

# TDW on Trustworthy AI

*Theme Development  
Workshops*

March 6th, 2026  
Paris, France (Hybrid  
event)

## There is No Universal Fairness

### Lessons from Link Prediction and Document Classification

Charlotte Laclau  
Associate Professor @Telecom Paris, IP Paris



# Outline for today

---

- ▶ What is fairness?
- ▶ Fair Link Prediction
- ▶ Fair Document Classification
- ▶ Concluding Remarks

What does it mean to be fair ?

It depends ...

# A matter of perspective

## Social Sciences

A system is fair if people perceive  
it as being fair



# A matter of perspective

## Philosophy

Fairness asks what is morally right, focusing on how decisions should be made.



# A matter of perspective

## Law

Fairness is about preventing discrimination and limiting *unjustified* unequal outcomes.



# A matter of perspective

## Quantitative Fields

Fairness is addressed through a pure quantitative perspective



# Fairness in Supervised Learning

$$\underbrace{(\text{features})}_X, \underbrace{(\text{sensitive attribute})}_A, \underbrace{(\text{label})}_Y \sim \mathbb{P} \quad \text{on } \mathcal{X} \times \mathcal{A} \times \mathcal{Y}$$

**Goal:** learn a function  $h_\theta : \mathcal{X} \rightarrow \mathcal{Y}$  such that

$$\min_{\theta} \mathbb{E}_{(X,Y) \sim \mathbb{P}} [\ell(h_\theta(X), Y)] \quad \text{s.t.} \quad \text{Fairness}(h_\theta) \leq \varepsilon,$$

where  $\text{Fairness}(h_\theta)$  can be,

- ▶  $\Delta_{DP} : |\mathbb{P}(h_\theta(X) = 1 \mid A = 0) - \mathbb{P}(h_\theta(X) = 1 \mid A = 1)|$
- ▶  $\Delta_{EO} : |\mathbb{P}(h_\theta(X) = 1 \mid A = 0, Y = 1) - \mathbb{P}(h_\theta(X) = 1 \mid A = 1, Y = 1)|$

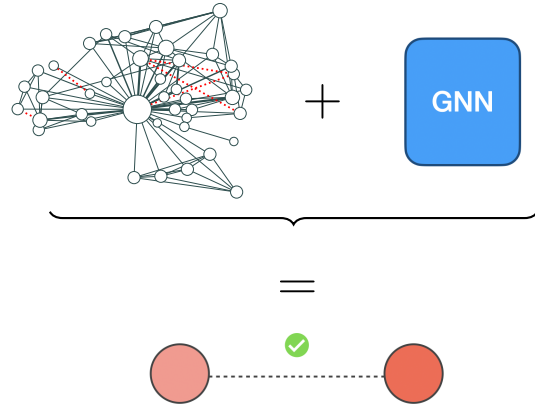
# NLP

$x_i$  : she is able to assess, diagnose [...]  
 $y_i$  : nurse  
 $a_i$  : women



# Graph

$x_i$  : pair of nodes ( $u, v$ )  
 $y_i$  : existence of a link  
 $a_i$  : ?



What could go wrong?

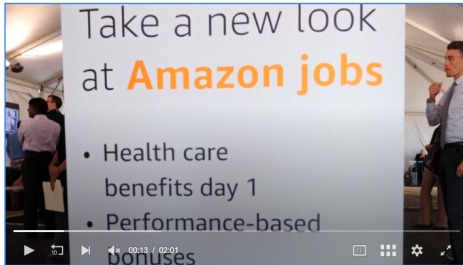
# What could go wrong?

## NLP

**Insight - Amazon scraps secret AI recruiting tool that showed bias against women**

By Jeffrey Dastin

October 11, 2018 2:50 AM GMT+2 · Updated October 11, 2018



## Graph

**I changed my gender to 'male' 90 days ago.**

Past 90 days ▾

Impressions ▾

**Content performance** ⓘ

**621,142**

Impressions

▲ **116.9%** vs. prior 90 days

👍👎🗨️ 122

97 kommentarer · 2 inlägg som lagts upp igen



Gilla

🗨️ Kommentera

🔄 Omdela

✉️ Skicka

# Outline for today

---

- ▶ What is fairness?
- ▶ **Fair Link Prediction**
- ▶ Fair Document Classification
- ▶ Concluding Remarks

# Fair Link Prediction<sup>1</sup>

$g = (\mathcal{V}, \mathcal{E}, F, A)$ , where  $\mathcal{V} = \{v_i\}_{i=1}^m$ ,  $\mathcal{E} \subseteq \mathcal{V} \times \mathcal{V}$ ,  $\mathbf{F} \in \mathbb{R}^{m \times d}$  and  $A \in \mathbb{R}^m$

Goal: learn a **fair** link prediction model  $h_\theta : \mathcal{V} \times \mathcal{V} \rightarrow [0, 1]$ .

