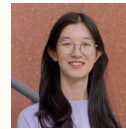


# Improving Fairness in Graph Neural Networks via Mitigating Sensitive Attribute Leakage

Yu Wang<sup>1</sup>



Yuying Zhao<sup>1</sup>



Yushun Dong<sup>2</sup>



Huiyuan Chen<sup>3</sup>



Jundong Li<sup>2</sup>



Tyler Derr<sup>1</sup>







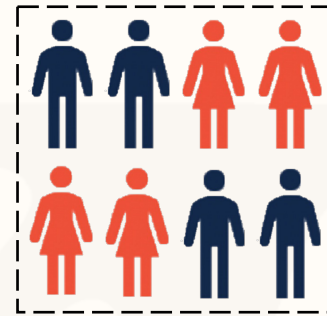
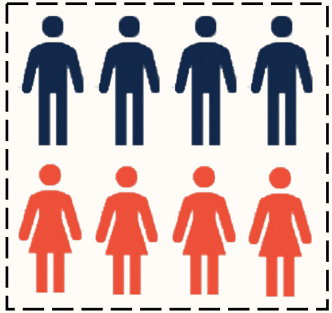
1. Network and Data Science Lab, Vanderbilt University

2. University of Virginia

3. Case Western Reserve University

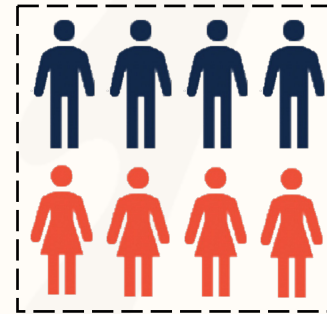
# Background – Group Fairness

s: group    
 $\hat{y}$ : label  



**Fair!**

$$\Delta_{sp} = 0$$



**Bias!**

$$\Delta_{sp} = 1$$

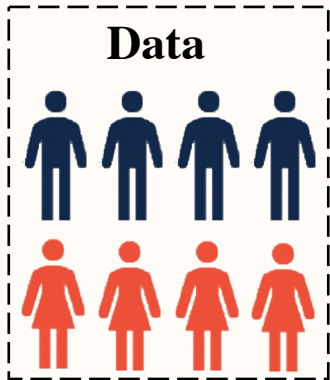
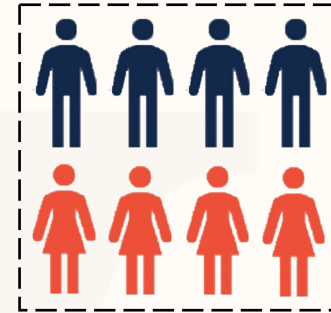
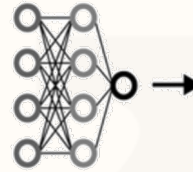
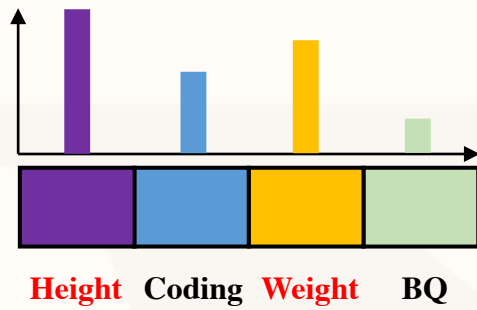


$$\Delta_{sp} = |P(\hat{y} = 1 | s = 0) - P(\hat{y} = 1 | s = 1)|$$

**s: sensitive attribute**

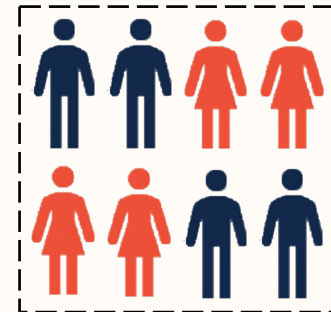
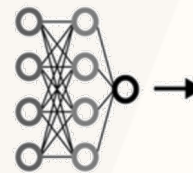
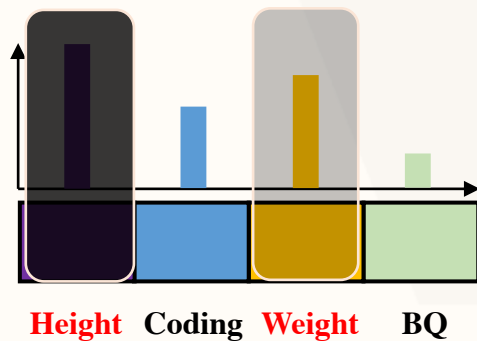
**Discriminative and unfair decision!**

