



Fairness and Explainability: Bridging the Gap Towards Fair Model Explanations



Yuying Zhao



Yu Wang

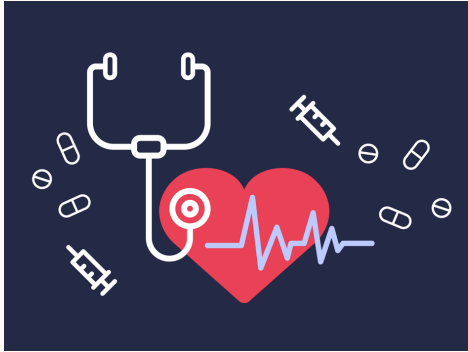


Tyler Derr

yuying.zhao@vanderbilt.edu
<https://yuyingzhao.github.io/>

Background: Bias and Fairness in ML

Healthcare



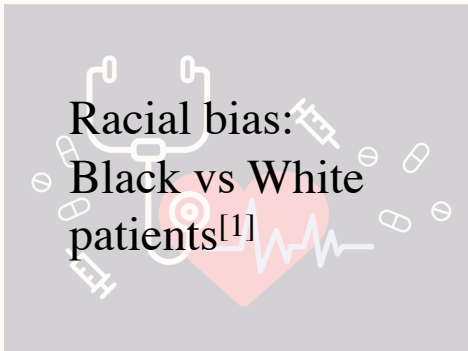
Education



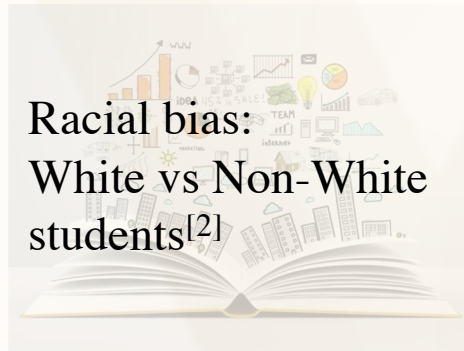
Finance



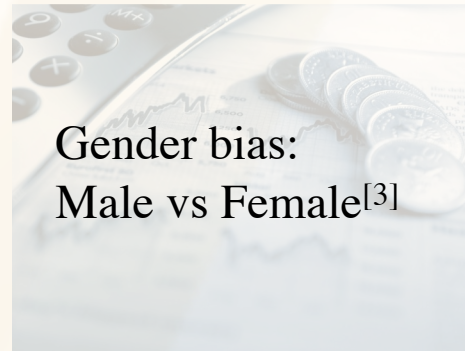
Healthcare



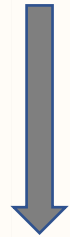
Education



Finance



Utility
Performance



Fairness
Concerns

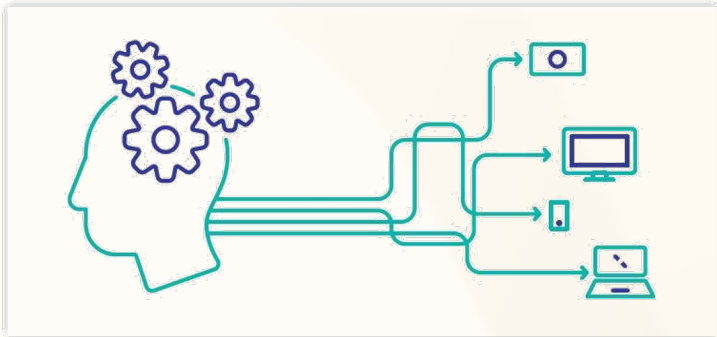
[1] Obermeyer, Ziad, et al. "Dissecting racial bias in an algorithm used to manage the health of populations." *Science*. 2019.

[2] Anderson, Henry, et al. "Assessing the Fairness of Graduation Predictions." *EDM*. 2019.

[3] Zhang, Yukun, et al. "Fairness assessment for artificial intelligence in financial industry." *NeurIPS*. 2019.

Background: Model Explainability

Explainability



- Why should I trust the model?
- Why did a model make a certain decision?

Healthcare



Education



Business perspective:

- Trust before deployment
- Find justification

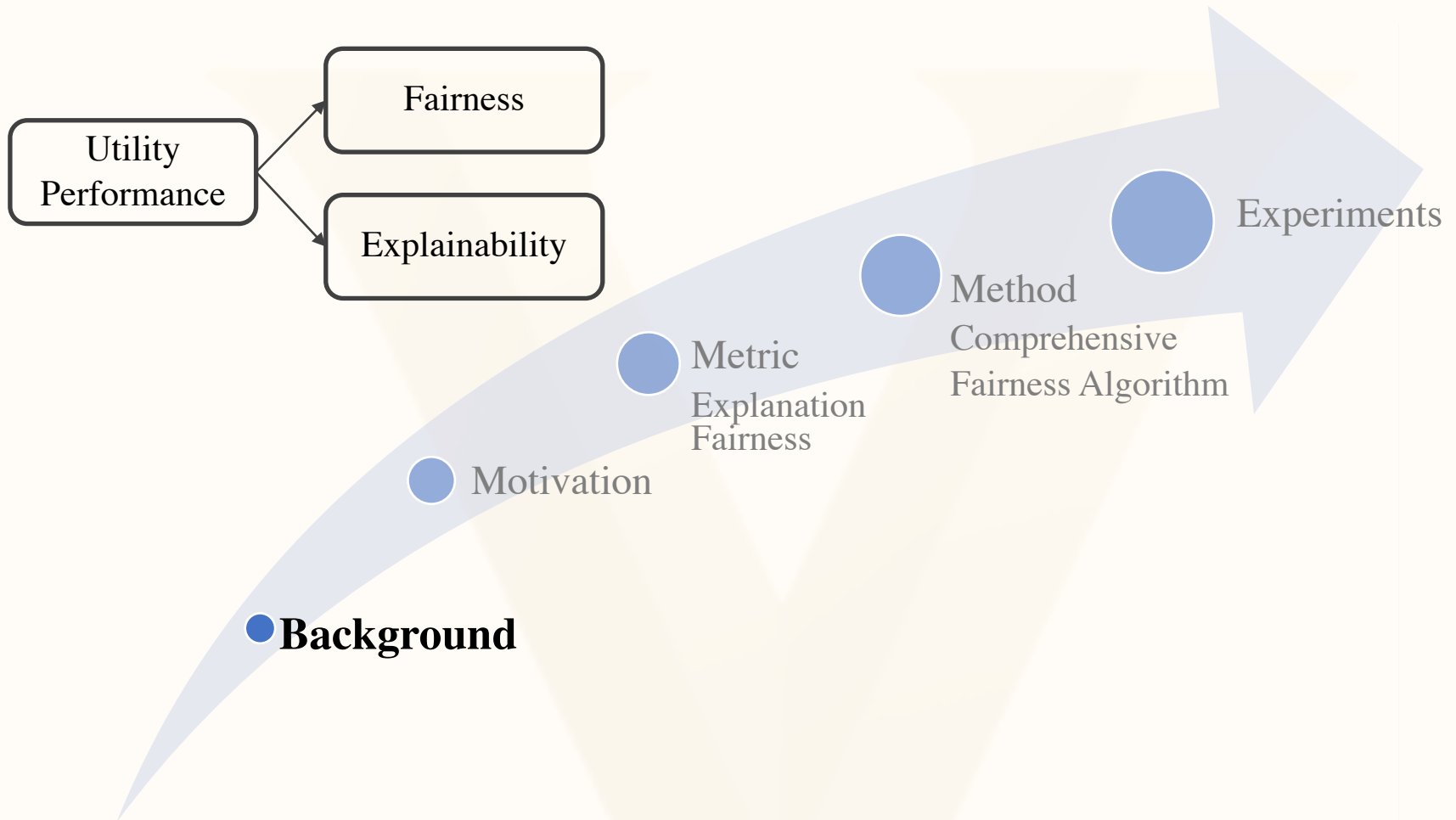
Model perspective:

- Debug model (mis)predictions
- Improve/verify ML models

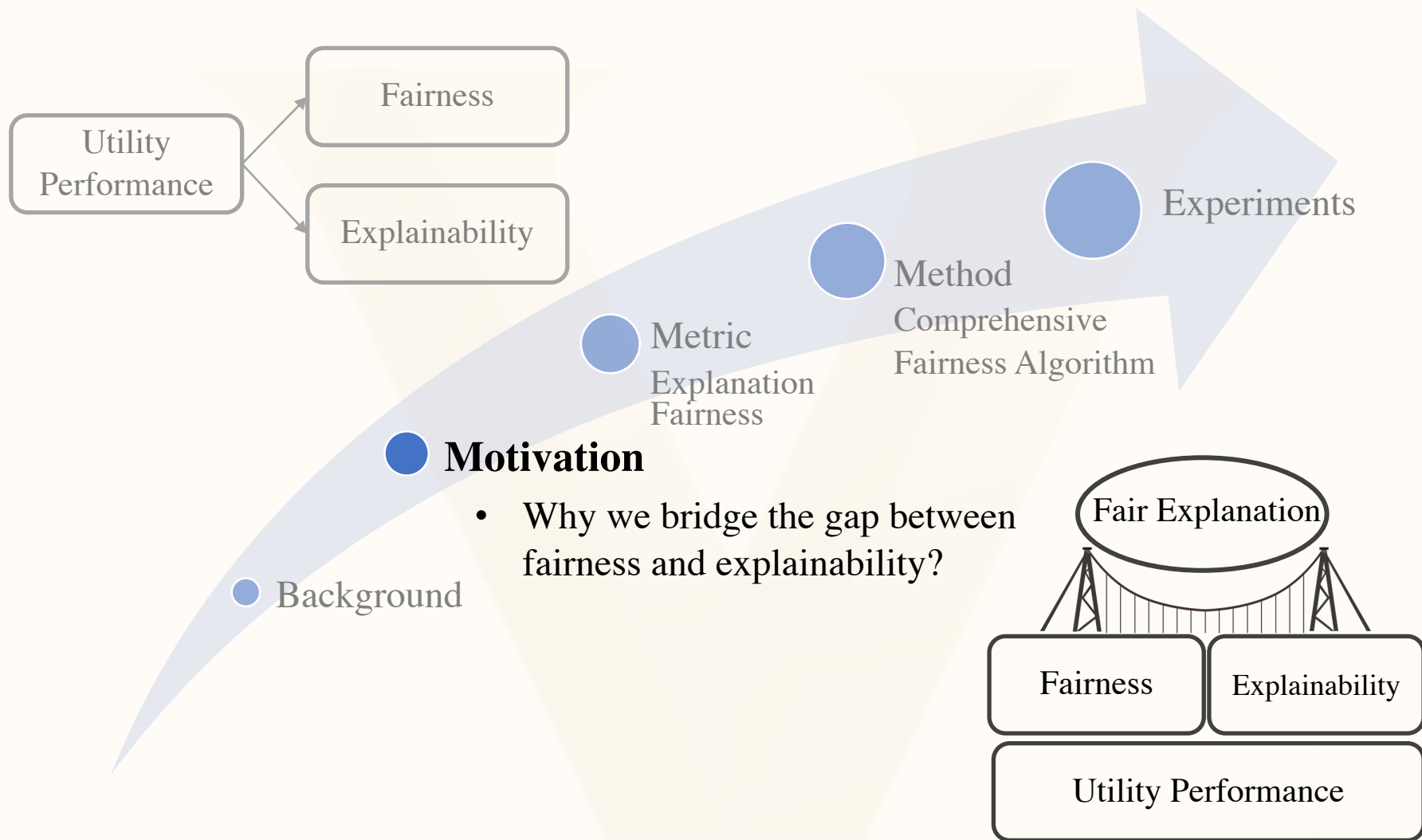
Regulatory perspective:

- GDPR: Article 22 empowers individuals with the right to demand an explanation of how an automated system made a decision that affects them.

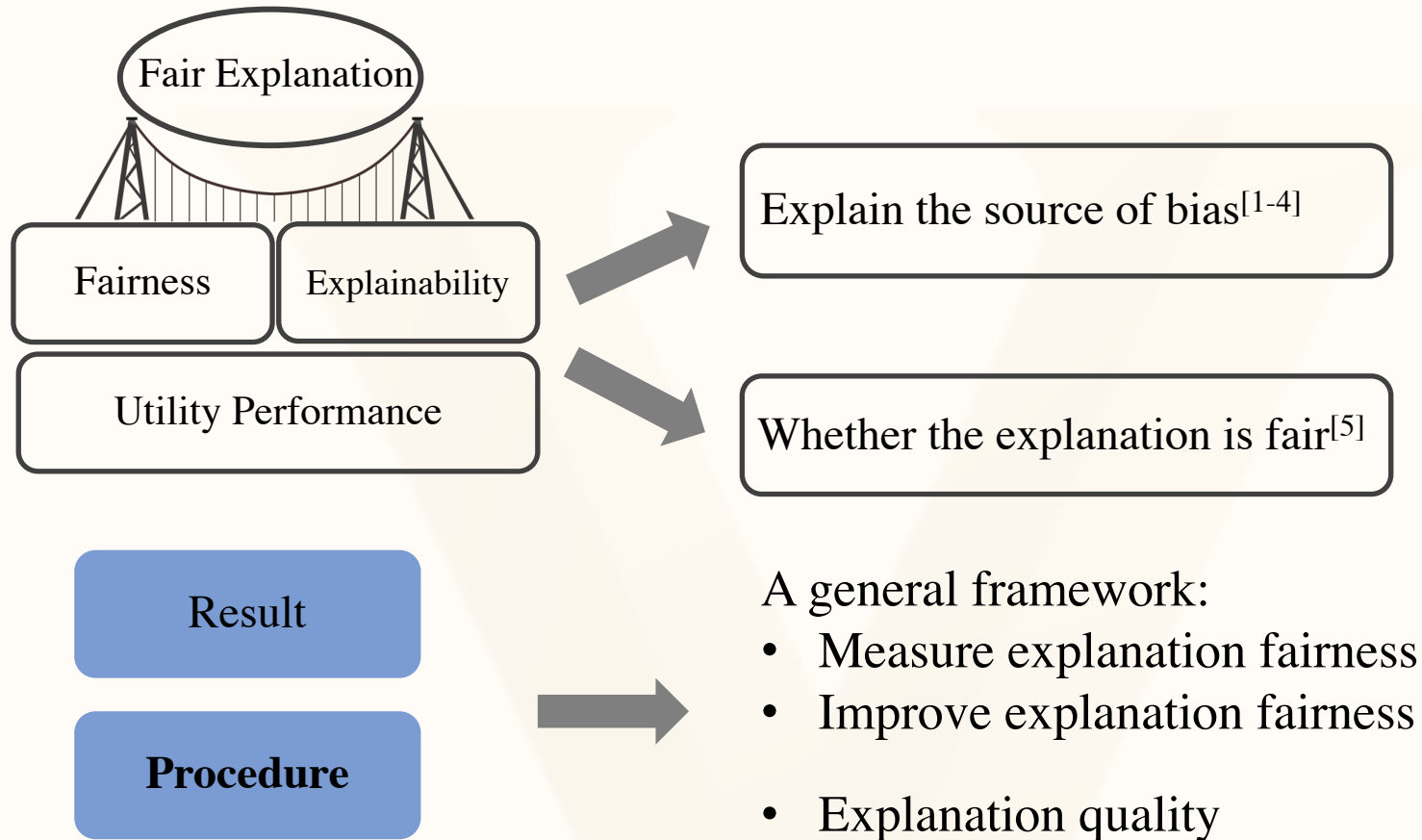
Content



Content



Existing Work



[1] Lundberg, S. M. "Explaining Quantitative Measures of Fairness." Fair & Responsible AI Workshop. 2020.

[2] Begley, Tom, et al. "Explainability for fair machine learning". arXiv. 2020.

[3] Chiappa, S. "Path-specific counterfactual fairness." AAAI. 2019.

[4] Pan, Weishen, et al. "Explaining algorithmic fairness through fairness-aware causal path decomposition". KDD. 2021.

[5] Fu, Zuohui, et al. "Fairness-aware explainable recommendation over knowledge graphs." SIGIR. 2020.

Motivation: Fairness and Explainability



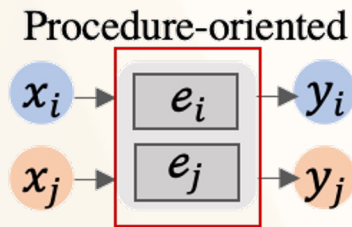
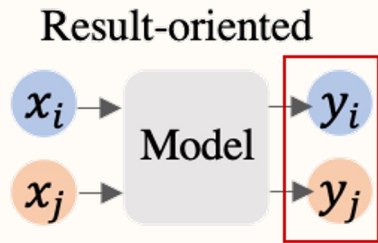
Motivation:

most fairness metrics: result-oriented
hide the potential bias during the procedure

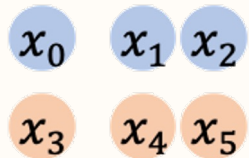
\hat{y} : predictions
 y : ground truth
 s : sensitive features

Statistical Parity:

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

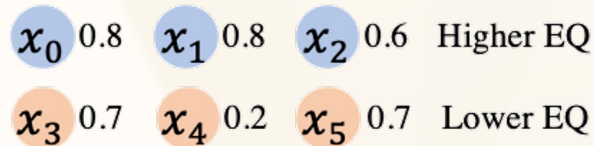


Fair Predictions



Hired Not

Unfair Explanations



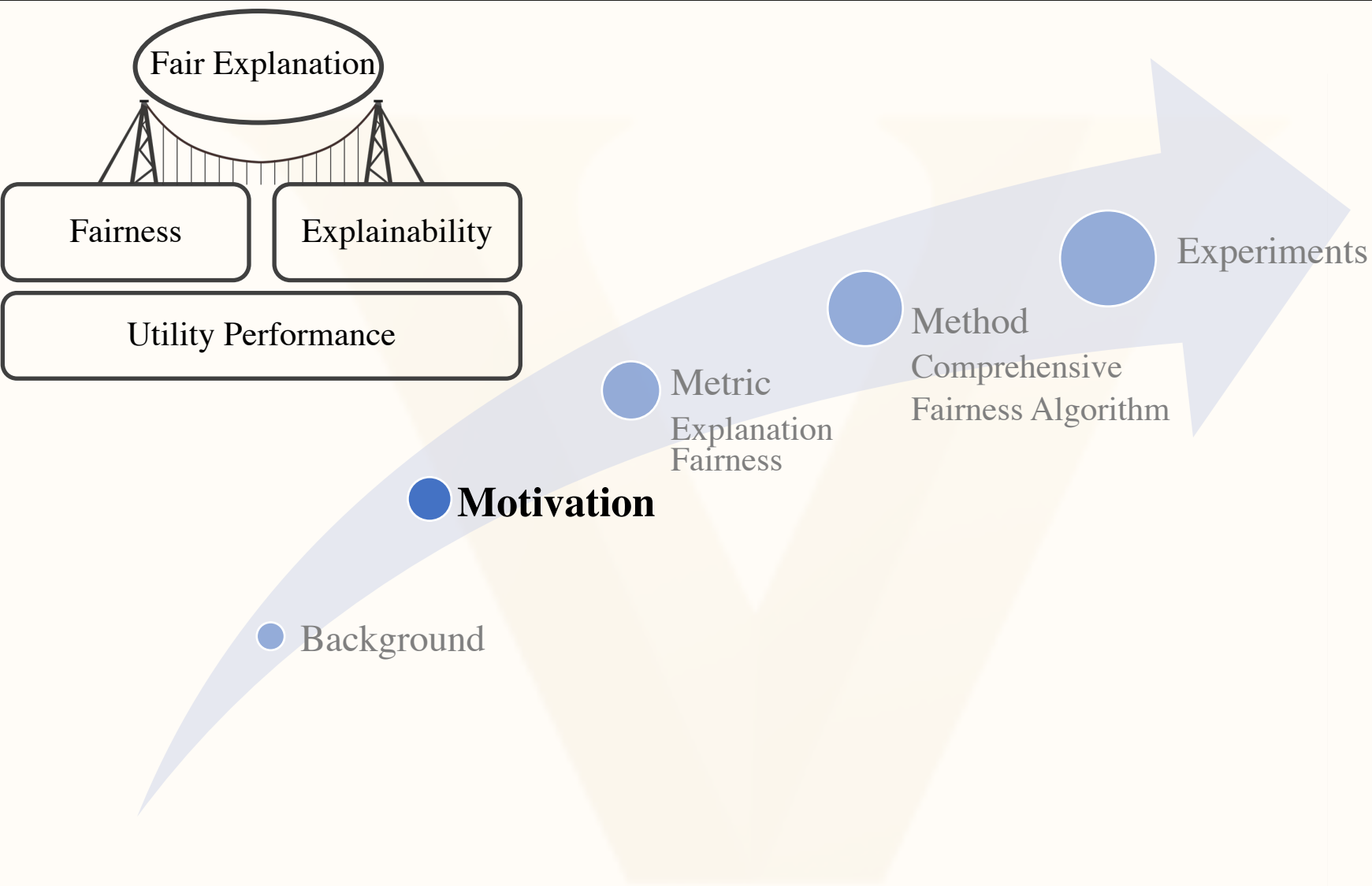
Explanation Quality (EQ)