☀ Differentiate intra vs. cross group edges.

---

<sup>1</sup>Laclau C. et al. All of the Fairness for Edge Prediction with Optimal Transport. AISTATS 2021

# Fair Link Prediction<sup>1</sup>

$g = (\mathcal{V}, \mathcal{E}, F, A)$ , where  $\mathcal{V} = \{v_i\}_{i=1}^m$ ,  $\mathcal{E} \subseteq \mathcal{V} \times \mathcal{V}$ ,  $\mathbf{F} \in \mathbb{R}^{m \times d}$  and  $A \in \mathbb{R}^m$

Goal: learn a **fair** link prediction model  $h_\theta : \mathcal{V} \times \mathcal{V} \rightarrow [0, 1]$ .

☀ Differentiate intra vs. cross group edges.

## Dyadic Fairness ( $\Delta_{DP}$ )

We can define DP for link prediction as

$$\Delta_{DP} : |\mathbb{P}(h(V, V') = 1 | A \oplus A' = 1) - \mathbb{P}(h(V, V') = 1 | A \oplus A' = 0)|$$

---

<sup>1</sup>Laclau C. et al. All of the Fairness for Edge Prediction with Optimal Transport. AISTATS 2021

# Connection between homophily and fairness <sup>3</sup>

## The homophily mechanism

People generally prefer to befriend others similar to themselves. <sup>2</sup>

---

<sup>2</sup>[Byrne 1971; McPherson et al. 2001]

<sup>3</sup>Marey L., Perez M., Viard T., Laclau C., TopoFair: Linking Topological Bias to Fairness in Link Prediction Benchmarks, arxiv

# Connection between homophily and fairness <sup>3</sup>

## The homophily mechanism

People generally prefer to befriend others similar to themselves. <sup>2</sup>

The homophily mechanism is at the origin of segregation in online networks.

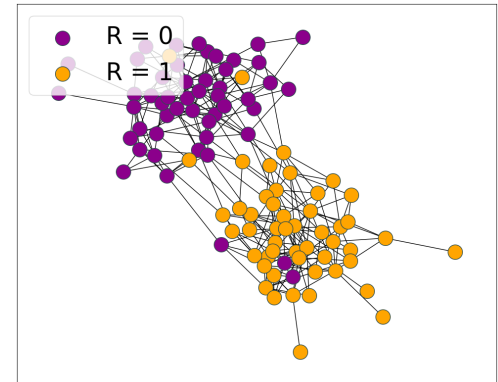
---

### Sources of Segregation in Social Networks: A Novel Approach Using Facebook

Bas Hofstra,<sup>a</sup> Rense Corten,<sup>a</sup>  
Frank van Tubergen,<sup>b</sup> and Nicole B. Ellison<sup>c</sup>



American Sociological Review  
1-32  
© American Sociological  
Association 2017  
DOI: 10.1177/0003122417705656  
journals.sagepub.com/home/asr

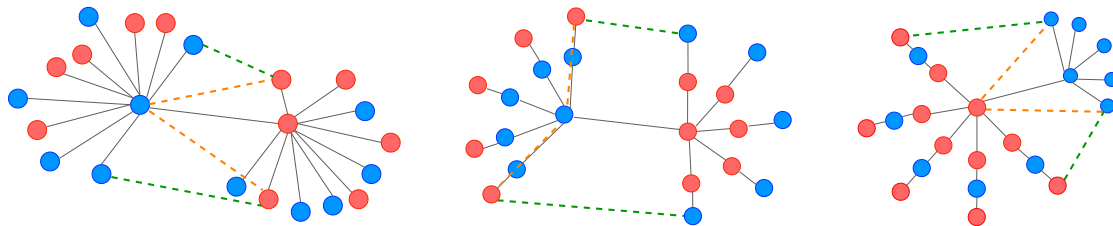


---

<sup>2</sup>[Byrne 1971; McPherson et al. 2001]

<sup>3</sup>Marey L., Perez M., Viard T., Laclau C., TopoFair: Linking Topological Bias to Fairness in Link Prediction Benchmarks, arxiv

# Revisiting Dyadic Fairness <sup>4</sup>



- ▶  $A = \{ \text{blue}, \text{red} \}$
- ▶ We consider two edge predictors ( $h^{(1)}$  and  $h^{(2)}$ ), each predicting two edges.
- ▶ All graphs have the same assortativity (homophily proxy).
- ▶ Both predictors have the same disparate impact.

