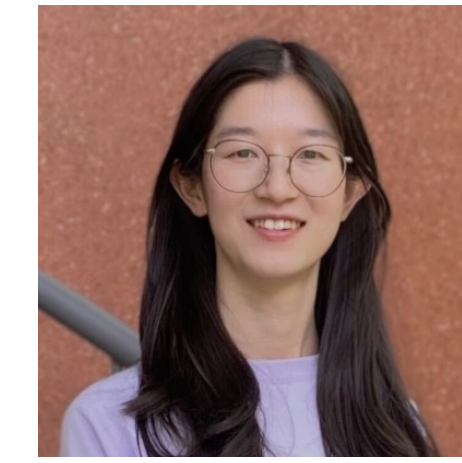# Fairness and Explainability: Bridging the Gap Towards Fair Model Explanations

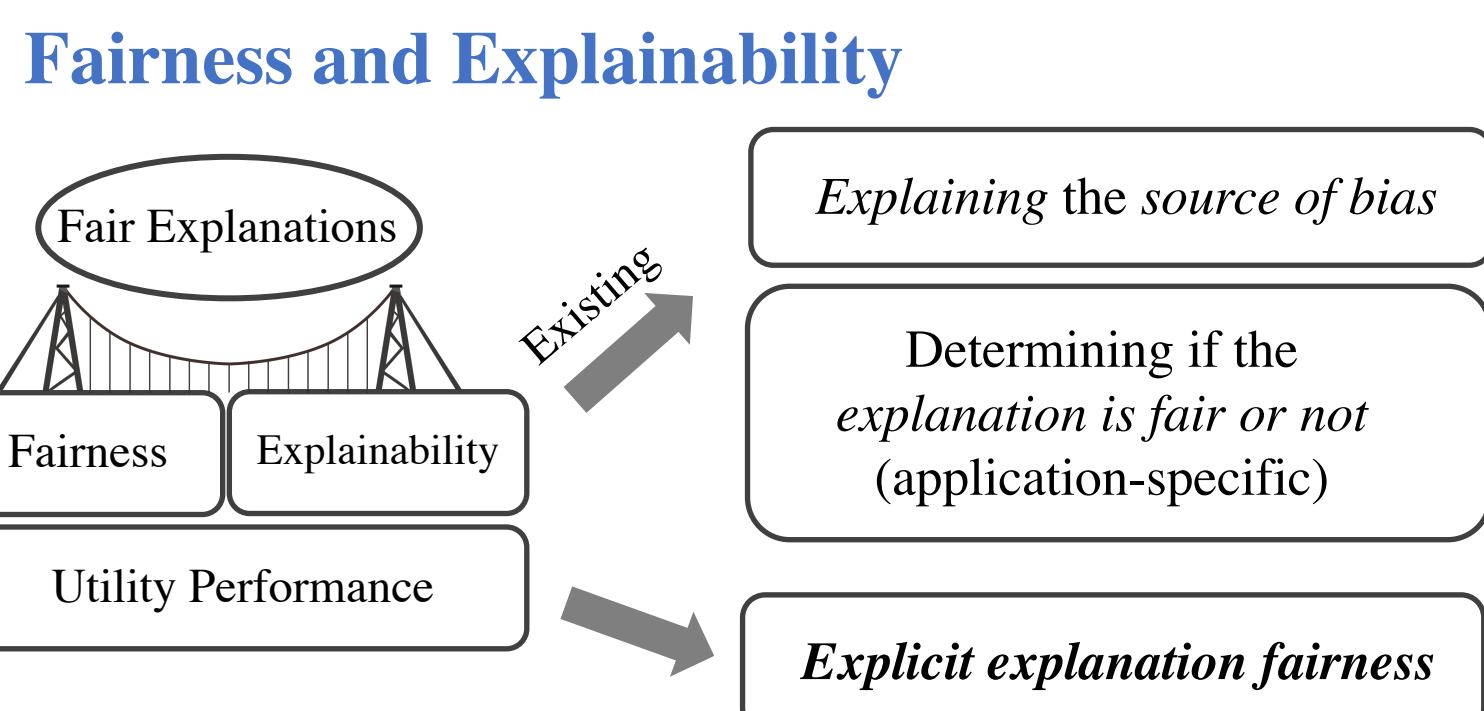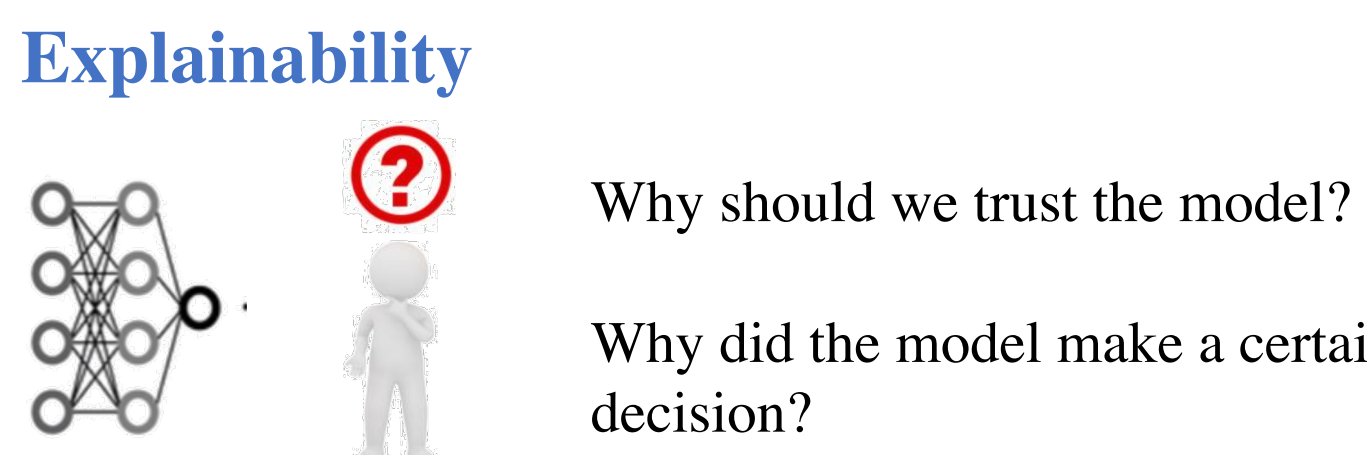**Yuying Zhao**   Yu Wang   Tyler Derr