# Background – Sensitive leakage

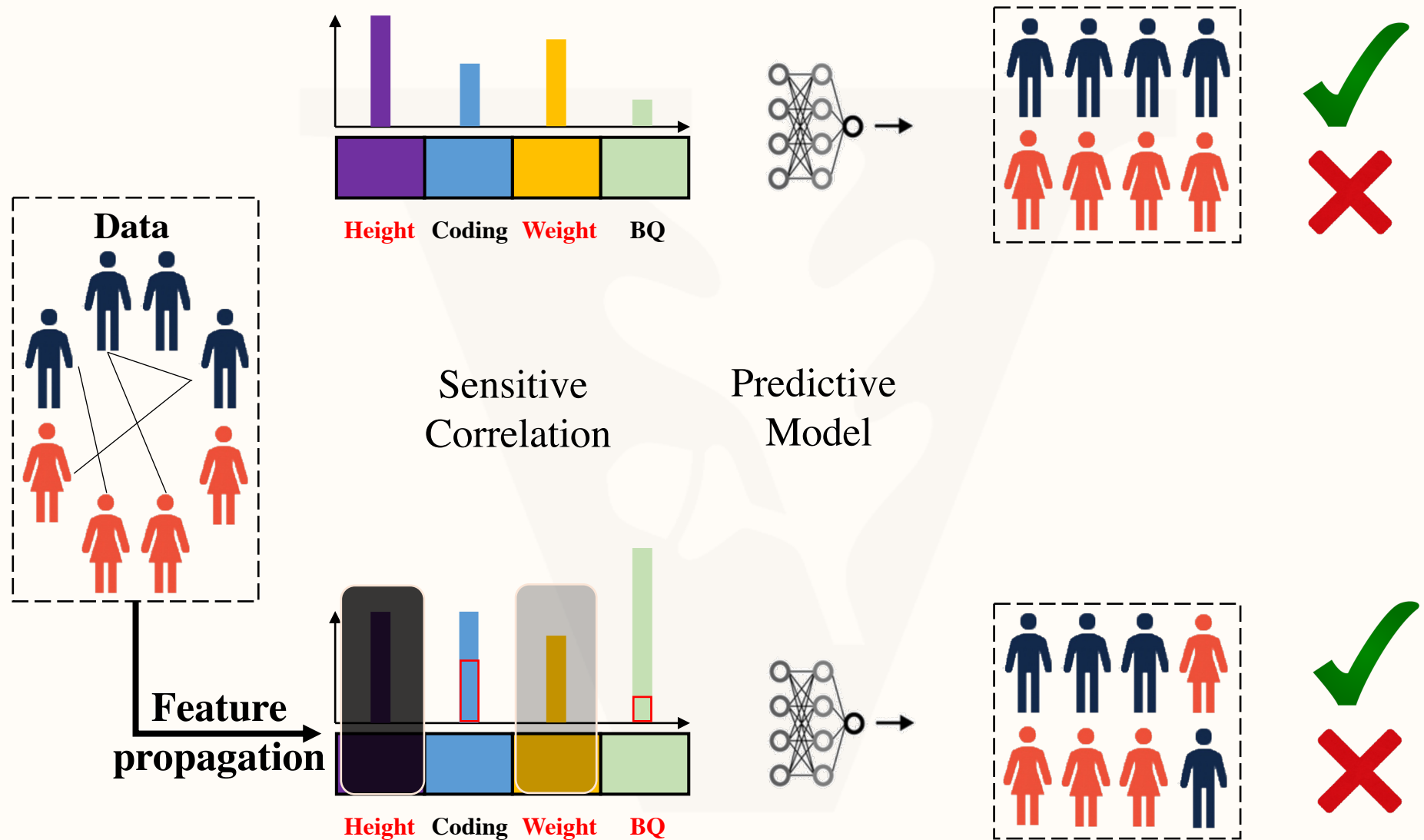


Sensitive  
Correlation

Predictive  
Model



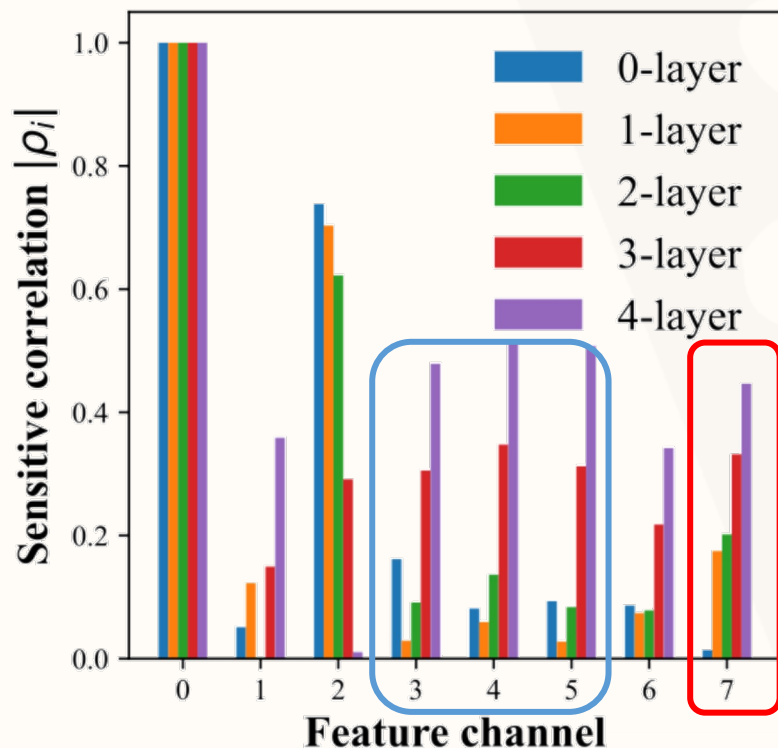
# Background – Correlation variation



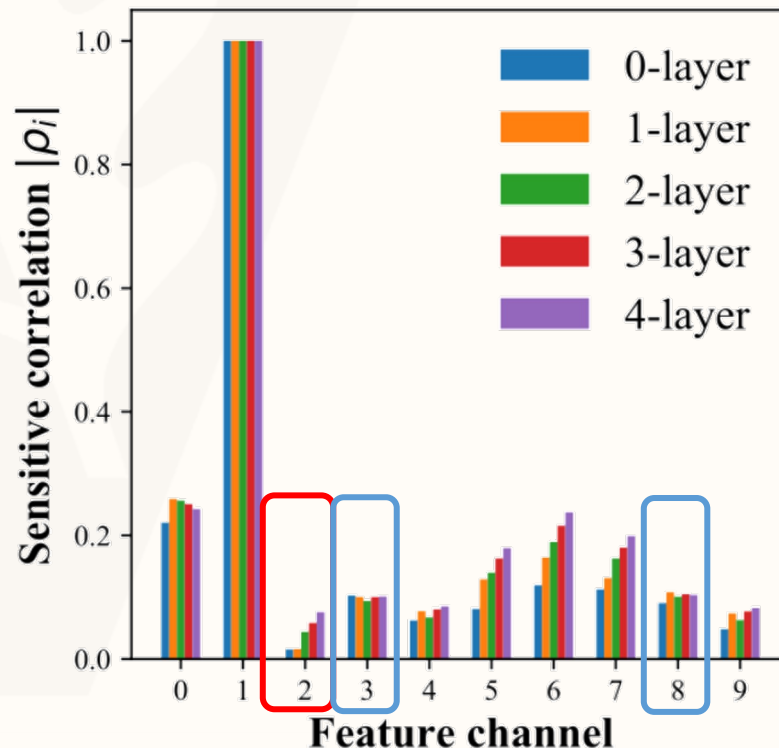
# Motivation – Correlation variation

Mask feature channels with higher correlation to the sensitive attributes

## German

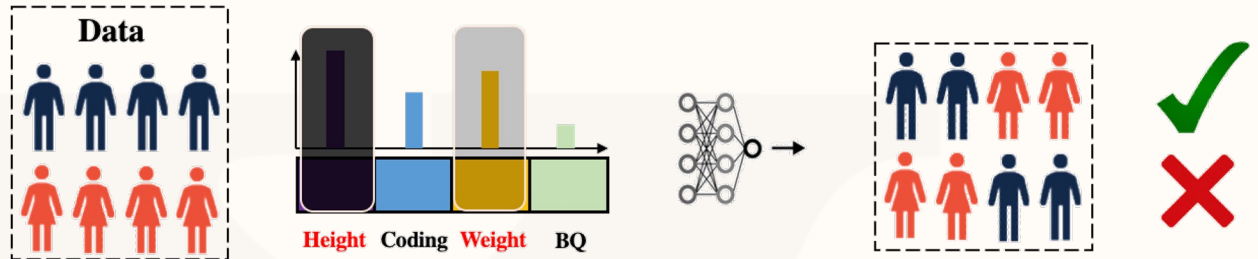