Unfairness:

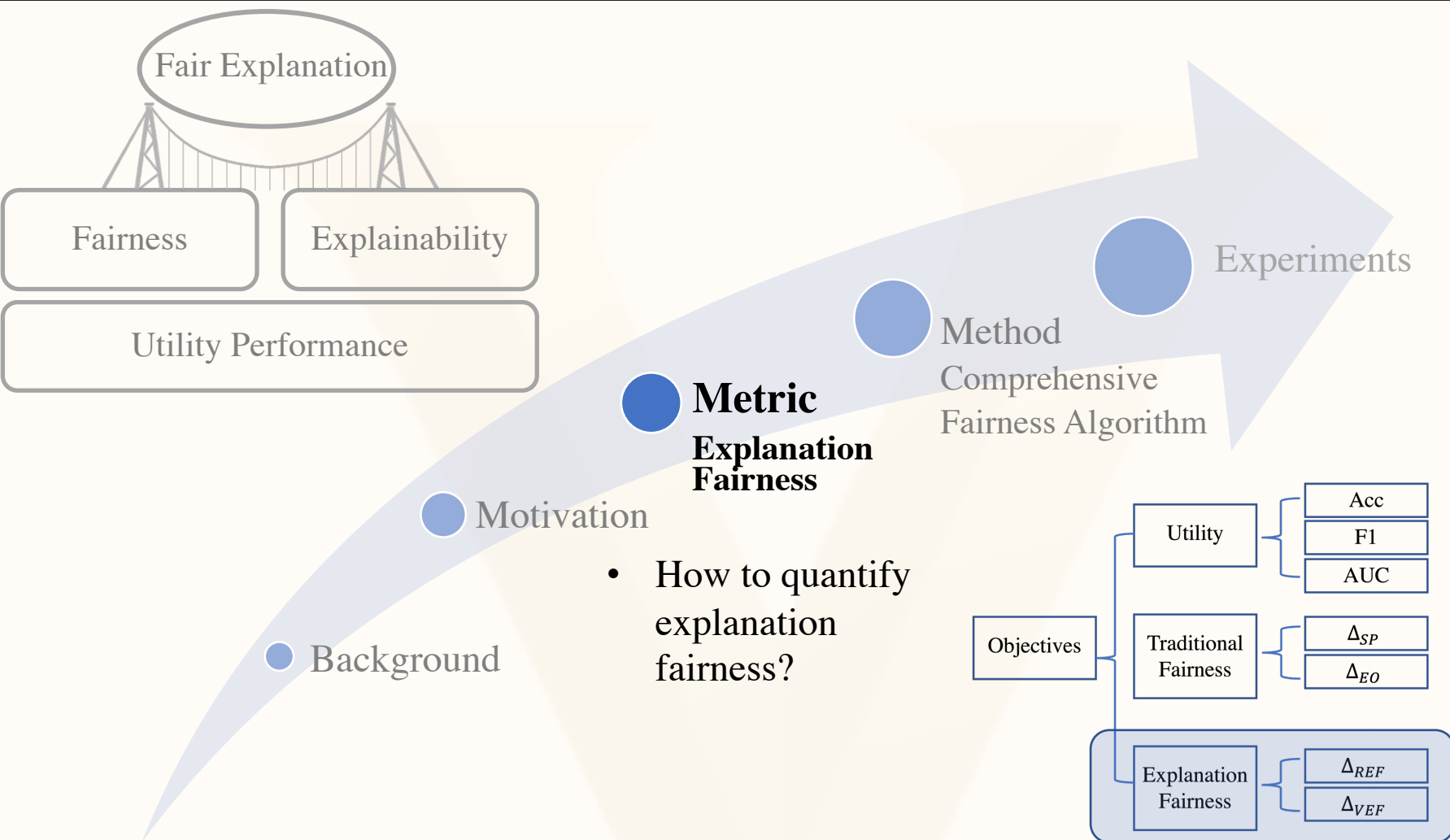
Better explanation for one group than the other

Example: Job hiring
Well-explained vs
Ambiguous explanation

Content

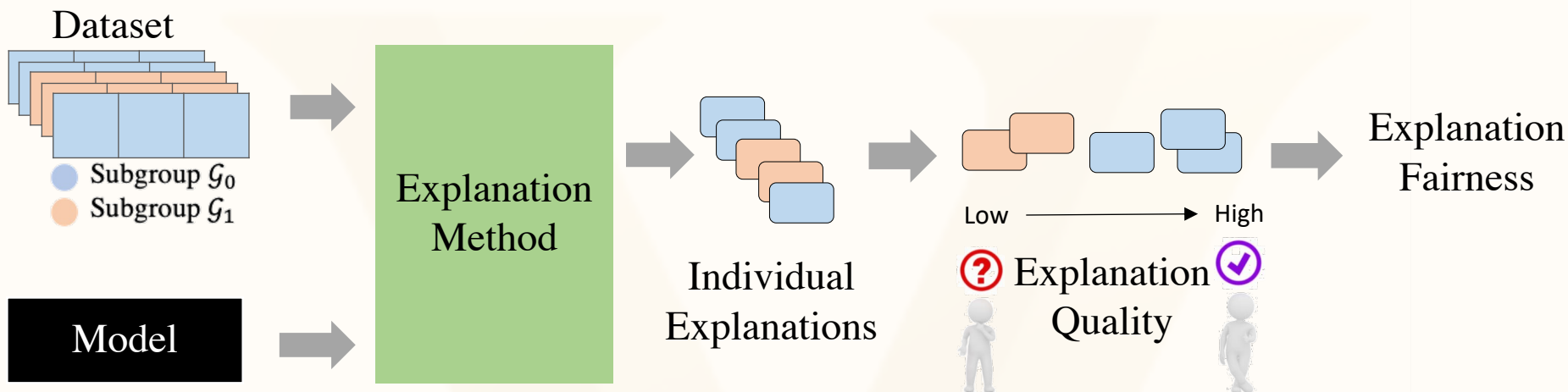


Content



Metric: High-level Idea of Explanation Fairness

Compare explanation quality from two subgroups

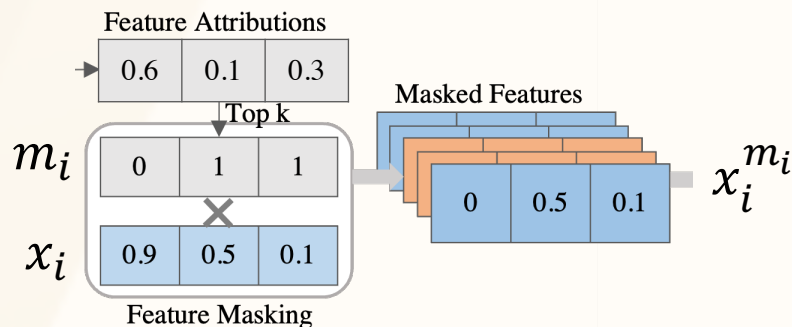


What makes a “good” explanation?

Fidelity:

$$\Delta_{F_i} = P(\hat{y}_i = y_i | x = x_i) - P(\hat{y}_i = y_i | x = x_i^{m_i})$$

How well does the explanation approximate the prediction of the black-box model?



Metric: Quantification of Explanation Fairness

Given explanation quality (EQ),
how to quantify explanation fairness?

(1) Ratio-based Fairness Δ_{REF}

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

Same opportunity of having positive prediction

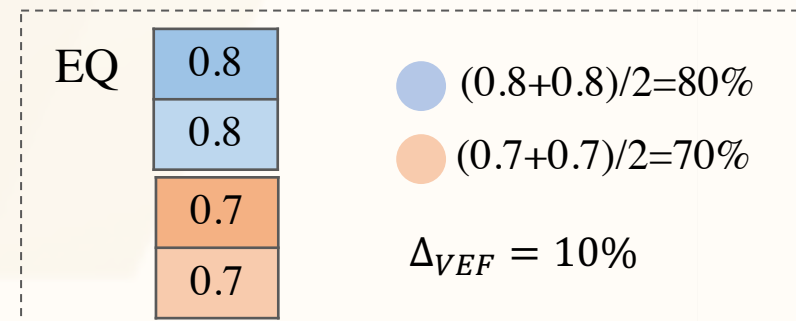
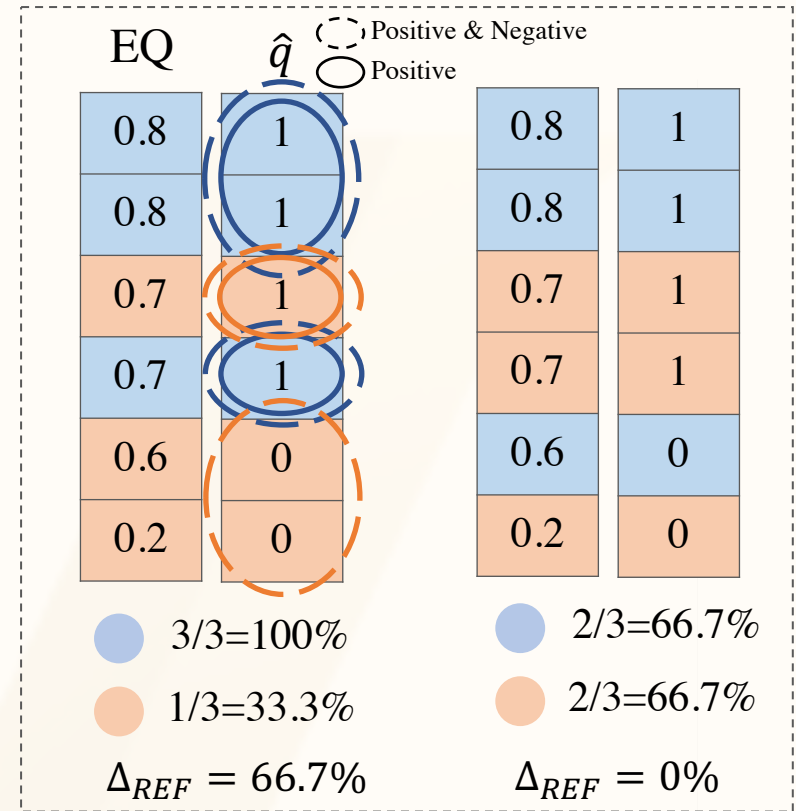
$$\Delta_{REF} = |P(\hat{q} = 1 | s = 0) - P(\hat{q} = 1 | s = 1)|$$