## Fairness and Explainability
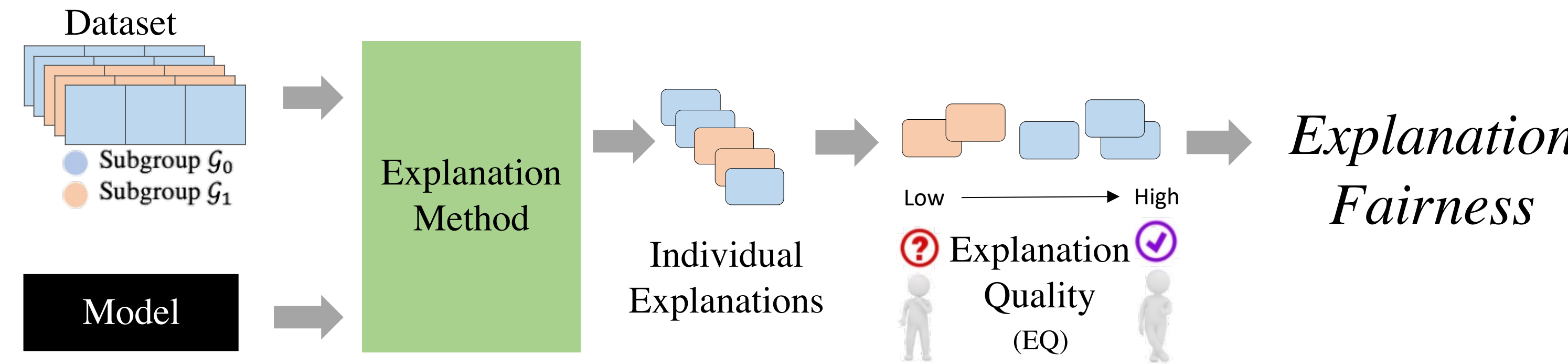
**Fairness** $s$: subgroup   $\hat{y}$: label

<u>Results</u>

Biased

Unbiased

**Explainability**

Why should we trust the model?

Why did the model make a certain decision?