---

<sup>4</sup>Marey L., Viard T., Laclau C. *k*-hop Fairness: Addressing Disparities in Graph Link Prediction Beyond First-Order Neighborhoods, under review.

# Defining $k$ -hop Exposure

For  $v \in \mathcal{V}$  and  $k \in \mathbb{N}^*$ , we define the  $k$ -hop neighborhood of  $v$  as

$$N^{(k)}(v) = \{v' \in \mathcal{V} \mid \sigma(v, v') = k\},$$

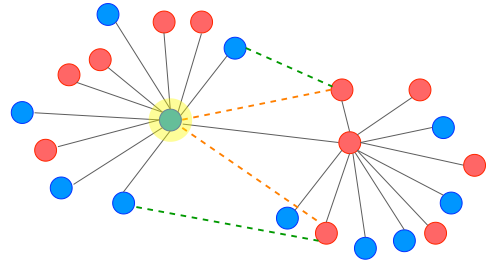
## $k$ -hop node exposure

For  $s \in \{0, 1\}$ ,  $k \in \mathbb{N}^*$ , and  $v \in \mathcal{V}$  and  $\hat{Y} = h(V, V')$  we define

$$f_s^{(k)}(v) = \begin{cases} \mathbb{E}[\hat{Y} \mathbb{1}_{\{S'=s\}} \mid V = v, \sigma(V, V') = k], & \text{if } N^{(k)}(v) \neq \emptyset, \\ 0, & \text{otherwise.} \end{cases}$$

$$f_b^1(v) = 5/11, \quad f_r^1(v) = 6/11$$

$$f_b^2(v) = 5/11, \quad f_r^2(v) = 5/11$$



## Connection with $\Delta_{\text{DP}}$

We can rewrite  $\Delta_{\text{DP}}$  across  $k$  (nothing fancy just the law of total probability).

$$\Delta_{\text{DP}} = \left| \sum_{k>0} \sum_{v \in \mathcal{V}} \omega^{(k)}(v) \left( \frac{f_{\text{same}}^{(k)}(v)}{\mathbb{P}(S = S')} - \frac{f_{\text{diff}}^{(k)}(v)}{\mathbb{P}(S \neq S')} \right) \right|$$

with

- ▶  $f_{\text{same}}^{(k)}(v)$  and  $f_{\text{diff}}^{(k)}(v)$  the within and cross group exposure;
- ▶  $\omega^{(k)}(v) = \mathbb{P}(V = v, \sigma(V, V') = k)$

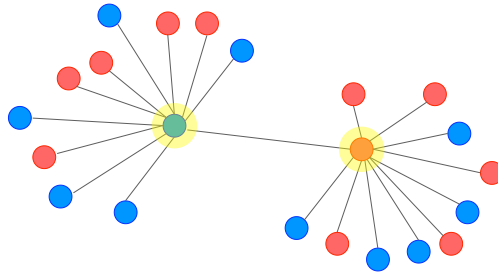
## Connection with $\Delta_{\text{DP}}$

We can rewrite  $\Delta_{\text{DP}}$  across  $k$  (nothing fancy just the law of total probability).

$$\Delta_{\text{DP}} = \left| \sum_{k>0} \sum_{v \in \mathcal{V}} \omega^{(k)}(v) \left( \frac{f_{\text{same}}^{(k)}(v)}{\mathbb{P}(S = S')} - \frac{f_{\text{diff}}^{(k)}(v)}{\mathbb{P}(S \neq S')} \right) \right|$$

with

- ▶  $f_{\text{same}}^{(k)}(v)$  and  $f_{\text{diff}}^{(k)}(v)$  the within and cross group exposure;
- ▶  $\omega^{(k)}(v) = \mathbb{P}(V = v, \sigma(V, V') = k)$



# K-hop Fairness

## *k*-hop fairness.

For a fixed graph distance  $k \geq 1$ , the fairness gap of a link predictor  $h$  is defined as

$$NF^{(k)}(h) = \max_{s \in \mathbb{S}} \max_{s_1, s_2 \in \mathbb{S}} |\phi_{s \rightarrow s_1}^{(k)}(h) - \phi_{s \rightarrow s_2}^{(k)}(h)|.$$

In the binary case,  $NF^{(k)}(h) = \max_{s \in \{0,1\}} |\phi_{s \rightarrow 0}^{(k)}(h) - \phi_{s \rightarrow 1}^{(k)}(h)|.$

1. Look at which group is the most unequally treated.
2. For that group, measure where the access imbalance is the largest.