Same opportunity of having high-quality explanations

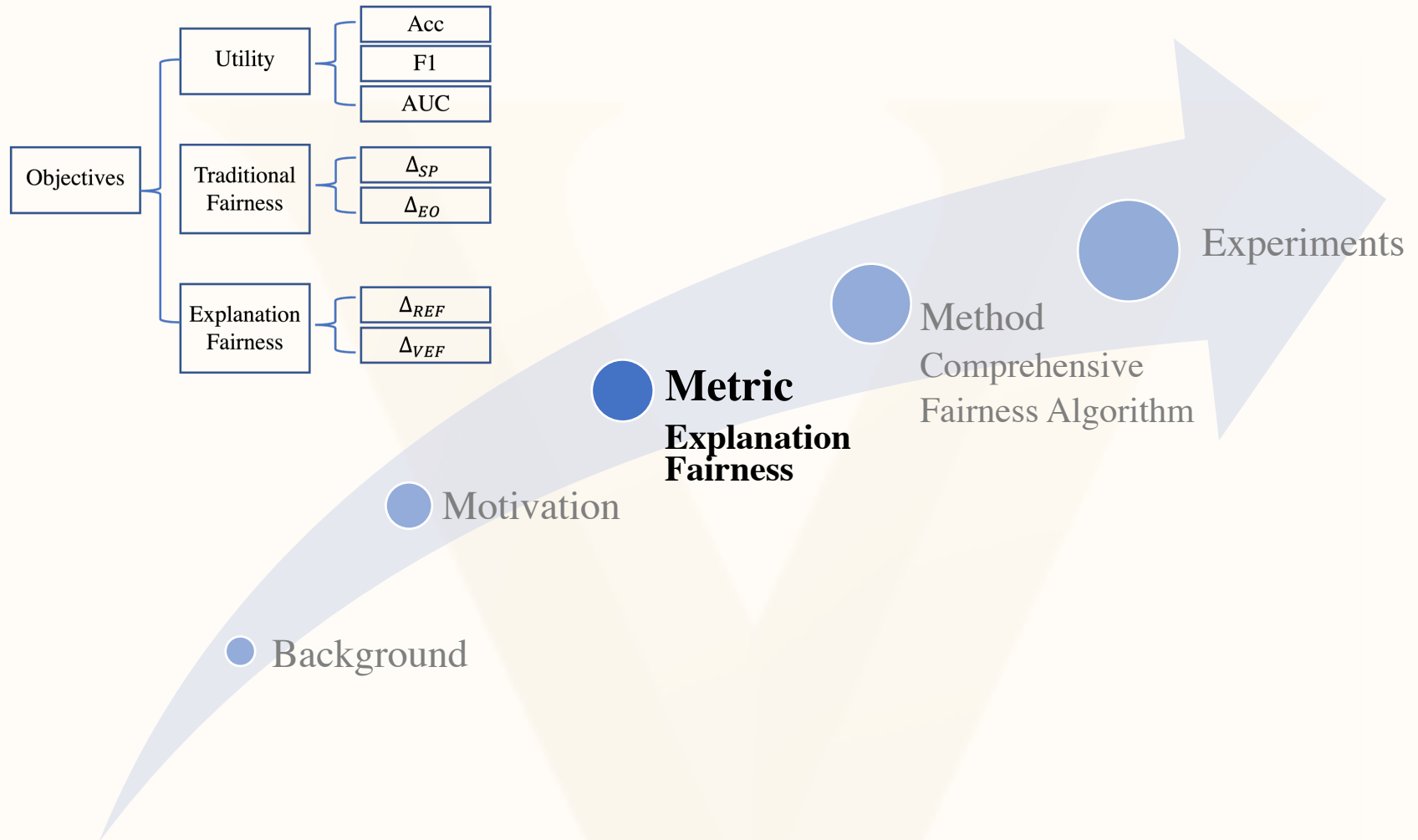
\hat{y} : prediction \hat{q} : explanation quality

(2) Value-based Fairness Δ_{VEF}

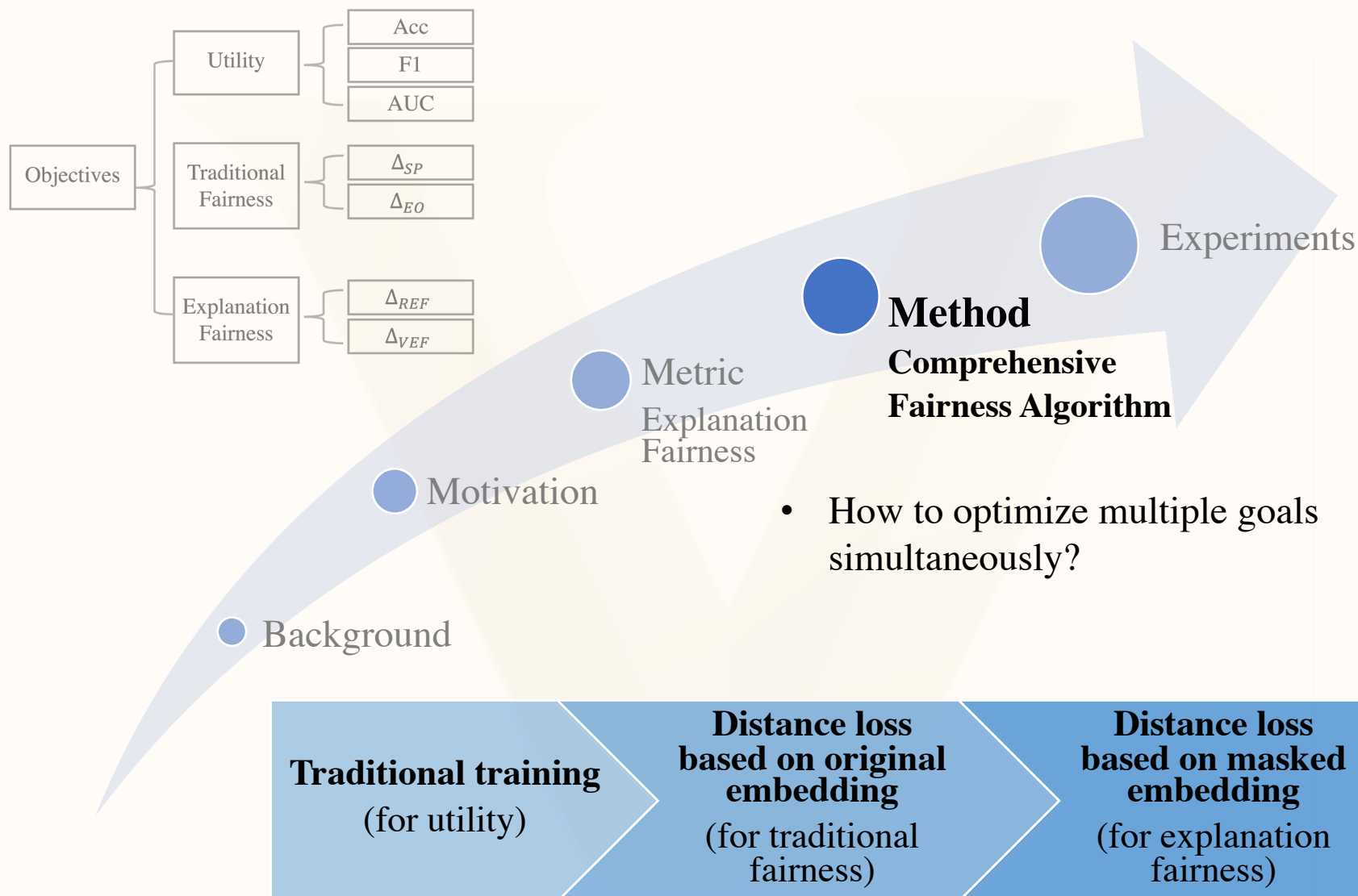
$$\Delta_{VEF} = \left| \frac{1}{|\mathcal{G}_0^K|} \sum_{i \in \mathcal{G}_0^K} EQ_i - \frac{1}{|\mathcal{G}_1^K|} \sum_{i \in \mathcal{G}_1^K} EQ_i \right|$$