**Fairness and Explainability**

Fair Explanations

Fairness | Explainability

Utility Performance

Existing →
- *Explaining* the *source of bias*
- Determining if the *explanation is fair or not* (application-specific)

→ ***Explicit explanation fairness***

## Motivation

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

**Result-oriented**

$x_i$ → Model → $y_i$
$x_j$ → Model → $y_j$

**Fair Predictions**

$x_0$ $x_1$ $x_2$ Hired
$x_3$ $x_4$ $x_5$ Not

**Procedure-oriented**

$x_i$ → $e_i$ → $y_i$
$x_j$ → $e_j$ → $y_j$

Model explanations

**Unfair Explanations**

$x_0$ 0.8  $x_1$ 0.8  $x_2$ 0.6  Higher EQ
$x_3$ 0.7  $x_4$ 0.2  $x_5$ 0.7  Lower EQ

Explanation Quality (EQ)

Low ——→ High

? Explanation Quality ✓

Explanations of Different Quality
↓
Different Levels of Trust
↓
Different Treatments
↓
Group Unfairness

🟦 higher-quality and clear explanation
🟧 lower-quality and ambiguous explanation

## Novel Procedure-Oriented Fairness Perspective

Dataset
- Subgroup $\mathcal{G}_0$
- Subgroup $\mathcal{G}_1$

→ Explanation Method → Individual Explanations → Explanation Quality (EQ) Low → High → *Explanation Fairness*

Model →

? Explanation Quality ✓

**(1) Ratio-based Fairness $\Delta_{REF}$**

$$\Delta_{SP} = |P(\hat{y}=1|s=0) - P(\hat{y}=1|s=1)|$$
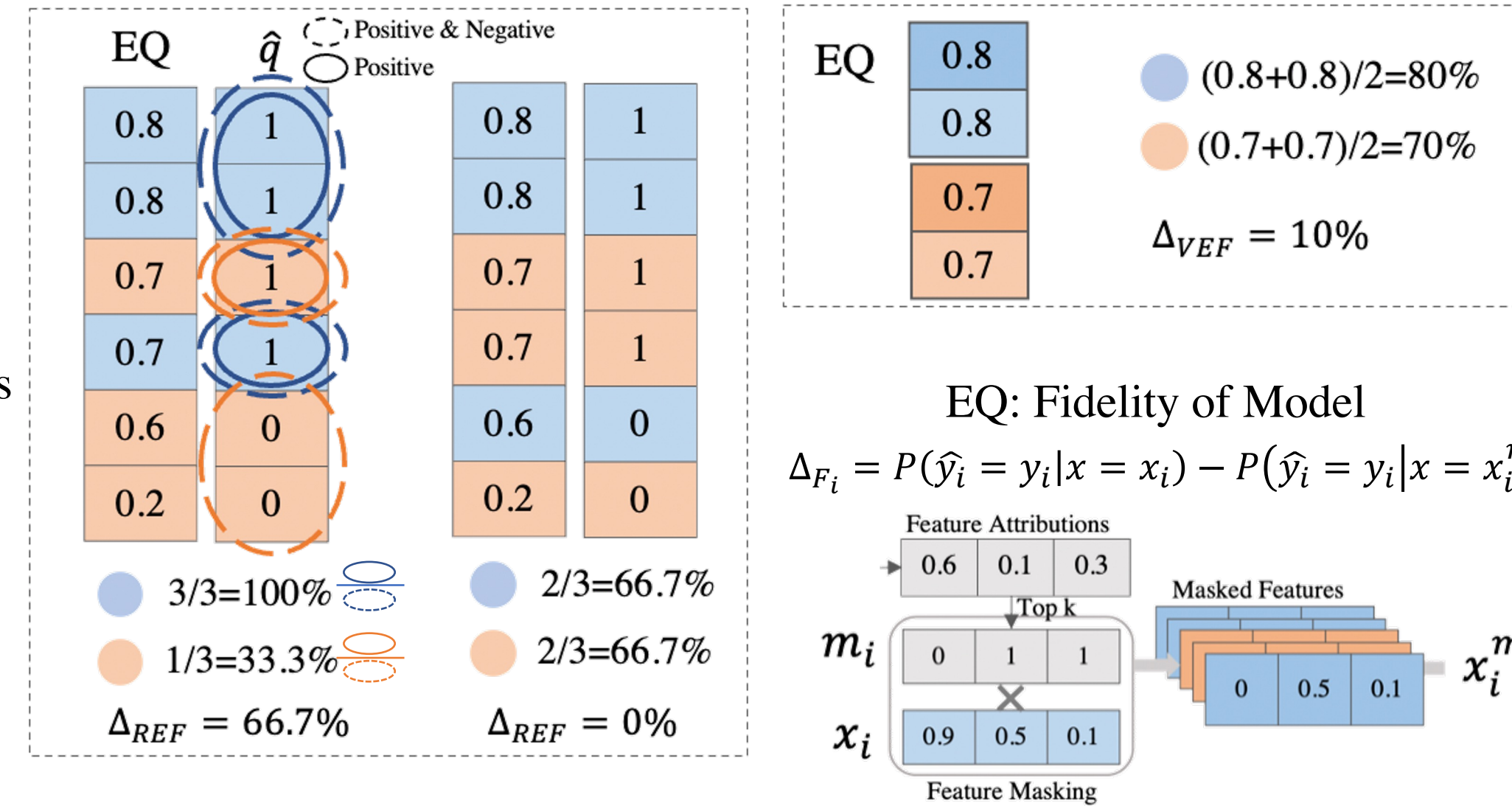
Same opportunity of having positive prediction