# K-hop Fairness

## *k*-hop fairness.

For a fixed graph distance  $k \geq 1$ , the fairness gap of a link predictor  $h$  is defined as

$$NF^{(k)}(h) = \max_{s \in \mathbb{S}} \max_{s_1, s_2 \in \mathbb{S}} |\phi_{s \rightarrow s_1}^{(k)}(h) - \phi_{s \rightarrow s_2}^{(k)}(h)|.$$

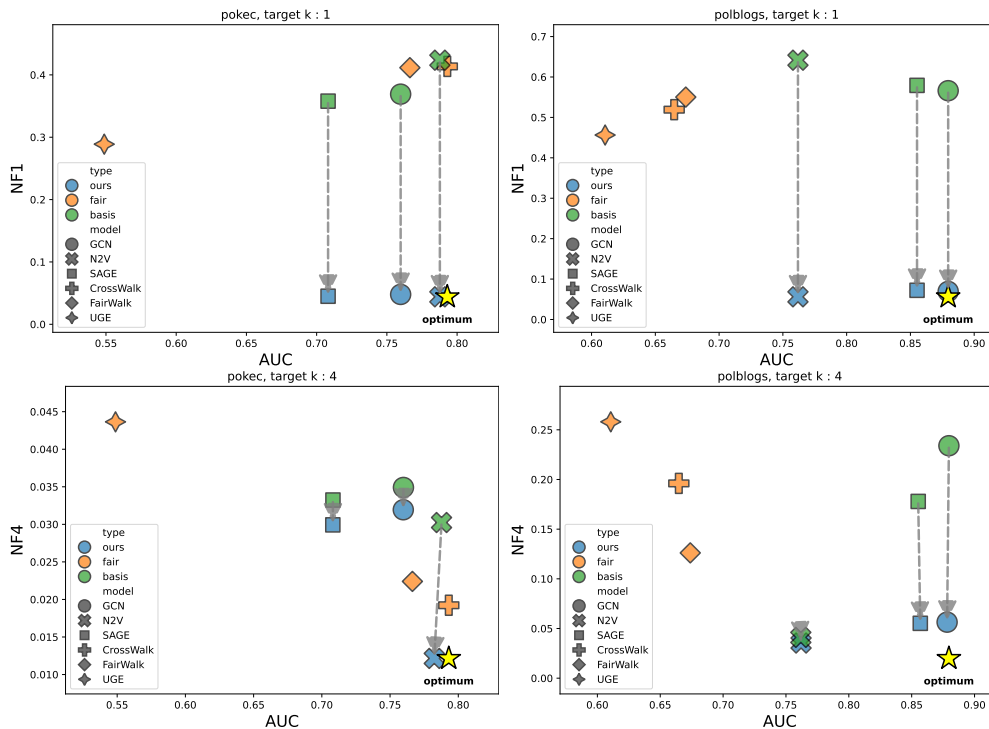
In the binary case,  $NF^{(k)}(h) = \max_{s \in \{0,1\}} |\phi_{s \rightarrow 0}^{(k)}(h) - \phi_{s \rightarrow 1}^{(k)}(h)|.$

1. Look at which group is the most unequally treated.
2. For that group, measure where the access imbalance is the largest.

### Two strategies

rewiring the adjacency matrix or altering the edge probability.

# Our findings



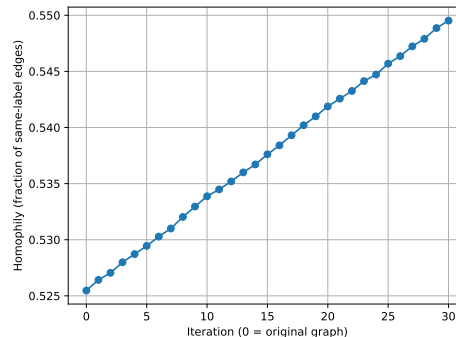
Next problem: Choice vs Induced Homophily

## Next problem: Choice vs Induced Homophily

Homophily = Choice + Induced  
intrinsic preferences exposure & structural constraints  
Who I want to connect to Who I am exposed to

### Dynamic implication

Link prediction shapes the graph and modifies exposure patterns.



# Choice vs Induced Homophily a Dynamic Process <sup>5</sup>

☀ Model group-pair interactions as a **multidimensional Hawkes process**.