Content

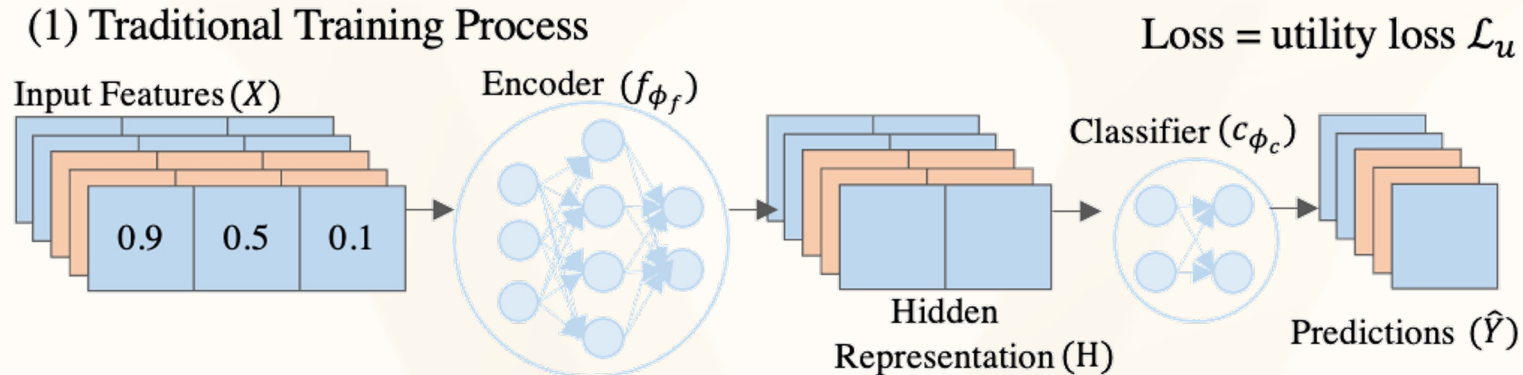


Content



Comprehensive Fairness Algorithm (CFA)

● Subgroup \mathcal{G}_0
● Subgroup \mathcal{G}_1

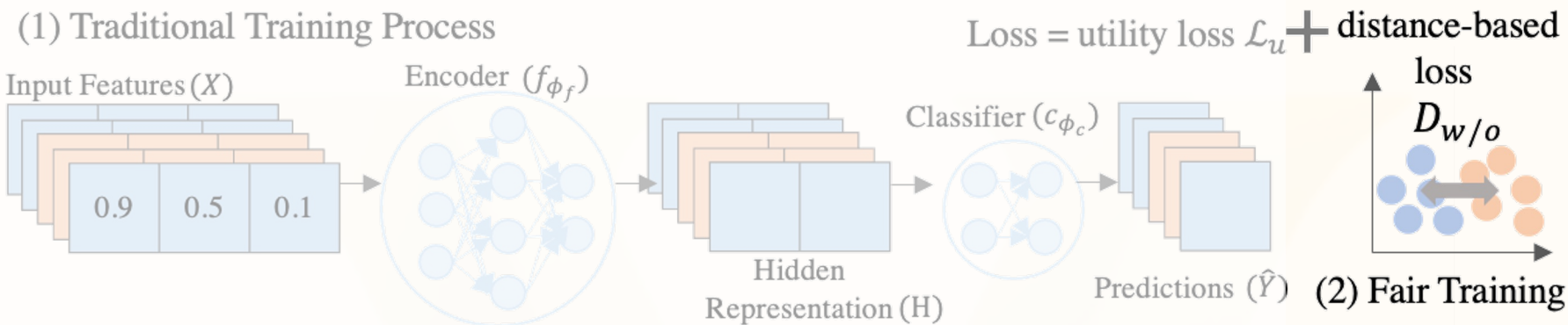


Utility loss: entropy loss (for binary classification)

$$\mathcal{L}_u = - \sum_{i=1}^{|\mathbf{Y}|} (y_i \log(p) + (1 - y_i) \log(1 - p))$$

CFA: Traditional Fairness Optimization

● Subgroup \mathcal{G}_0
● Subgroup \mathcal{G}_1



$$D_{w/o} = \mathcal{D}(\mathbf{H}_{\mathcal{G}_0}, \mathbf{H}_{\mathcal{G}_1})$$

w/o: without masking
based on the original feature

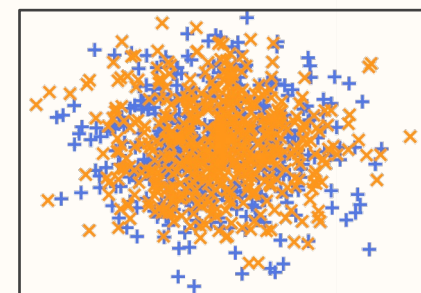
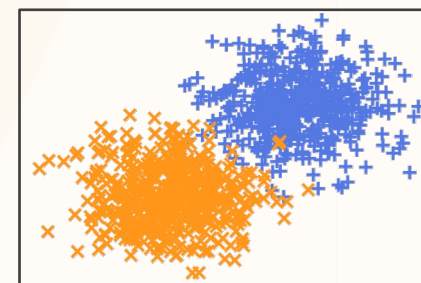
$$\Delta_{SP} = |P(\hat{y} = 1 | s = 0) - P(\hat{y} = 1 | s = 1)|$$

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

The predictions should be irrelevant to sensitive features

Requirements to the hidden representation

- (1) Encode sufficient information for prediction
- (2) Hide information related to sensitive features

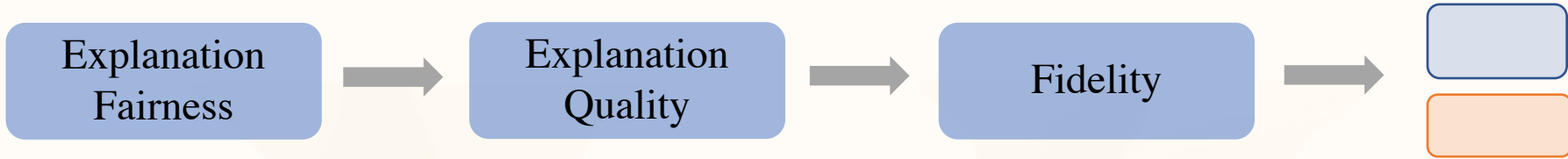


[1]

[1] Dong, Yushun, et al. "Edits: Modeling and mitigating data bias for graph neural networks." WWW. 2022

CFA: Explanation Fairness Optimization

● Subgroup \mathcal{G}_0
 ● Subgroup \mathcal{G}_1

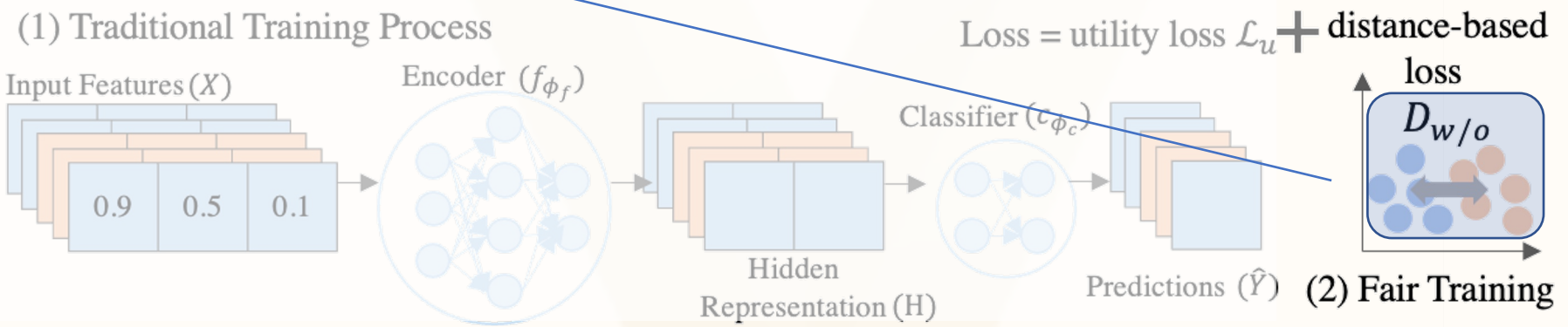


Fair explanation quality (measured by fidelity)

Fidelity: $P(\hat{y}_i = y_i | x = x_i)$ - $P(\hat{y}_i = y_i | x = x_i^{m_i})$

Original feature

Masked feature



CFA: Explanation Fairness Optimization

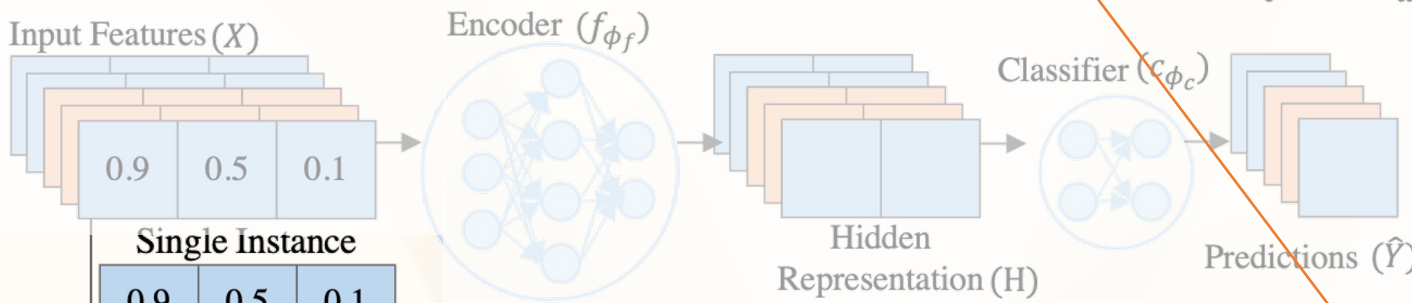
● Subgroup \mathcal{G}_0
● Subgroup \mathcal{G}_1

Fidelity:

$$P(\hat{y}_i = y_i | x = x_i) - P(\hat{y}_i = y_i | x = x_i^{m_i})$$

Masked feature

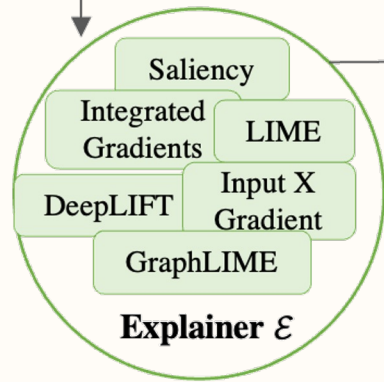
(1) Traditional Training Process



Loss = utility loss \mathcal{L}_u + distance-based



(2) Fair Training



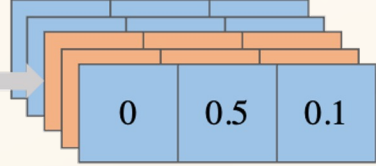
Feature Attributions
0.6 0.1 0.3

Top k

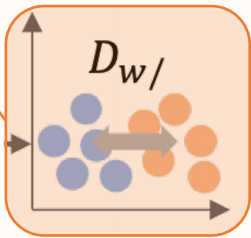
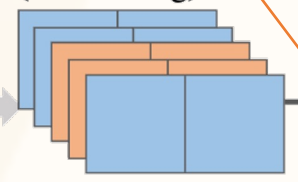
0	1	1
0.9	0.5	0.1