## Credit

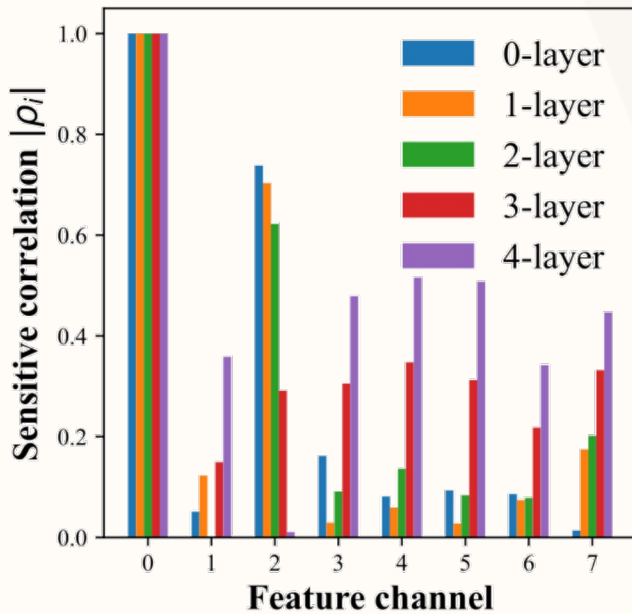


# FairVGNN – Fair View Graph Neural Network

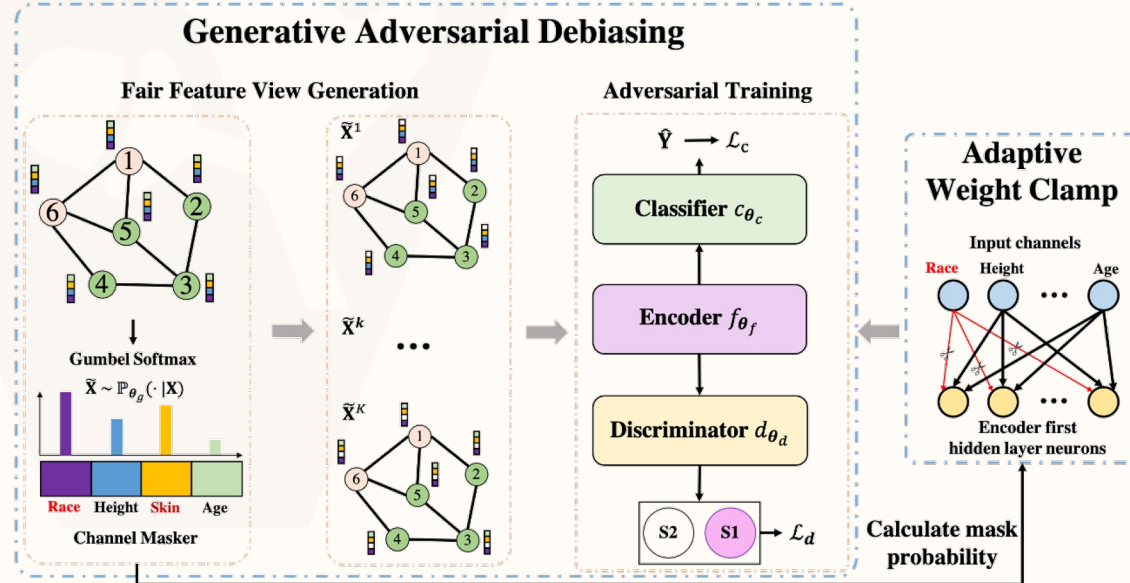
What we desire



Graph Domain Challenge

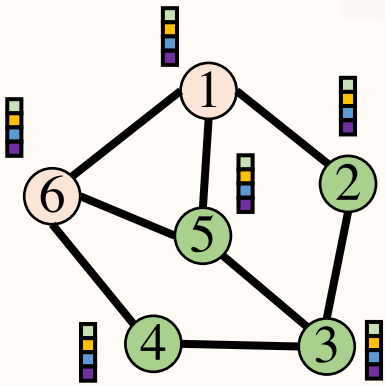


Solution: FairVGNN

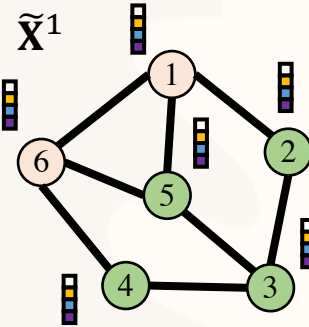


# Fair VGNN – Fair View Graph Neural Network

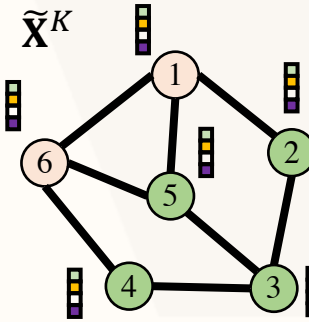
## Fair Feature View Generation



Gumbel Softmax



$\tilde{\mathbf{X}}^k$   
...

