

# Fairness under Graph Uncertainty: Achieving Interventional Fairness with Partially Known Causal Graphs over Clusters of Variables



uai2026

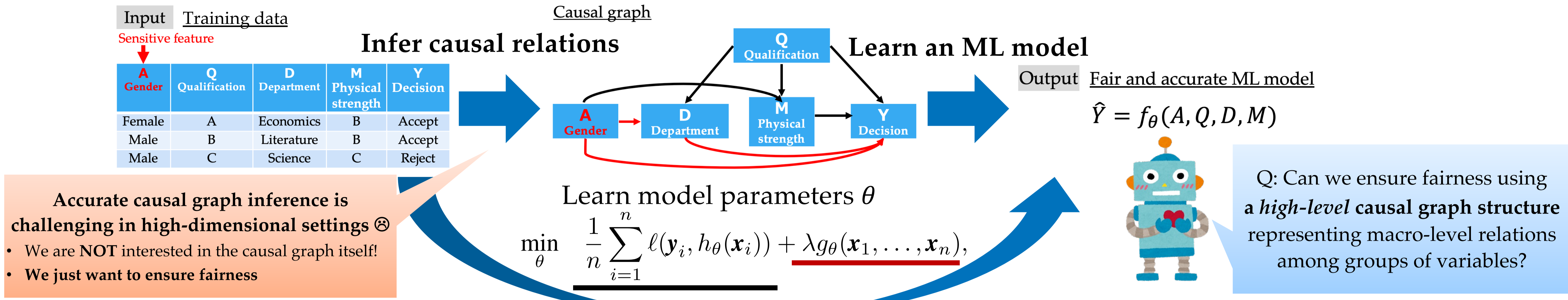


Yoichi Chikahara<sup>1</sup>.

<sup>1</sup>Communication Science Laboratories, NTT, Inc.

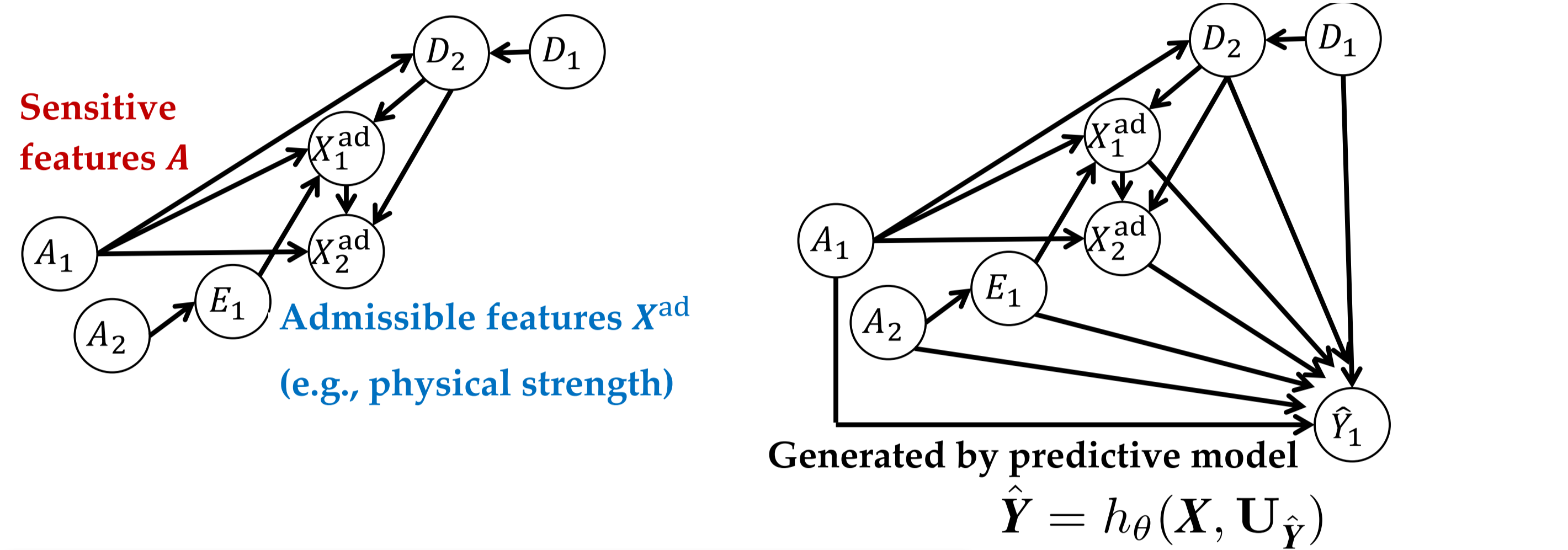
This work was supported by JST ACT-X (JPMJAX23CF).