Feature Masking

Masked Features (X^m)



Hidden Representation (H^m) (w/ masking)



(3) Fair Explanation Training

w/: with masking
based on the masked feature

$$\mathcal{L}_u = - \sum_{i=1}^{|\mathbf{Y}|} (y_i \log(p) + (1 - y_i) \log(1 - p))$$

$$D_{w/o} = \mathcal{D}(\mathbf{H}_{\mathcal{G}_0}, \mathbf{H}_{\mathcal{G}_1})$$

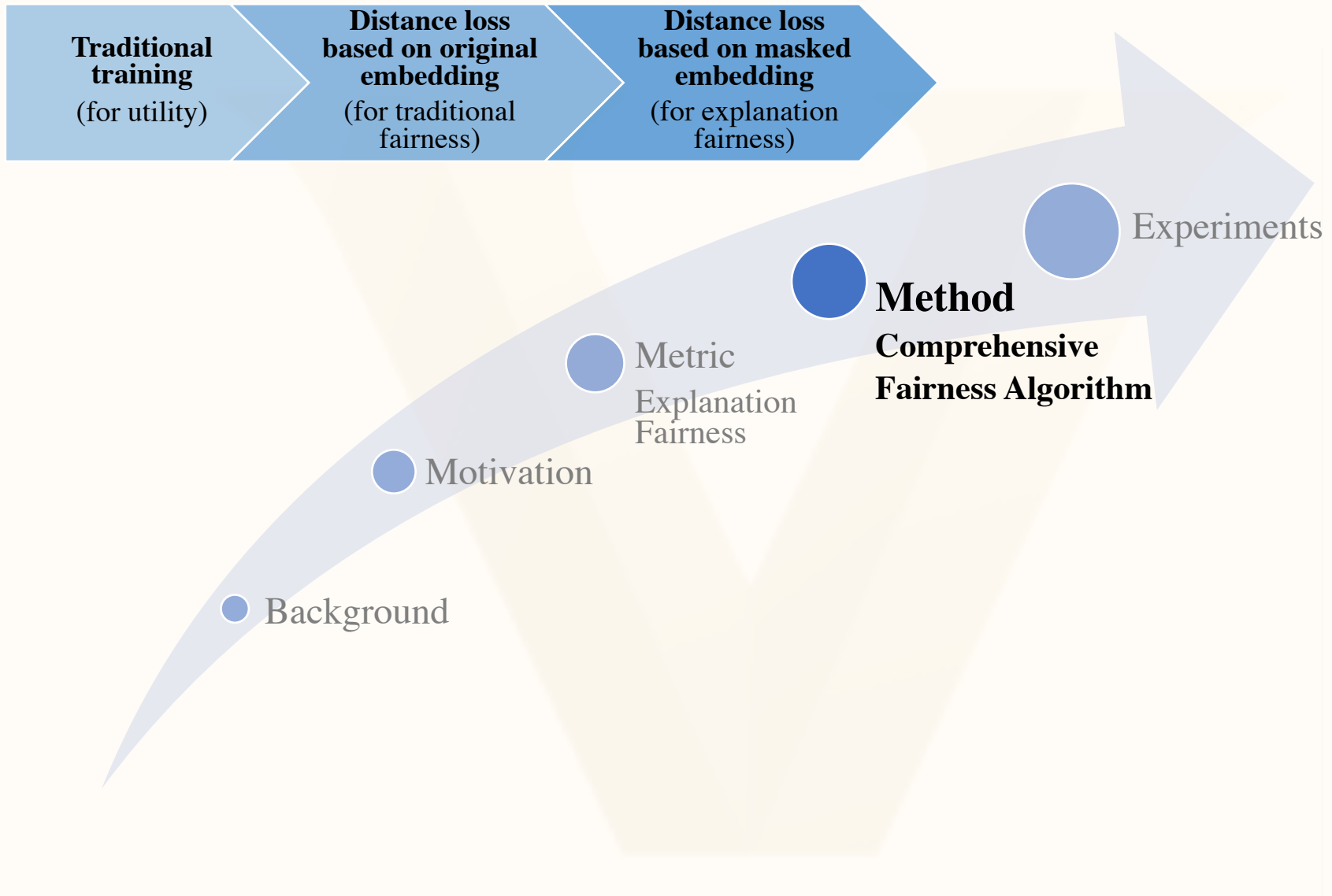
$$D_{w/} = \mathcal{D}(\mathbf{H}_{\mathcal{G}_0}^m, \mathbf{H}_{\mathcal{G}_1}^m)$$

$$\mathcal{L}_{exp} = D_{w/o} + D_{w/}$$

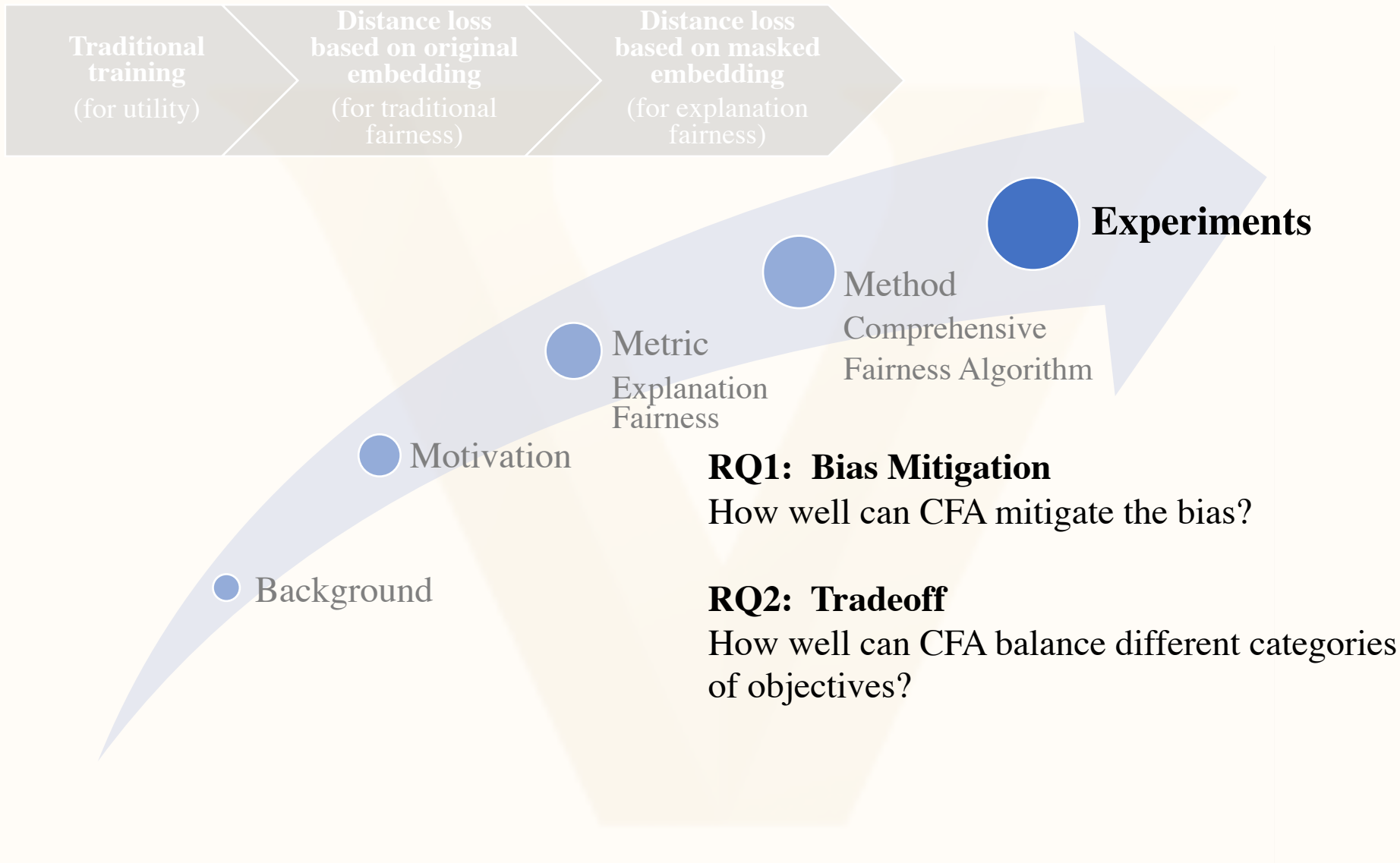
Traditional fairness

Explanation fairness

Content



Content



Experimental Setting

Dataset	Dataset	German	Recidivism	Math	Por
	# Nodes	1,000	18,876	649	649
	# Features	27	18	33	33
	Sens.	Gender	Race	Gender	Gender
	Label	Credit Risk	Recidivism	Grade	Grade

Evaluation metrics

- Utility(\uparrow): accuracy, F1, AUC
- Traditional fairness (result-oriented, \downarrow): Δ_{SP} and Δ_{EO}
- Explanation fairness (procedure-oriented, \downarrow): Δ_{VEF} and Δ_{REF}
- Overall score: $\frac{AUC+F1+ACC}{3} - \frac{\Delta_{SP}+\Delta_{EO}}{2} - \frac{\Delta_{VEF}+\Delta_{REF}}{2}$ (model selection)

Baselines

- (1) Reweight^[1]: [reweighing-based] reweight the training loss
- (2) Reduction^[2]: [constraint-based] optimization under fairness constraints

$$\begin{aligned}\Delta_{SP} &= |P(\hat{y} = 1 | s = 0) - P(\hat{y} = 1 | s = 1)| & \Delta_{REF} &= |P(\hat{q} = 1 | s = 0) - P(\hat{q} = 1 | s = 1)| \\ \Delta_{EO} &= |P(\hat{y} = 1 | y = 1, s = 0) - P(\hat{y} = 1 | y = 1, s = 1)| & \Delta_{VEF} &= \left| \frac{1}{|\mathcal{G}_0^K|} \sum_{i \in \mathcal{G}_0^K} EQ_i - \frac{1}{|\mathcal{G}_1^K|} \sum_{i \in \mathcal{G}_1^K} EQ_i \right|\end{aligned}$$

[1] Jiang, Heinrich, et al. "Identifying and correcting label bias in machine learning." AISTATS, 2020.