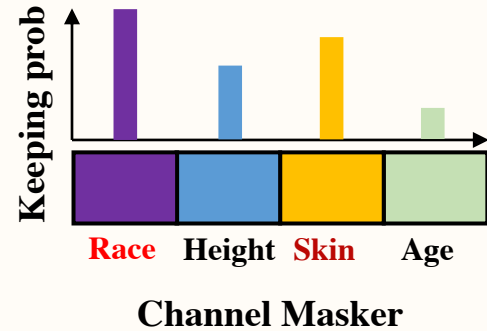


$$\tilde{\mathbf{X}}^k = \mathbf{X} \odot \mathbf{m}^k = [\mathbf{X}_1^T \odot \mathbf{m}, \dots, \mathbf{X}_n^T \odot \mathbf{m}]$$

$\mathbf{m}_i, \mathbf{m}_i$  are independent

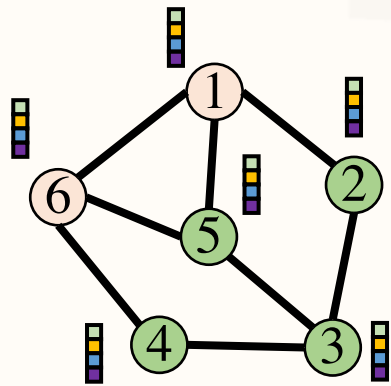
$$\mathbf{m}_i^k \sim \text{Bernoulli}(1 - p_i), \forall i \in \{1, 2, \dots, d\}$$

Gumbel Softmax

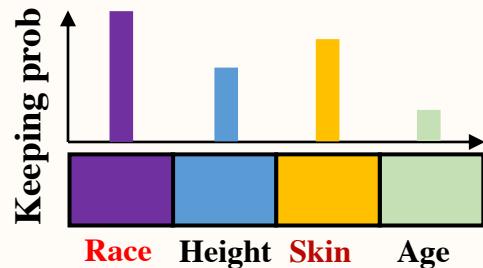


# Fair VGNN – Fair View Graph Neural Network

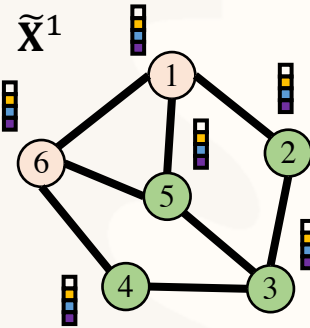
## Fair Feature View Generation



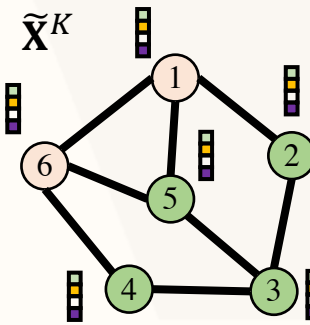
Gumbel Softmax



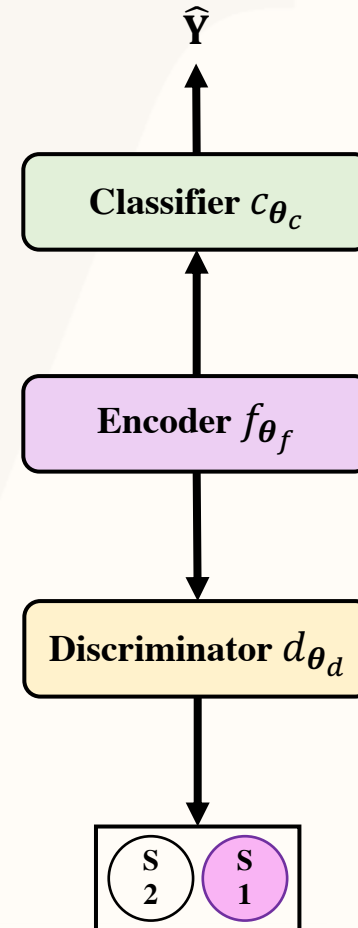
Channel Masker



$\tilde{X}^k$   
...

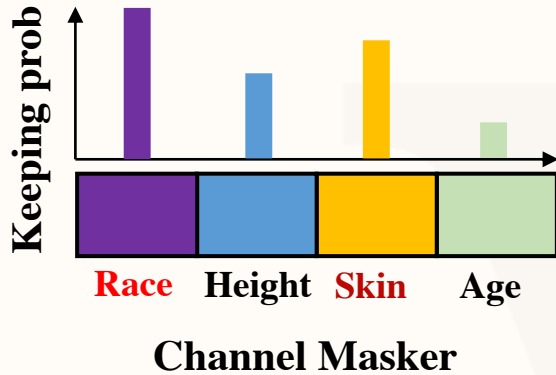


## Adversarial debiasing and classification





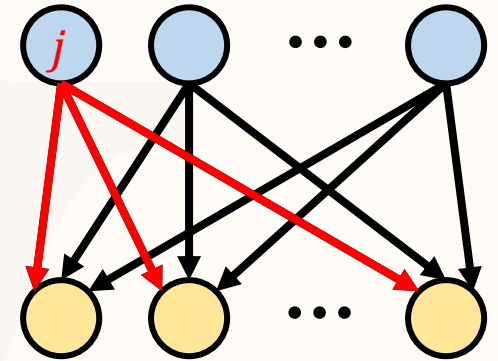
# Fair VGNN – Fair View Graph Neural Network



$$\mathbf{p} = \sum_{k=1}^K \mathbf{m}^k \in \mathbb{R}^d$$

Input channels

Encoder first hidden layer neurons



$$\mathbf{W}_{ij}^{f,1} = \begin{cases} \mathbf{W}_{ij}^{f,1}, & |\mathbf{W}_{ij}^{f,1}| \leq \epsilon * \mathbf{p}_j \\ \text{sign}(\mathbf{W}_{ij}^{f,1}) * \epsilon * \mathbf{p}_j, & |\mathbf{W}_{ij}^{f,1}| > \epsilon * \mathbf{p}_j \end{cases}$$

