.016</math></b>
F2 Score	<b><u><math>0.674 \pm 0.019</math></u></b>	<b><math>0.661 \pm 0.016</math></b>	<b><math>0.669 \pm 0.016</math></b>	<b><math>0.569 \pm 0.027</math></b>	<b><math>0.546 \pm 0.024</math></b>	<b><math>0.591 \pm 0.025</math></b>
AUC	<b><u><math>0.927 \pm 0.003</math></u></b>	<b><u><math>0.930 \pm 0.009</math></u></b>	$0.910 \pm 0.008$	<b><math>0.907 \pm 0.005</math></b>	<b><math>0.903 \pm 0.005</math></b>	$0.871 \pm 0.011$

# Evaluation of Bias Detection models

**Table 5: Performance Metrics and 95%-CIs for Fine-Tuned Models against Baseline (\*Salavati et al., 2024). Best results and statistical ties shown in bold.**

Metric	RoBERTa	TinyLlama	Baseline*
Precision	0.613 $\pm$ 0.015	<b>0.675 <math>\pm</math> 0.008</b>	0.504 $\pm$ 0.054
Recall	0.692 $\pm$ 0.024	0.548 $\pm$ 0.030	<b>0.812 <math>\pm</math> 0.069</b>
F1 Score	<b>0.650 <math>\pm</math> 0.014</b>	0.604 $\pm$ 0.021	0.615 $\pm$ 0.022
F2 Score	0.674 $\pm$ 0.019	0.569 $\pm$ 0.027	<b>0.717 <math>\pm</math> 0.027</b>
AUC	<b>0.927 <math>\pm</math> 0.003</b>	0.907 $\pm$ 0.005	<b>0.923 <math>\pm</math> 0.004</b>

# Evaluation of Bias Detection models

**Table 6: Performance Metrics and 95%-CIs for Prompting GPT-4o mini. Best results for each metric shown in bold. AUC was ommitted as it cannot be computed for binary outputs.**

Metric	Zero-Shot	Few-Shot
Precision	<b>0.367</b> $\pm$ <b>0.071</b>	0.259 $\pm$ 0.019
Recall	0.260 $\pm$ 0.029	<b>0.610</b> $\pm$ <b>0.026</b>
F1 Score	<b>0.303</b> $\pm$ <b>0.040</b>	<b>0.363</b> $\pm$ <b>0.023</b>
F2 Score	0.274 $\pm$ 0.032	<b>0.480</b> $\pm$ <b>0.025</b>

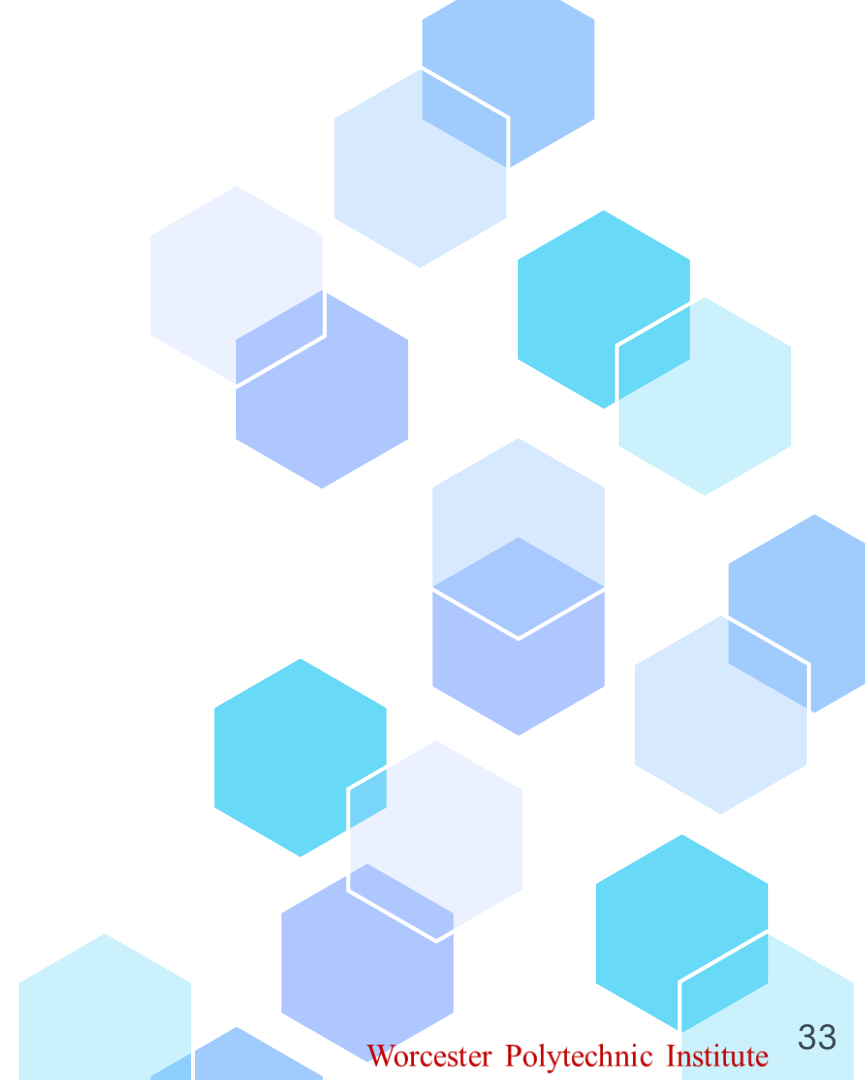


---

# 06

# Conclusion

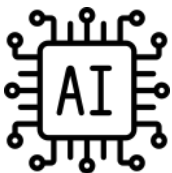
Key Findings and Future Implications



# Conclusion

- **Health-related applications** and **recommender systems** are prone to biases
- Developed a framework to **detect and diagnose bias** in medical curricula by an emphasized focus on data over model
- **WSD** models were **highly effective** at distinguishing biased excerpts from non-biased ones
- While **prompt engineering** of **LLMs** can handle many tasks, they are **not well-suited** for health related bias classification

# Discussion



Further explore the potential of  
ChatGPT-4o (or **other future**  
**OpenAI models**)



Use of **other bias categories**  
(E.g. geography)



Use case in **other domains**  
(Crucial role of tone in  
determining word meaning esp.  
in social media)



**Challenges with LLMs**  
(Computational cost, time  
constraints, accessibility  
issues)



# WPI

---

## Thank You. Questions?



Link to paper