[2] Agarwal, Alekh, et al. "A reductions approach to fair classification." ICML, 2018.

RQ1: Bias Mitigation

Dataset	Metric	MLP	Reduction	Reweight	CFA
Recidivism	AUC↑	86.12 ± 1.91	81.17 ± 0.00	89.24 ± 0.00	89.02 ± 0.86
	F1↑	76.54 ± 2.52	<u>76.69 ± 0.00</u>	72.99 ± 0.00	81.28 ± 1.35
	Acc↑	83.48 ± 1.53	<u>84.66 ± 0.00</u>	83.70 ± 0.00	87.17 ± 0.84
	Δ_{SP} ↓	6.07 ± 2.18	<u>2.04 ± 0.00</u>	4.27 ± 0.00	1.16 ± 0.49
	Δ_{EO} ↓	<u>3.19 ± 0.73</u>	4.66 ± 0.00	3.37 ± 0.00	1.14 ± 0.39
	Δ_{REF} ↓	4.45 ± 2.96	0.53 ± 0.00	<u>1.34 ± 0.91</u>	1.98 ± 1.23
	Δ_{VEF} ↓	<u>2.1 ± 1.38</u>	2.06 ± 0.00	3.22 ± 0.00	2.70 ± 0.78
Score ↑	74.15 ± 2.03	<u>76.19 ± 0.00</u>	75.88 ± 0.00	82.33 ± 0.62	
Por	AUC↑	<u>90.86 ± 0.35</u>	67.64 ± 0.00	89.07 ± 0.00	91.30 ± 0.55
	F1↑	<u>58.41 ± 4.10</u>	51.43 ± 0.00	51.43 ± 0.00	60.55 ± 4.73
	Acc↑	<u>89.57 ± 0.78</u>	<u>89.57 ± 0.00</u>	<u>89.57 ± 0.00</u>	89.82 ± 1.00
	Δ_{SP} ↓	2.08 ± 0.75	<u>1.93 ± 0.00</u>	<u>1.93 ± 0.00</u>	1.00 ± 0.72
	Δ_{EO} ↓	32.35 ± 7.07	20.59 ± 0.00	20.59 ± 0.00	<u>27.65 ± 5.44</u>
	Δ_{REF} ↓	8.68 ± 3.18	1.37 ± 0.00	8.68 ± 0.00	<u>4.66 ± 3.76</u>
	Δ_{VEF} ↓	<u>4.44 ± 2.22</u>	0.00 ± 0.00	7.69 ± 0.00	4.70 ± 3.67
	Score ↑	55.83 ± 3.97	<u>57.60 ± 0.00</u>	57.25 ± 0.00	61.55 ± 3.26

Takes up largest proportion of bold/underline

Bold text: best performance

Underline text: second best performance

RQ1: Bias Mitigation

Dataset	Metric	MLP	Reduction	Reweight	CFA	
Recidivism	AUC↑	86.12 ± 1.91	81.17 ± 0.00	89.24 ± 0.00	89.02 ± 0.86	
	F1↑	76.54 ± 2.52	76.69 ± 0.00	72.99 ± 0.00	81.28 ± 1.35	
	Acc↑	83.48 ± 1.53	84.66 ± 0.00	83.70 ± 0.00	87.17 ± 0.84	
	Δ_{SP} ↓	6.07 ± 2.18	2.04 ± 0.00	4.27 ± 0.00	1.16 ± 0.49	
	Δ_{EO} ↓	3.19 ± 0.73	4.66 ± 0.00	3.37 ± 0.00	1.14 ± 0.39	
	Δ_{REF} ↓	4.45 ± 2.96	0.53 ± 0.00	1.34 ± 0.91	1.98 ± 1.23	
	Δ_{VEF} ↓	2.1 ± 1.38	2.06 ± 0.00	3.22 ± 0.00	2.70 ± 0.78	
	Score ↑	74.15 ± 2.03	76.19 ± 0.00	75.88 ± 0.00	82.33 ± 0.62	
	Por	AUC↑	90.86 ± 0.35	67.64 ± 0.00	89.07 ± 0.00	91.30 ± 0.55
		F1↑	58.41 ± 4.10	51.43 ± 0.00	51.43 ± 0.00	60.55 ± 4.73
Acc↑		89.57 ± 0.78	89.57 ± 0.00	89.57 ± 0.00	89.82 ± 1.00	
Δ_{SP} ↓		2.08 ± 0.75	1.93 ± 0.00	1.93 ± 0.00	1.00 ± 0.72	
Δ_{EO} ↓		32.35 ± 7.07	20.59 ± 0.00	20.59 ± 0.00	27.65 ± 5.44	
Δ_{REF} ↓		8.68 ± 3.18	1.37 ± 0.00	8.68 ± 0.00	4.66 ± 3.76	
Δ_{VEF} ↓		4.44 ± 2.22	0.00 ± 0.00	7.69 ± 0.00	4.70 ± 3.67	
Score ↑	55.83 ± 3.97	57.60 ± 0.00	57.25 ± 0.00	61.55 ± 3.26		

Utility Performance

Traditional Fairness

Explanation Fairness

Comparable or better than baselines

RQ1: Bias Mitigation

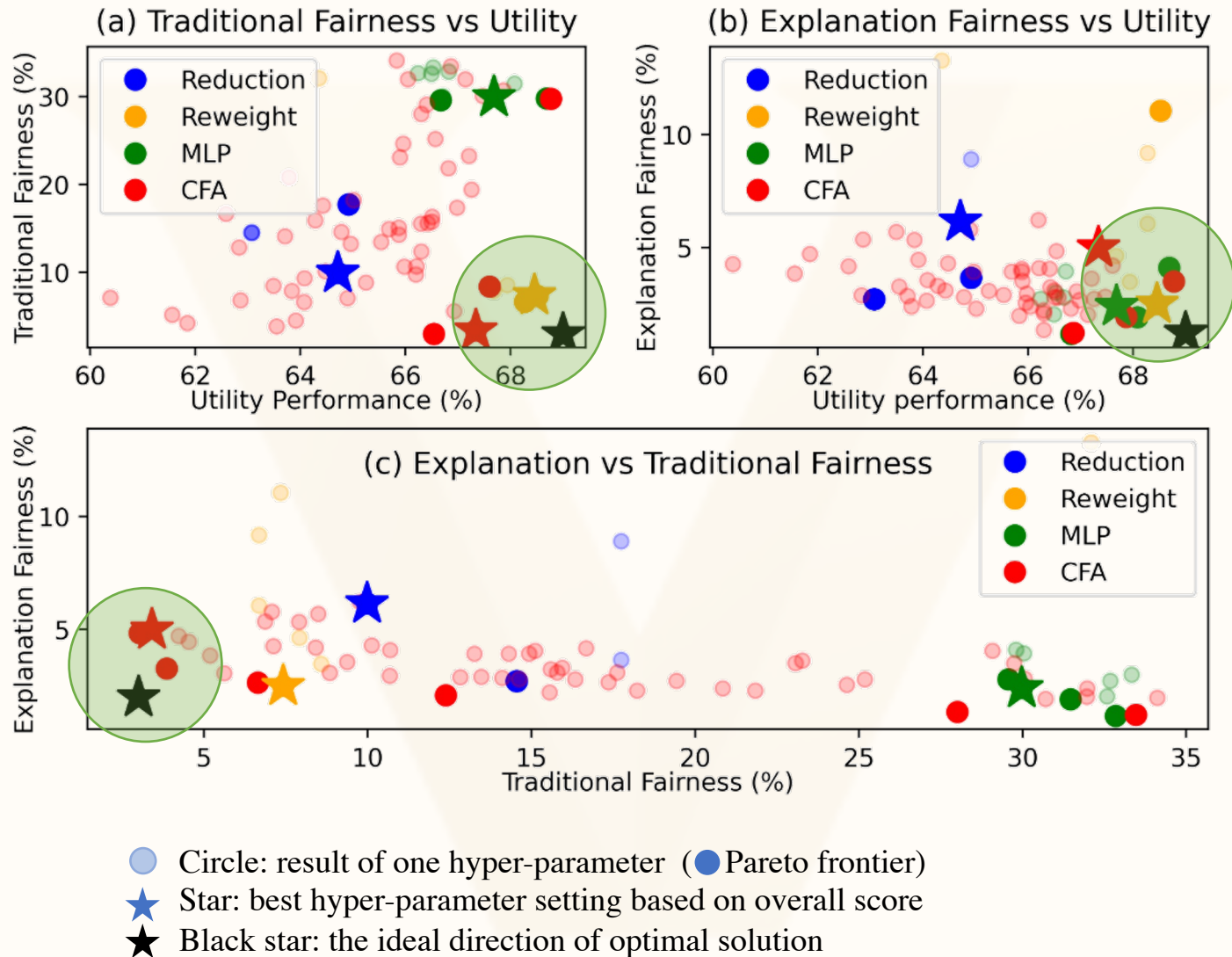
Dataset	Metric	MLP	Reduction	Reweight	CFA
Recidivism	AUC↑	86.12 ± 1.91	81.17 ± 0.00	89.24 ± 0.00	89.02 ± 0.86
	F1↑	76.54 ± 2.52	76.69 ± 0.00	72.99 ± 0.00	81.28 ± 1.35
	Acc↑	83.48 ± 1.53	84.66 ± 0.00	83.70 ± 0.00	87.17 ± 0.84
	Δ_{SP} ↓	6.07 ± 2.18	2.04 ± 0.00	4.27 ± 0.00	1.16 ± 0.49
	Δ_{EO} ↓	3.19 ± 0.73	4.66 ± 0.00	3.37 ± 0.00	1.14 ± 0.39
	Δ_{REF} ↓	4.45 ± 2.96	0.53 ± 0.00	1.34 ± 0.91	1.98 ± 1.23
	Δ_{VEF} ↓	2.1 ± 1.38	2.06 ± 0.00	3.22 ± 0.00	2.70 ± 0.78
Score ↑	74.15 ± 2.03	76.19 ± 0.00	75.88 ± 0.00	82.33 ± 0.62	
Por	AUC↑	90.86 ± 0.35	67.64 ± 0.00	89.07 ± 0.00	91.30 ± 0.55
	F1↑	58.41 ± 4.10	51.43 ± 0.00	51.43 ± 0.00	60.55 ± 4.73
	Acc↑	89.57 ± 0.78	89.57 ± 0.00	89.57 ± 0.00	89.82 ± 1.00
	Δ_{SP} ↓	2.08 ± 0.75	1.93 ± 0.00	1.93 ± 0.00	1.00 ± 0.72
	Δ_{EO} ↓	32.35 ± 7.07	20.59 ± 0.00	20.59 ± 0.00	27.65 ± 5.44
	Δ_{REF} ↓	8.68 ± 3.18	1.37 ± 0.00	8.68 ± 0.00	4.66 ± 3.76
	Δ_{VEF} ↓	4.44 ± 2.22	0.00 ± 0.00	7.69 ± 0.00	4.70 ± 3.67
Score ↑	55.83 ± 3.97	57.60 ± 0.00	57.25 ± 0.00	61.55 ± 3.26	