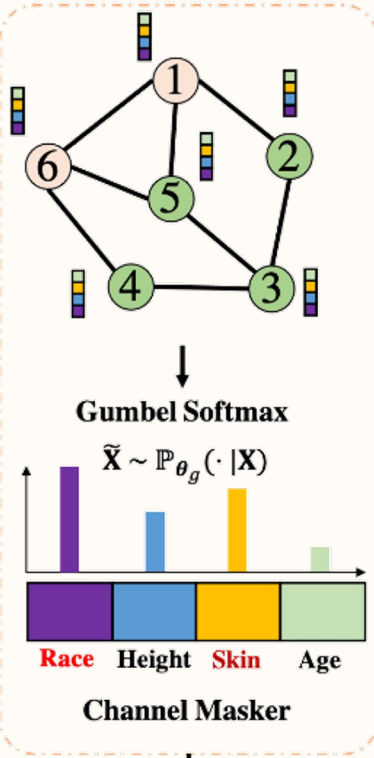
Adaptive Weight Clamp

$$\|\boldsymbol{\mu}\|_2 = \|(2\chi - 1)\mathbf{W}^{f,1}\Delta\boldsymbol{\mu}\|_2 \leq (2\chi - 1) \left( \sum_{i=1}^{d_1} \left( \sum_{r \in \mathcal{S}} \epsilon \mathbf{p}_r \Delta\boldsymbol{\mu}_r + \sum_{k \in \mathcal{N}\mathcal{S}} \epsilon \mathbf{p}_k \Delta\boldsymbol{\mu}_k \right)^2 \right)^{0.5}$$

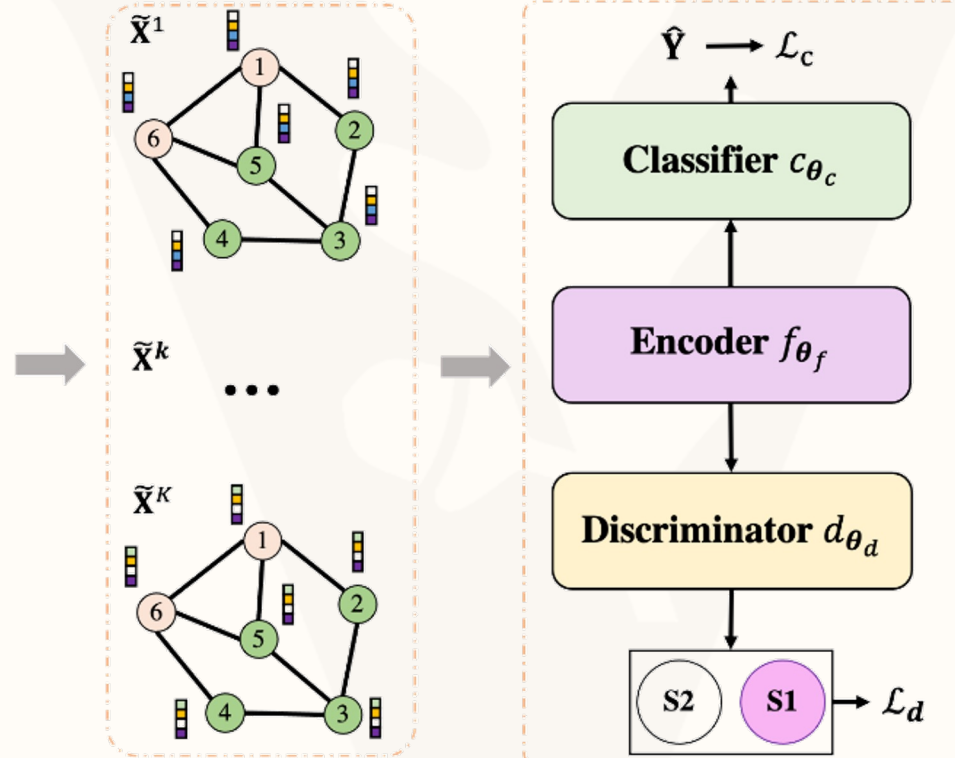
# Fair VGNN – Fair View Graph Neural Network

## Generative Adversarial Debiasing

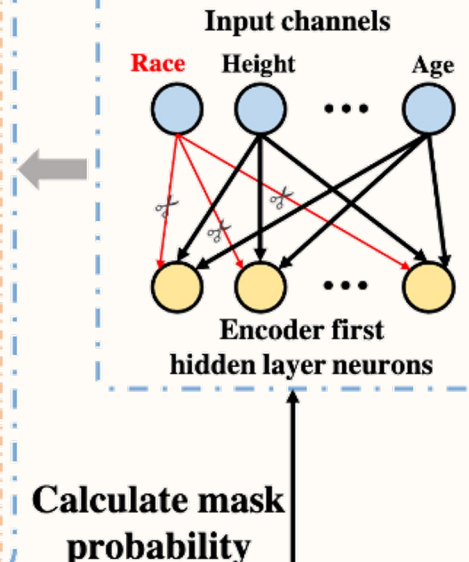
### Fair Feature View Generation



### Adversarial Training



### Adaptive Weight Clamp



# FairVGNN – Datasets and Baselines

**Table 2: Basic dataset statistics.**

<b>Dataset</b>	<b>German</b>	<b>Credit</b>	<b>Bail</b>
#Nodes	1000	30,000	18,876
#Edges	22,242	1,436,858	321,308
#Features	27	13	18
Sens.	Gender	Age	Race
Label	Good/bad Credit	Default/no default Payment	Bail/no bail