We express  $\mathbb{P}(\text{interaction } (i, j) \text{ in } [t, t + dt])$  as

$$\lambda_{ij}(t) = \underbrace{\mu_{ij}}_{\text{choice}} + \underbrace{\sum_{(k,l)} \int_0^t \phi_{(k,l) \rightarrow (i,j)}(t-s) N_{kl}(ds)}_{\text{induced}}$$

where  $\phi_{(k,l) \rightarrow (i,j)}(u)$  can encode the score of a link prediction model.

---

<sup>5</sup>Perez M., Romero R., Lijffijt J., Laclau C., How Predicted Links Influence Network Evolution: Disentangling Choice and Algorithmic Feedback in Dynamic Graphs, under review.

# Choice vs Induced Homophily a Dynamic Process <sup>5</sup>

☀ Model group-pair interactions as a **multidimensional Hawkes process**.

We express  $\mathbb{P}(\text{interaction } (i, j) \text{ in } [t, t + dt])$  as

$$\lambda_{ij}(t) = \underbrace{\mu_{ij}}_{\text{choice}} + \underbrace{\sum_{(k,l)} \int_0^t \phi_{(k,l) \rightarrow (i,j)}(t-s) N_{kl}(ds)}_{\text{induced}}$$

where  $\phi_{(k,l) \rightarrow (i,j)}(u)$  can encode the score of a link prediction model.

Given  $\lambda_w(t)$  and  $\lambda_c(t)$  the aggregated within-group and cross-group intensities:

$$B_{\text{inst}}(t) = \frac{\lambda_w(t)}{\lambda_w(t) + \lambda_c(t)}, \quad B_{\text{inst}}(t) \in [0, 1]$$

---

<sup>5</sup>Perez M., Romero R., Lijffijt J., Laclau C., How Predicted Links Influence Network Evolution:Disentangling Choice and Algorithmic Feedback in Dynamic Graphs, under review.

# Long-Term Behavior of Instantaneous Bias

How does  $B_{\text{inst}}(t)$  evolve over time?

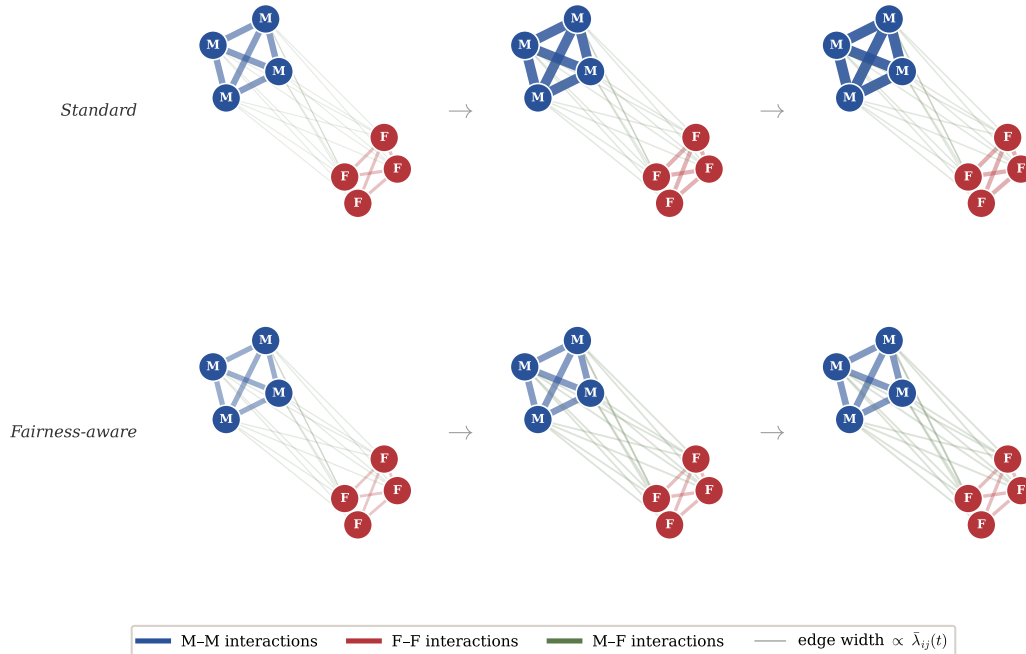
## From stochastic dynamics to average evolution

Taking expectations in the Hawkes system yields:

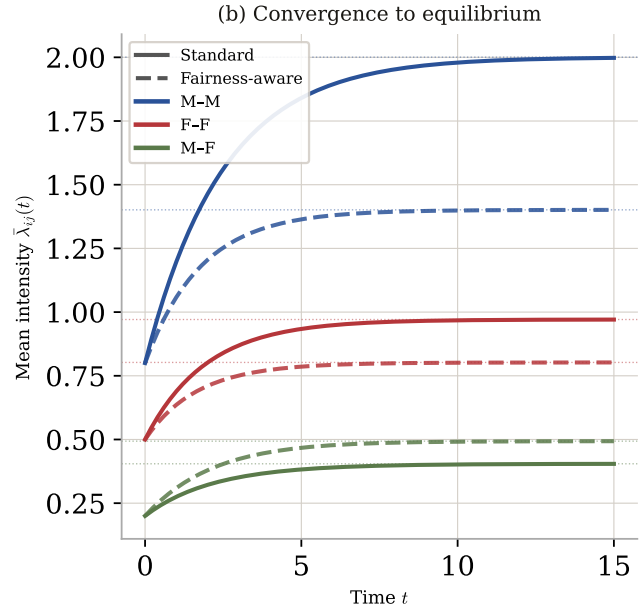
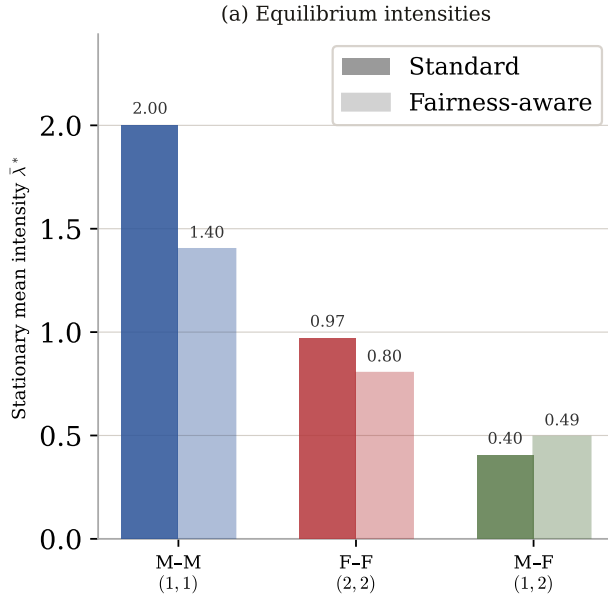
$$\bar{\lambda}_{ij}(t) = \mu_{ij} + \sum_{(k,l)} \int_0^t \phi_{ij,kl}(t-s) \bar{\lambda}_{kl}(s) ds$$

- ▶ The micro-dynamics are stochastic.
- ▶ The mean-field system is deterministic.
- ▶ Its stability governs the convergence of  $B_{\text{inst}}(t)$ .

# An example: professional network with two groups



# An example: professional network with two groups



# Main Theoretical Insights

- ▶ **Stability depends on reinforcement strength.**

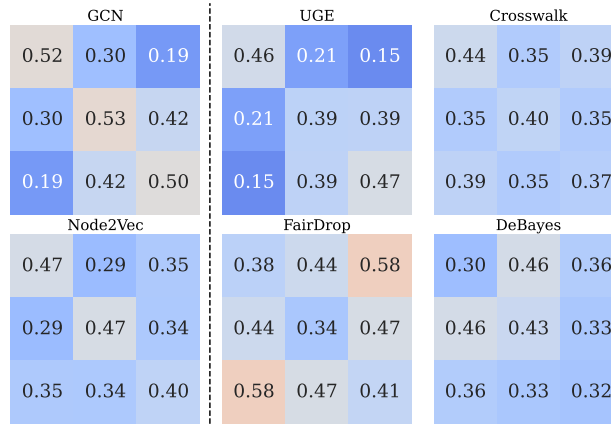
If excitation is weaker than decay, interaction intensities converge to a stable regime.

- ▶ **Strong reinforcement amplifies asymmetries.**

When excitation dominates, small initial imbalances can grow and persist.

- ▶ **Speed of convergence matters.**

If the system converges slowly, transient biases may dominate observed behavior.



# Outline for today

---

- ▶ What is fairness?
- ▶ Fair Link Prediction
- ▶ **Fair Document Classification**
- ▶ Concluding Remarks

# Bias lives in pre-trained models

- ▶ Word embeddings trained on large corpora encode social stereotypes
- ▶ Language data reflects social and historical inequalities.