Overall Score

The highest for all datasets

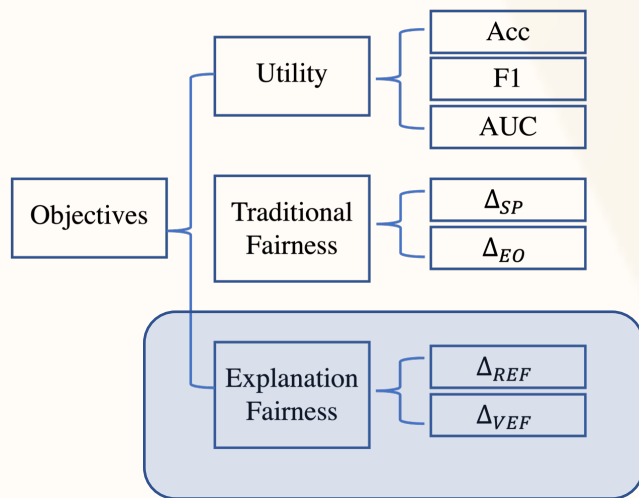
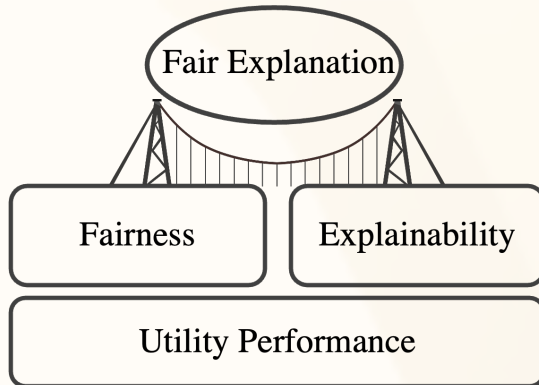
$$\text{Overall score: } \frac{\text{AUC} + \text{F1} + \text{ACC}}{3} - \frac{\Delta_{SP} + \Delta_{EO}}{2} - \frac{\Delta_{VEF} + \Delta_{REF}}{2} \text{ (model selection)}$$

RQ2: Tradeoff

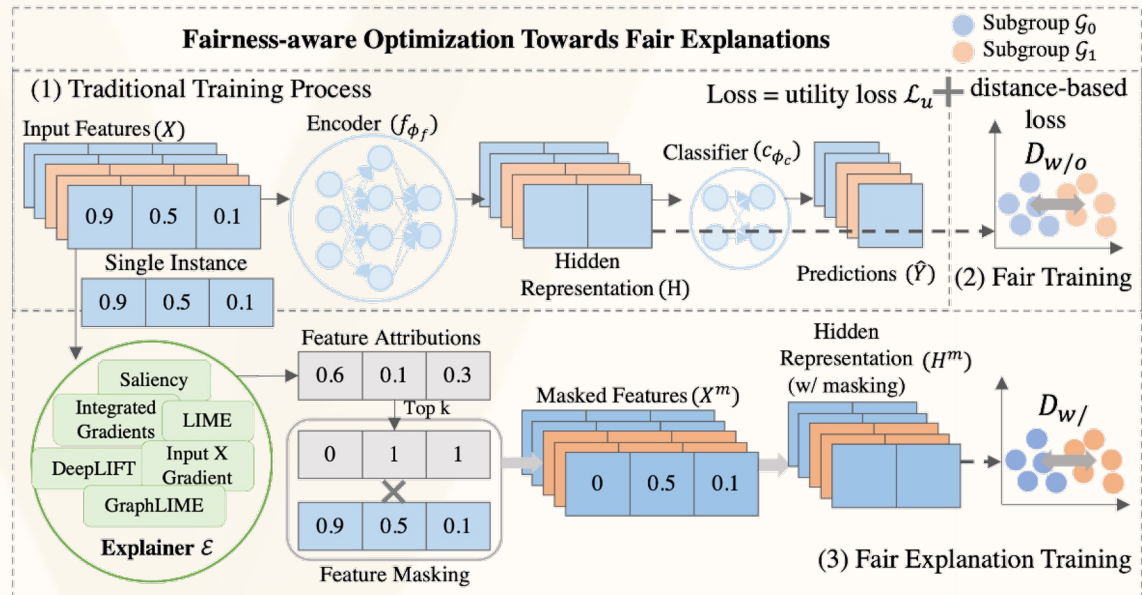
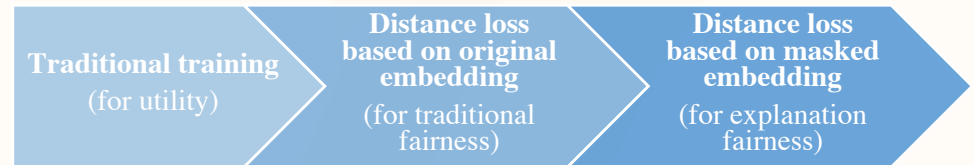


Summary

Novel Fairness Perspective

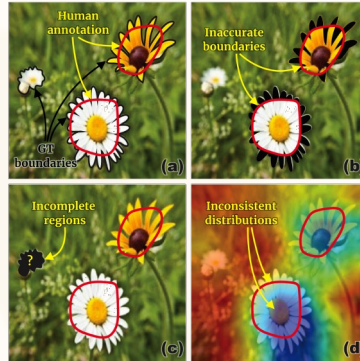


Comprehensive Fairness Algorithm

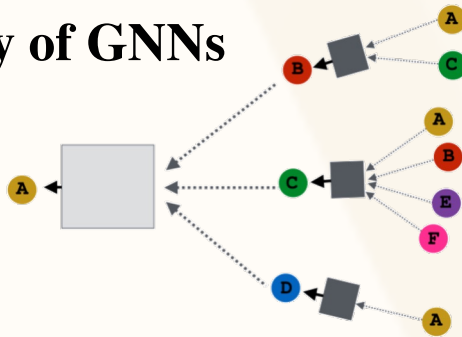


Future Directions

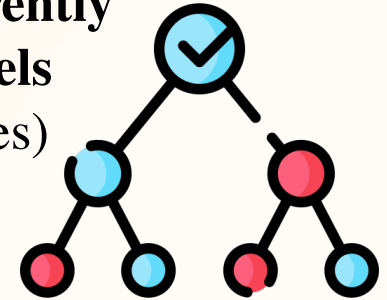
Extending CFA
towards fair model
explanations in
other data types
(e.g., images)



**Improved Fairness and
Explainability of GNNs**



**Defining novel fair explanation
metrics for inherently
explainable models**
(e.g., decision trees)

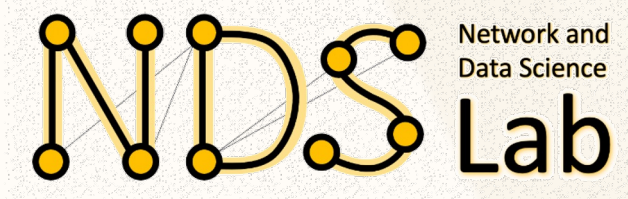
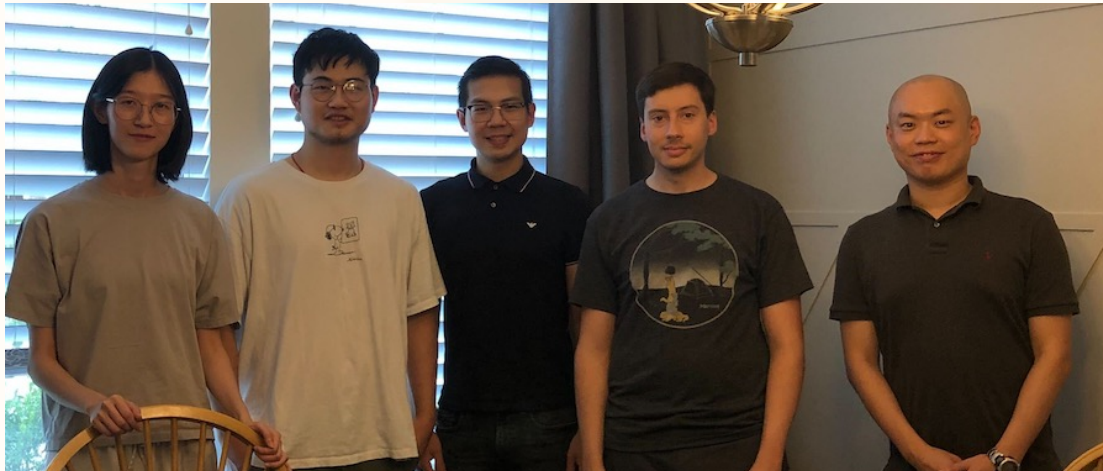


<https://yuyingzhao.github.io/>



Please see my website for
other work

Acknowledgement



AAI-23 Student Scholar,
Diversity and Inclusion Scholar,
& Volunteer Program