**Augmentation-based:** NIFTY, EDITS

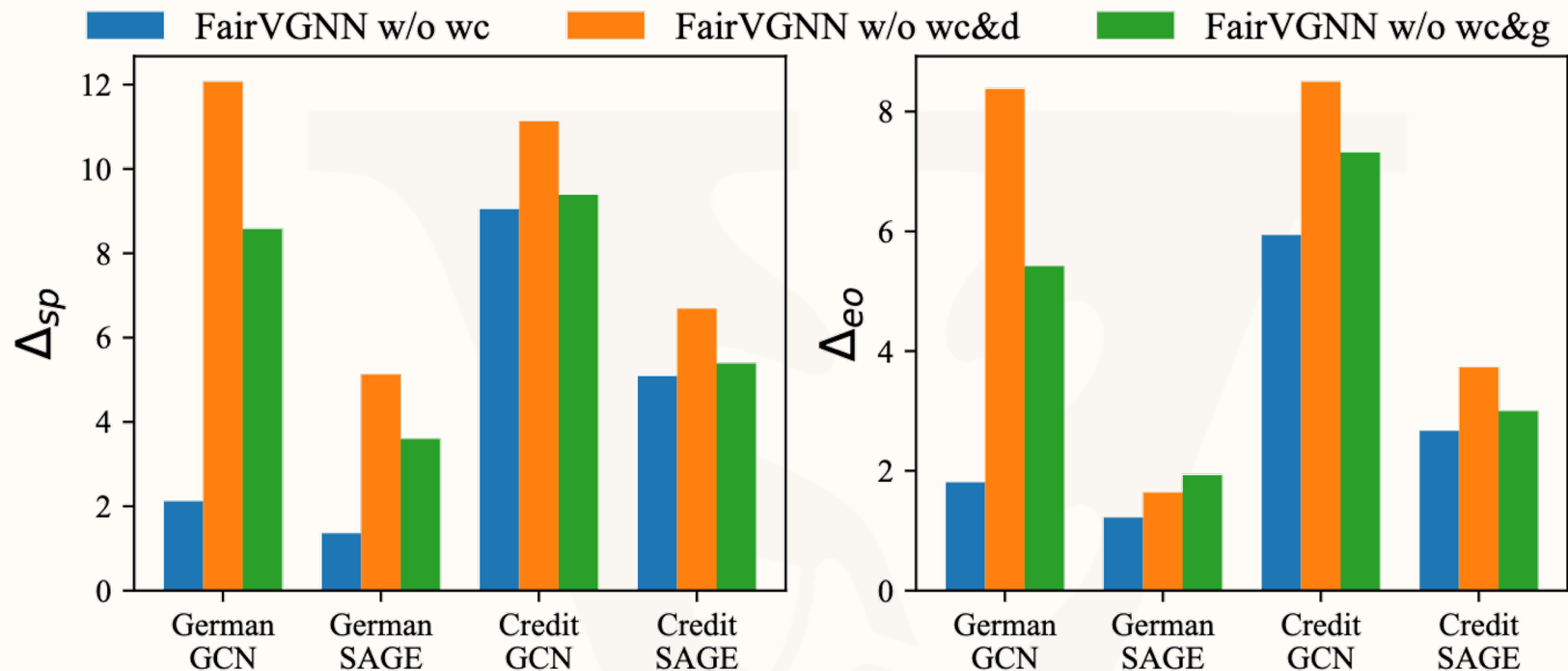
**Adversarial-based:** FairGNN

# FairVGNN – Experiments

Encoder	Method	German				
		AUC (↑)	F1 (↑)	ACC (↑)	$\Delta_{sp}$ (↓)	$\Delta_{eo}$ (↓)
GCN	Vanilla	74.11±0.37	82.46±0.89	73.44±1.09	35.17±7.27	25.17±5.89
	NIFTY	68.78±2.69	81.40±0.54	69.92±1.14	5.73±5.25	5.08±4.29
	EDITS	69.41±2.33	81.55±0.59	71.60±0.89	4.05±4.48	3.89±4.23
	FairGNN	67.35±2.13	82.01±0.26	69.68±0.30	3.49±2.15	3.40±2.15
	FairVGNN	72.41±2.10	82.14±0.42	70.16±0.86	1.71±1.68	0.88±0.58
Encoder	Method	Credit				
		AUC (↑)	F1 (↑)	ACC (↑)	$\Delta_{sp}$ (↓)	$\Delta_{eo}$ (↓)
GIN	Vanilla	74.36±0.21	82.28±0.64	74.02±0.73	14.48±2.44	12.35±2.86
	NIFTY	70.90±0.24	84.05±0.82	75.59±0.66	7.09±4.62	6.22±3.26
	EDITS	72.35±1.11	82.47±0.85	74.07±0.98	14.11±14.45	15.40±15.76
	FairGNN	68.66±4.48	79.47±5.29	70.33±5.50	4.67±3.06	3.94±1.49
	FairVGNN	71.36±0.72	87.44±0.23	78.18±0.20	2.85±2.01	1.72±1.80

- (1) Compared with vanilla GNN model, the bias-mitigating model can achieve lower bias
- (2) Compared with other baselines, FairVGNN can achieve even better trade-off between fairness and utility performance