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

Same opportunity of having high-quality explanations

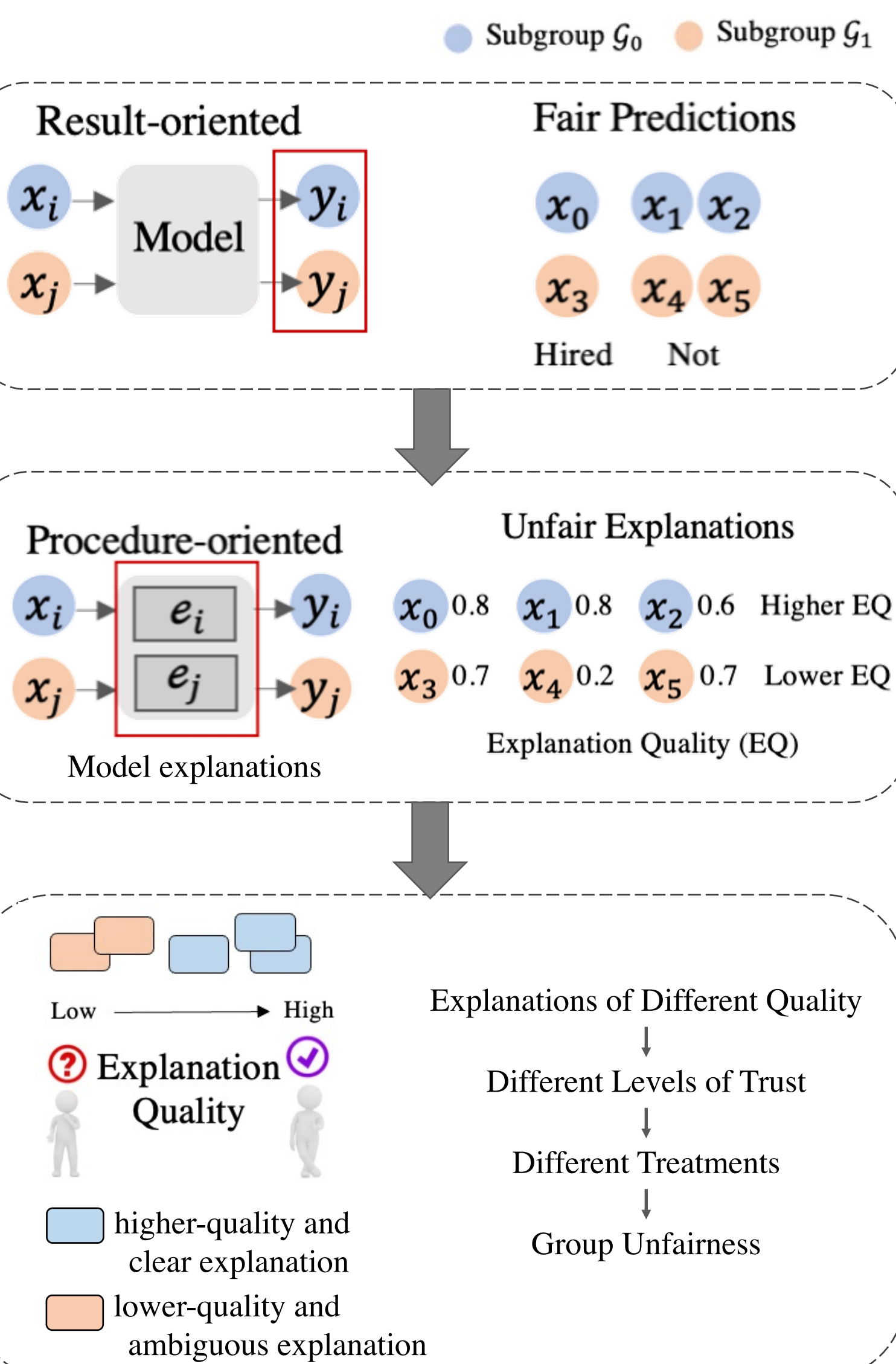$\hat{y}$: prediction   $\hat{q}$: explanation quality   $s$: sensitive labels
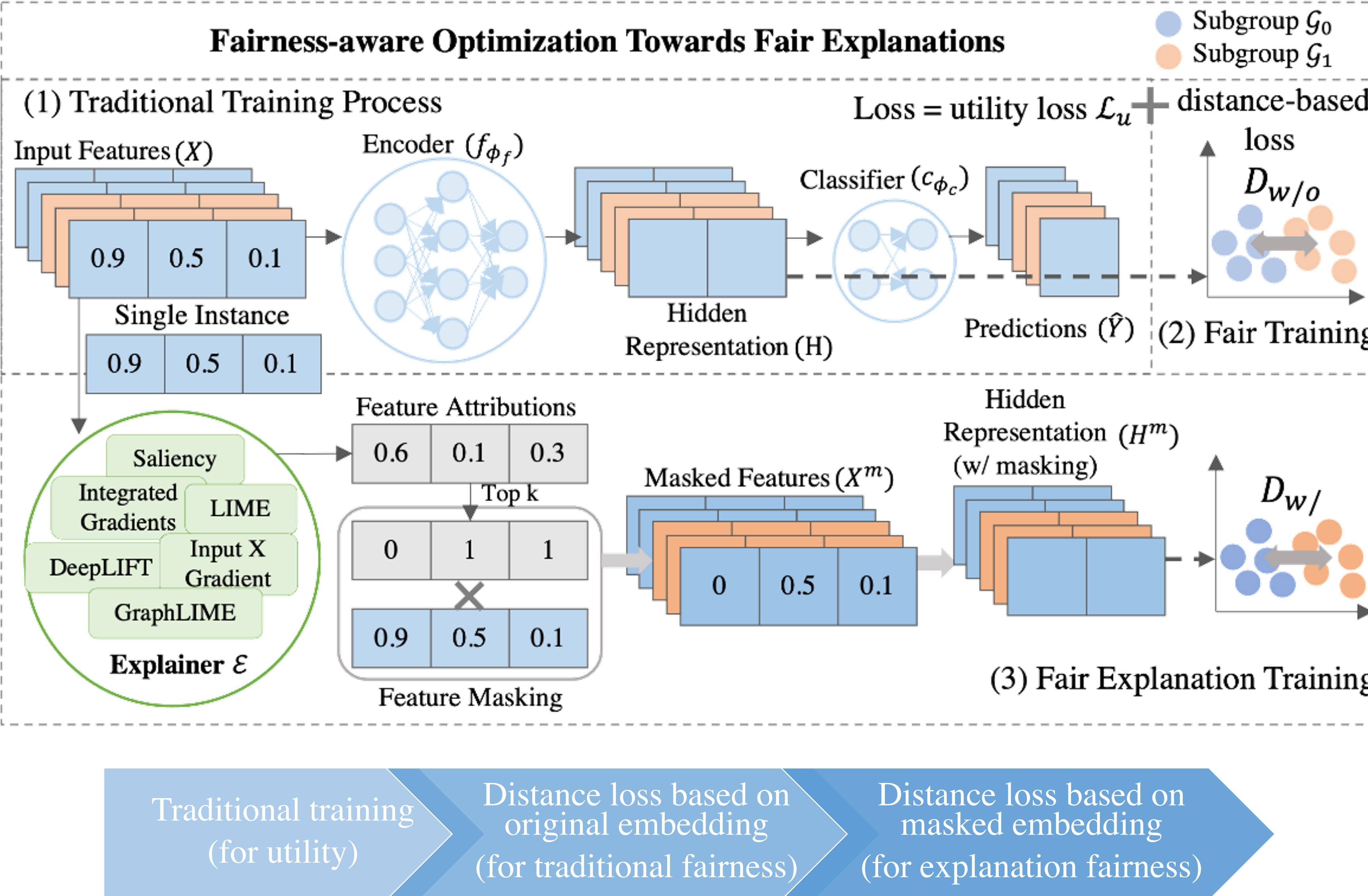
**(2) Value-based Fairness $\Delta_{VEF}$**