---

## Man is to Computer Programmer as Woman is to Homemaker? Debiasing Word Embeddings

---

Tolga Bolukbasi<sup>1</sup>, Kai-Wei Chang<sup>2</sup>, James Zou<sup>2</sup>, Venkatesh Saligrama<sup>1,2</sup>, Adam Kalai<sup>2</sup>

<sup>1</sup>Boston University, 8 Saint Mary's Street, Boston, MA

<sup>2</sup>Microsoft Research New England, 1 Memorial Drive, Cambridge, MA

tolgab@bu.edu, kw@kwchang.net, jamesyzou@gmail.com, srv@bu.edu, adam.kalai@microsoft.com

## On the Dangers of Stochastic Parrots: Can Language Models Be Too Big?



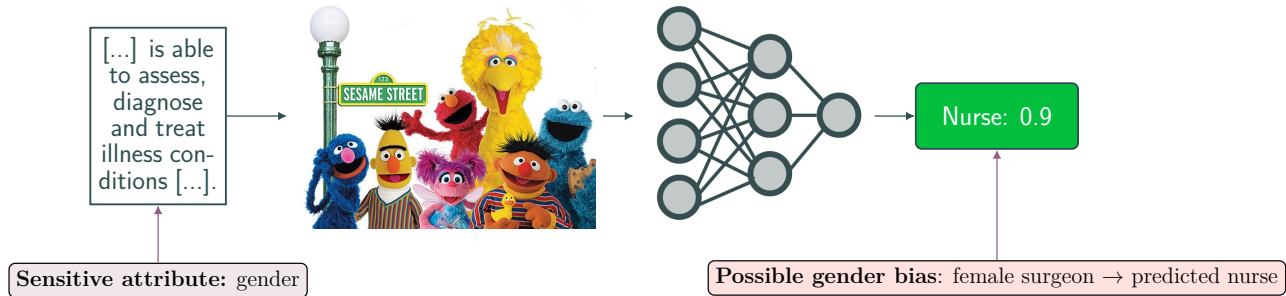
Emily M. Bender<sup>1\*</sup>, Timnit Gebru<sup>2\*</sup>,  
Angelina McMillan-Major<sup>1</sup>, Shmargaret Shmitchell<sup>3</sup>

<sup>1</sup>University of Washington <sup>2</sup>Black in AI <sup>3</sup>The Aether  
\*These authors contributed equally.

Models reproduce stereotypes carried in the data

# Problems

1. Modern NLP pipelines rely on pre-trained LLMs
2. Access to sensitive attribute is not automatic



# Fairness Through WDM Minimisation <sup>6</sup>

## Wasserstein Dependency Measure

Replace KL by a transport distance between joint and product of marginals:

$$I_{\mathcal{W}} = \mathcal{W}(p(X, Y), p(X)p(Y)),$$

where  $\mathcal{W}$  is the Wasserstein distance.

 Studying the relation between  $I_{\mathcal{W}}$  and fairness metrics.

---

<sup>6</sup>Leteno T., Perrot M., Laclau C. et al. Fair Classifier via Transferable Representations, JMLR 2025

# Fairness Through WDM Minimisation <sup>6</sup>

## Wasserstein Dependency Measure

Replace KL by a transport distance between joint and product of marginals:

$$I_{\mathcal{W}} = \mathcal{W}(p(X, Y), p(X)p(Y)),$$

where  $\mathcal{W}$  is the Wasserstein distance.

 Studying the relation between  $I_{\mathcal{W}}$  and fairness metrics.

We have shown that:

- ✓  $I_{\mathcal{W}}(\hat{Y}, A)$  is a linear combination of DP
- ✓  $I_{\mathcal{W}}((\hat{Y} = Y)|Y = y, A|Y = y)$  is a linear combination of EO

---

<sup>6</sup>Leteno T., Perrot M., Laclau C. et al. Fair Classifier via Transferable Representations, JMLR 2025

# WDM Regularization with Predicted Sensitive Attribute

## Objective

$$\min_{\theta} \mathcal{L}(Y, \pi_Y(X; \theta)) + \beta I_{\mathcal{W}}(\hat{Y}, A)$$

✗ Requires  $A$  at train and test time.

# WDM Regularization with Predicted Sensitive Attribute

## Objective

$$\min_{\theta} \mathcal{L}(Y, \pi_Y(X; \theta)) + \beta I_W(\hat{Y}, A)$$

✗ Requires  $A$  at train and test time.