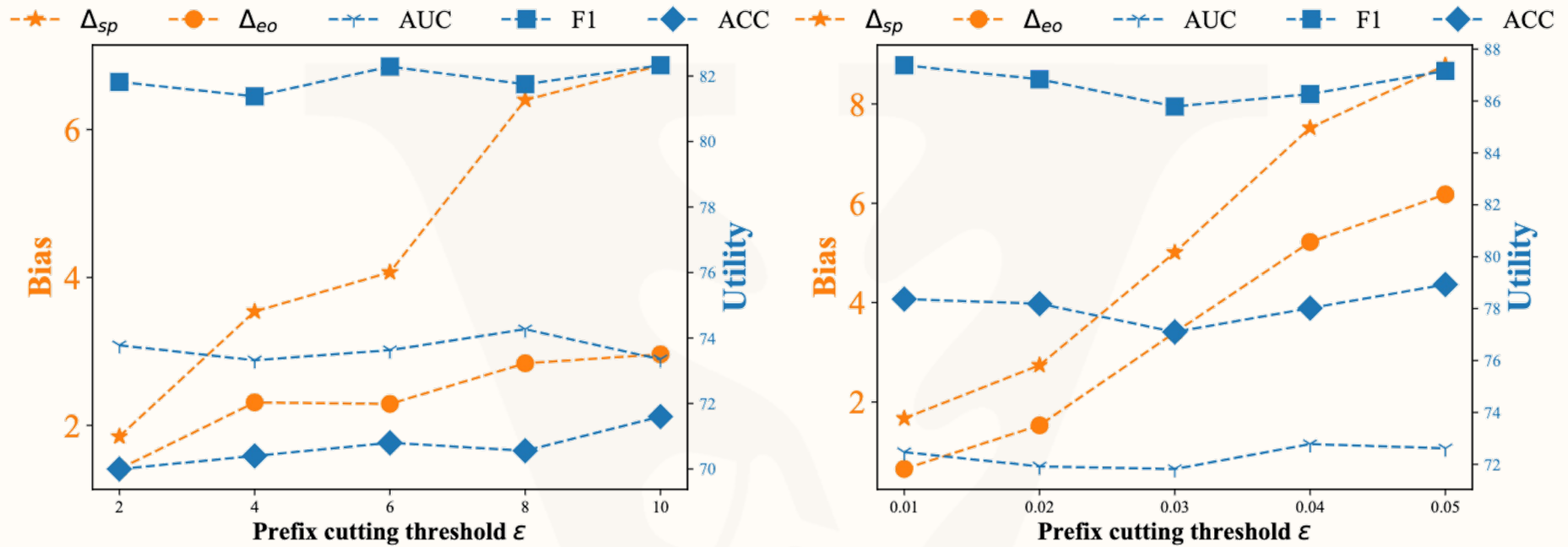
# FairVGNN – Adversarial training



(1) Removing either the generator or the discriminator would lower the fairness

(2) Removing the discriminator causes the highest bias

# Fair VGNN – Adaptive weight clamping

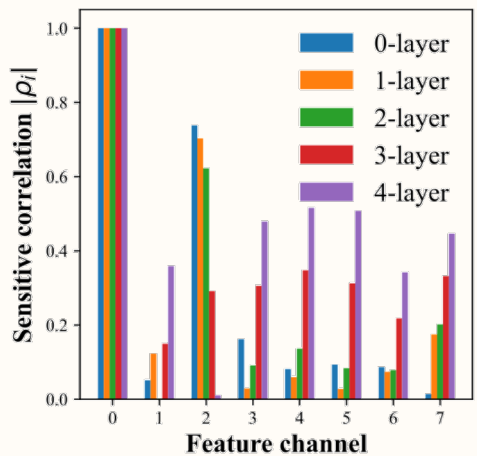
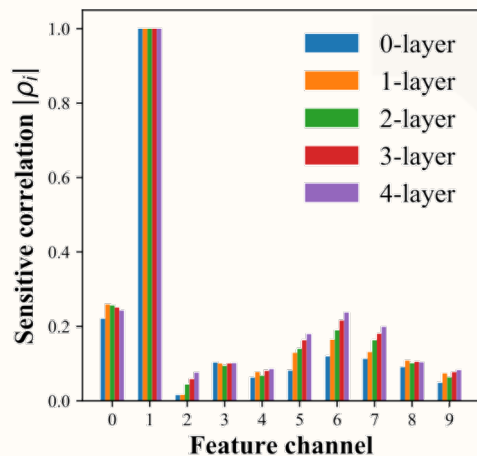


(a) German

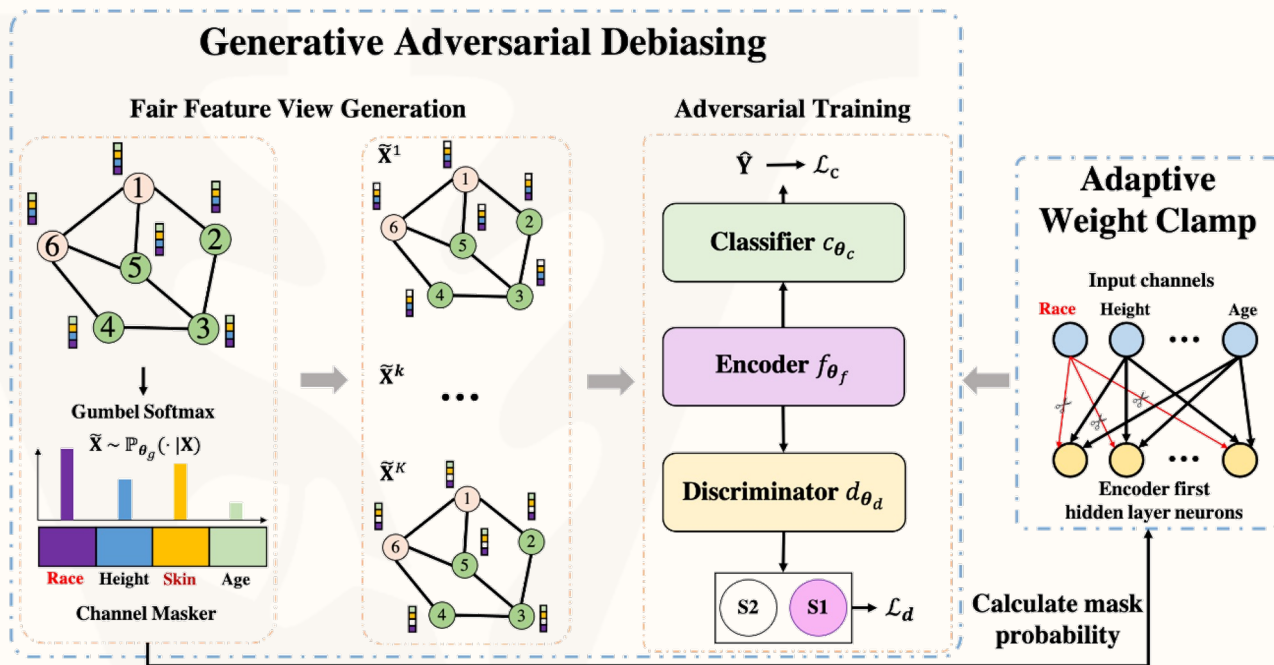
(b) Credit

# Contributions

## Novel problem



## Solution: FairVGNN

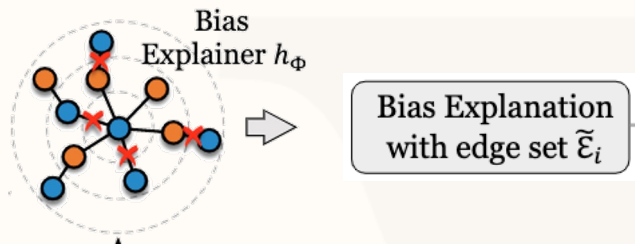


## New Finding: bias and homophily

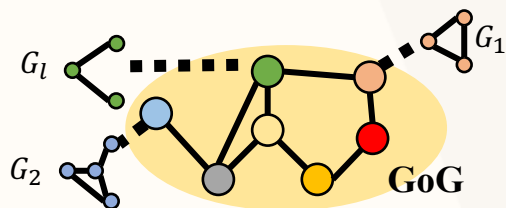
$$\|\mu\|_2 = \|(2\chi - 1)W^{f,1}\Delta\mu\|_2 \leq (2\chi - 1) \left( \sum_{i=1}^{d_1} \left( \sum_{r \in S} \epsilon_{pr} \Delta\mu_r + \sum_{k \in NS} \epsilon_{pk} \Delta\mu_k \right)^2 \right)^{0.5}$$