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

EQ   $\hat{q}$   ⬭ Positive & Negative ◯ Positive

| 0.8 | 1 |
| 0.8 | 1 |
| 0.7 | 1 |
| 0.7 | 1 |
| 0.6 | 0 |
| 0.2 | 0 |

3/3=100%   2/3=66.7%
1/3=33.3%  2/3=66.7%

$\Delta_{REF}$ = 66.7%   $\Delta_{REF}$ = 0%

EQ 0.8 → (0.8+0.8)/2=80%
EQ 0.8
EQ 0.7 → (0.7+0.7)/2=70%
EQ 0.7

$\Delta_{VEF}$ = 10%

**EQ: Fidelity of Model**

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

Feature Attributions
| 0.6 | 0.1 | 0.3 |
Top k
| 0 | 1 | 1 |  Masked Features
× → | 0 | 0.5 | 0.1 | $x_i^{m_i}$
| 0.9 | 0.5 | 0.1 |
Feature Masking

## Comprehensive Fairness Algorithm (CFA)

**Fairness-aware Optimization Towards Fair Explanations**

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

**(1) Traditional Training Process**

Input Features ($X$) → Encoder ($f_{\phi_f}$) → Hidden Representation ($H$) → Classifier ($c_{\phi_c}$) → Predictions ($\hat{Y}$)

| 0.9 | 0.5 | 0.1 |
Single Instance
| 0.9 | 0.5 | 0.1 |

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

$D_{w/o}$

**(2) Fair Training**

**Explainer $\mathcal{E}$**: Saliency, Integrated Gradients, LIME, DeepLIFT, Input X Gradient, GraphLIME

→ Feature Attributions
| 0.6 | 0.1 | 0.3 |
Top k
| 0 | 1 | 1 |
×
| 0.9 | 0.5 | 0.1 |
Feature Masking

→ Masked Features ($X^m$)
| 0 | 0.5 | 0.1 |

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

$D_{w/}$

**(3) Fair Explanation Training**

Traditional training (for utility) → Distance loss based on original embedding (for traditional fairness) → Distance loss based on masked embedding (for explanation fairness)

## Experiments

### 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** |
| | $\Delta_{SP}$↓ | 6.07 ± 2.18 | 2.04 ± 0.00 | 4.27 ± 0.00 | **1.16 ± 0.49** |
| | $\Delta_{EO}$↓ | 3.19 ± 0.73 | 4.66 ± 0.00 | 3.37 ± 0.00 | **1.14 ± 0.39** |
| | $\Delta_{REF}$↓ | 4.45 ± 2.96 | **0.53 ± 0.00** | 1.34 ± 0.91 | 1.98 ± 1.23 |
| | $\Delta_{VEF}$↓ | 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** |
| | $\Delta_{SP}$↓ | 2.08 ± 0.75 | 1.93 ± 0.00 | 1.93 ± 0.00 | **1.00 ± 0.72** |
| | $\Delta_{EO}$↓ | 32.35 ± 7.07 | **20.59 ± 0.00** | **20.59 ± 0.00** | 27.65 ± 5.44 |
| | $\Delta_{REF}$↓ | 8.68 ± 3.18 | **1.37 ± 0.00** | 8.68 ± 0.00 | 4.66 ± 3.76 |
| | $\Delta_{VEF}$↓ | 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
🟨 Overall Score

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

### RQ2: Multi-objective Tradeoff

(a) Traditional Fairness vs Utility
Legend: Reduction, Reweight, MLP, CFA

(b) Explanation Fairness vs Utility
Legend: Reduction, Reweight, MLP, CFA

(c) Explanation vs Traditional Fairness
Legend: Reduction, Reweight, MLP, CFA

## Future Directions

Extending CFA
- Data types
- EQ metrics

Fair and Explainable Graph Neural Networks

Fair Explanation for Inherently Explainable Models (e.g., Decision Tree)