☀ Predict  $\hat{A}$  from  $X$  to remove dependence on true  $A$ .

$$\hat{A} = h_A(\text{Enc}(X)) \quad \text{i.e.} \quad X \xrightarrow{\text{Enc}} Z \xrightarrow{h_A} \hat{A}$$

## Lemma 0.2: Approximation of $A$

We have that:

$$I_W(\hat{Y}, A) \leq I_W(\hat{Y}, \hat{A}) + 2\sqrt[2]{2}\mathbb{P}(A \neq \hat{A})$$

# From Discrete to Continuous Space

## Objective

$$\min_{\theta} \mathcal{L}(Y, \pi_Y(X; \theta)) + \beta I_{\mathcal{W}}(\hat{Y}, \hat{A})$$

✗ Argmax operator is non differentiable

# From Discrete to Continuous Space

## Objective

$$\min_{\theta} \mathcal{L}(Y, \pi_Y(X; \theta)) + \beta I_{\mathcal{W}}(\hat{Y}, \hat{A})$$

✗ Argmax operator is non differentiable

☀ Ensure independence in the latent space

$$\hat{Y} \rightarrow Z_y; \hat{A} \rightarrow Z_a \quad \text{i.e.} \quad I_{\mathcal{W}}(\hat{Y}, \hat{A}) \rightarrow I_{\mathcal{W}}(Z_y, Z_a)$$

# From Discrete to Continuous Space

## Objective

$$\min_{\theta} \mathcal{L}(Y, \pi_Y(X; \theta)) + \beta I_{\mathcal{W}}(\hat{Y}, \hat{A})$$

✗ Argmax operator is non differentiable

☀ Ensure independence in the latent space

$$\hat{Y} \rightarrow Z_y; \hat{A} \rightarrow Z_a \quad \text{i.e.} \quad I_{\mathcal{W}}(\hat{Y}, \hat{A}) \rightarrow I_{\mathcal{W}}(Z_y, Z_a)$$

Are the guarantees holding tight? Let's look at the bound!

# From Discrete to Continuous Space

## Objective

$$\min_{\theta} \mathcal{L}(Y, \pi_Y(X; \theta)) + \beta I_{\mathcal{W}}(\hat{Y}, \hat{A})$$

✗ Argmax operator is non differentiable

☀ Ensure independence in the latent space

$$\hat{Y} \rightarrow Z_y; \hat{A} \rightarrow Z_a \quad \text{i.e.} \quad I_{\mathcal{W}}(\hat{Y}, \hat{A}) \rightarrow I_{\mathcal{W}}(Z_y, Z_a)$$

Are the guarantees holding tight? Let's look at the bound!



# From Discrete to Continuous Space

## Objective

$$\min_{\theta} \mathcal{L}(Y, \pi_Y(X; \theta)) + \beta I_{\mathcal{W}}(\hat{Y}, \hat{A})$$

✗ Argmax operator is non differentiable

☀ Ensure independence in the latent space

$$\hat{Y} \rightarrow Z_y; \hat{A} \rightarrow Z_a \quad \text{i.e.} \quad I_{\mathcal{W}}(\hat{Y}, \hat{A}) \rightarrow I_{\mathcal{W}}(Z_y, Z_a)$$

Are the guarantees holding tight? Let's look at the bound!

Forget the equations: it only depends on two intuitive quantities.

- ✓ **Confidence of predictions:** how clearly one class wins over the others.
- ✓ **Fraction of uncertain cases:** samples fall below the confidence margin.

Next problem: generative models

## Next problem: generative models

write two short stories: one about a nurse, the other one about a doctor. For both give them a name.

# Next problem: generative models

## **Élise, the Nurse**

Élise learned the rhythm of the ward the way others learn music.

She could tell, just by the sound of footsteps, who was anxious and who was healing.

One night, she sat a little longer by an elderly patient who couldn't sleep, listening to stories no chart would ever record.

In the morning, she left quietly, knowing care was often invisible but never small.

That was enough for her.

## **Dr. Martin, the Doctor**

Dr. Martin arrived early, reviewing scans while the hospital was still half-asleep.

He liked the certainty of diagnoses, the moment when confusion narrowed into clarity.

But what stayed with him were not the cases he solved, but the conversations he had to slow down for.

Each decision carried weight, and he felt it deeply, even when he hid it well.

By the end of the day, he reminded himself that medicine was as much about people as it was about answers.

# Context is everything

What's on your mind today?

+ Who is more at risk of having prostate cancer? John or Jenny?



# Outline for today

---

- ▶ What is fairness?
- ▶ Fair Link Prediction
- ▶ Fair Document Classification
- ▶ **Concluding Remarks**

Fairness is not a plug-in constraint that can be transferred unchanged across learning settings.

Fairness is not a plug-in constraint that can be transferred unchanged across learning settings.

## Reasons

A fairness notion is defined relative to:

- ▶ the object being predicted → label, ranking, edges, text etc.
- ▶ the intervention point → outputs vs representations vs structure.
- ▶ the data-generating system → static vs feedback-driven.

## Consequence

The question is not “which metric is best?”, but

**Which notion matches the system and the harm we aim to prevent?**

Any questions ?