# Concurrent and Future work

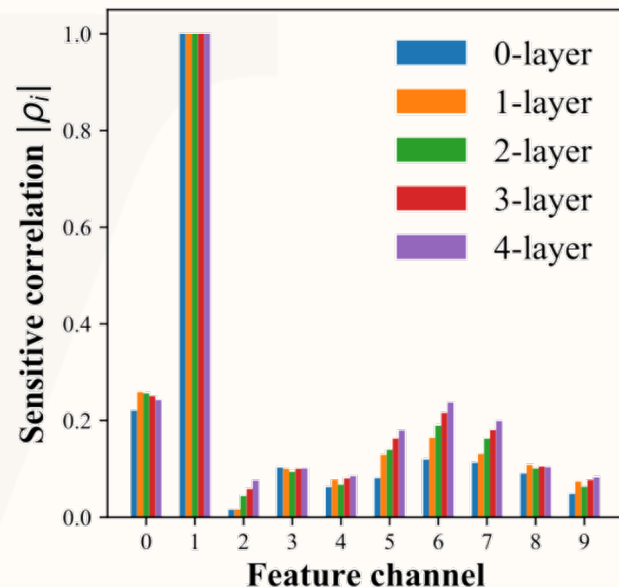
## Structural explanation for bias (KDD 22)



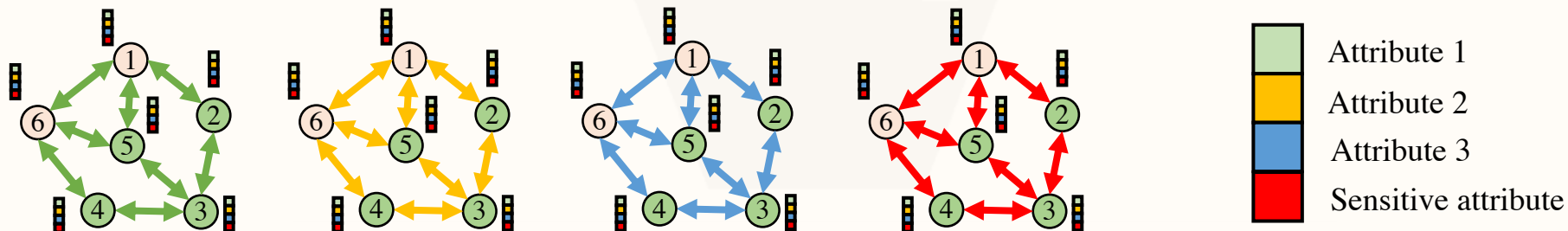
## Mitigating Labeling bias (CIKM 22)



## Correlation variation

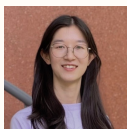


## Network channel homophily, propagation and fairness





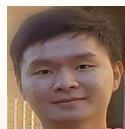
# Acknowledgement



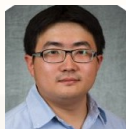
Yuying  
Zhao



Yushun  
Dong



Huiyuan  
Chen



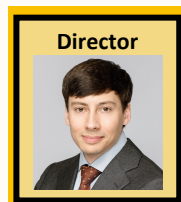
Jundong  
Li



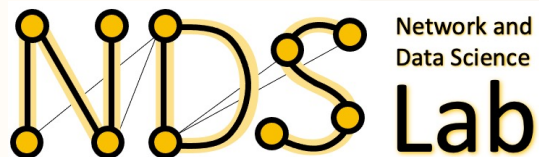
Tyler  
Derr



Association for  
Computing Machinery

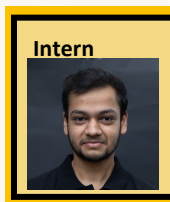


Director

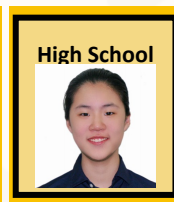


Network and  
Data Science

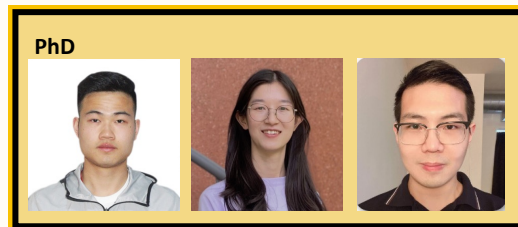
Lab



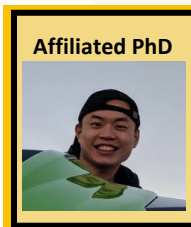
Intern



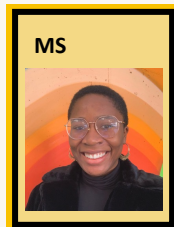
High School



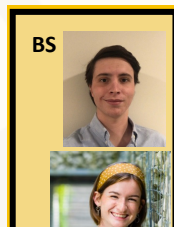
PhD



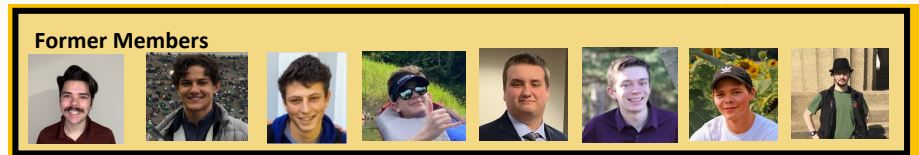
Affiliated PhD



MS



BS



Former Members



More about me

<https://yuwvandy.github.io/>

<https://nds-vu.github.io/>