## Problem Setup & Motivation: Can we achieve fairness using a cluster causal graph?

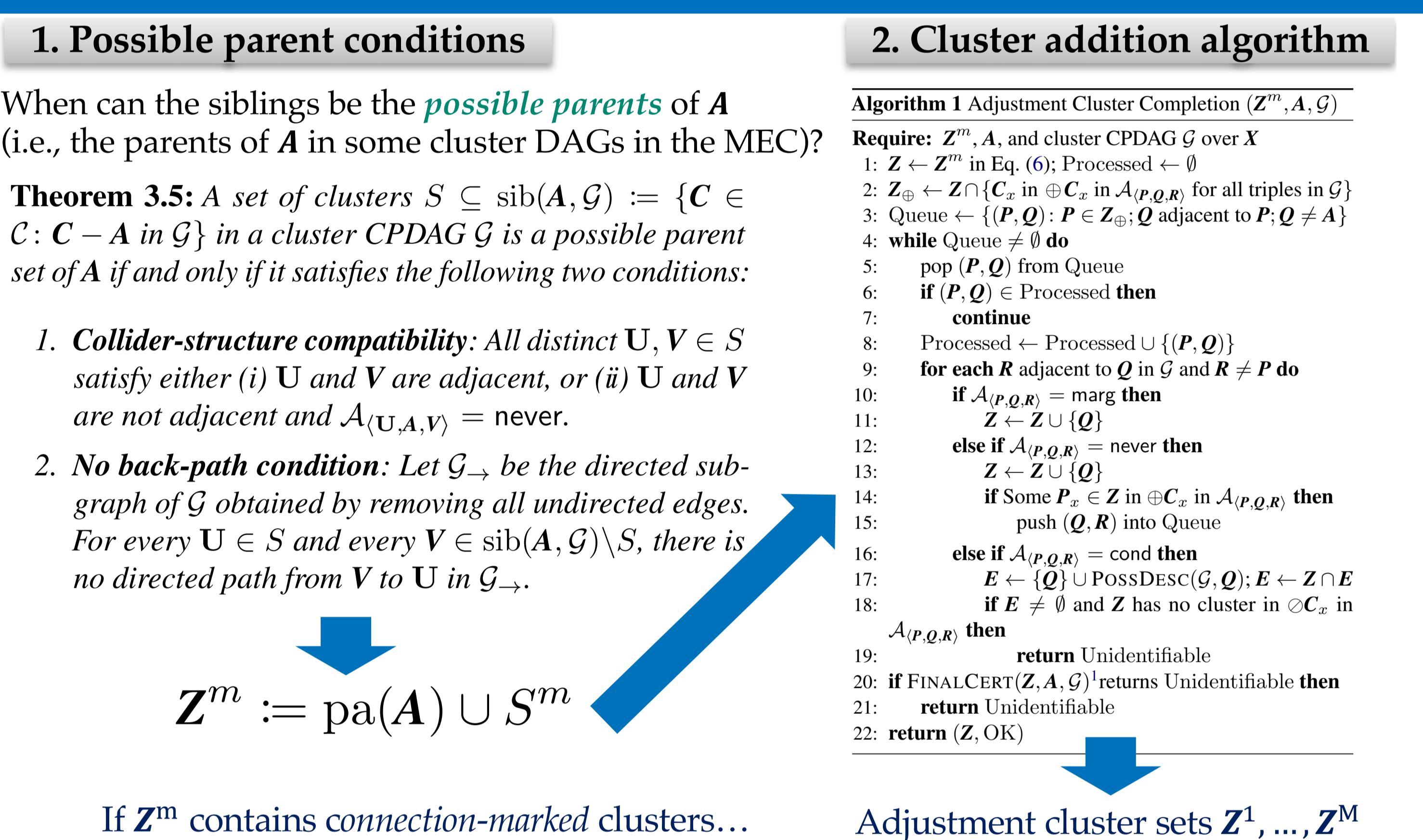


## Interventional Fairness

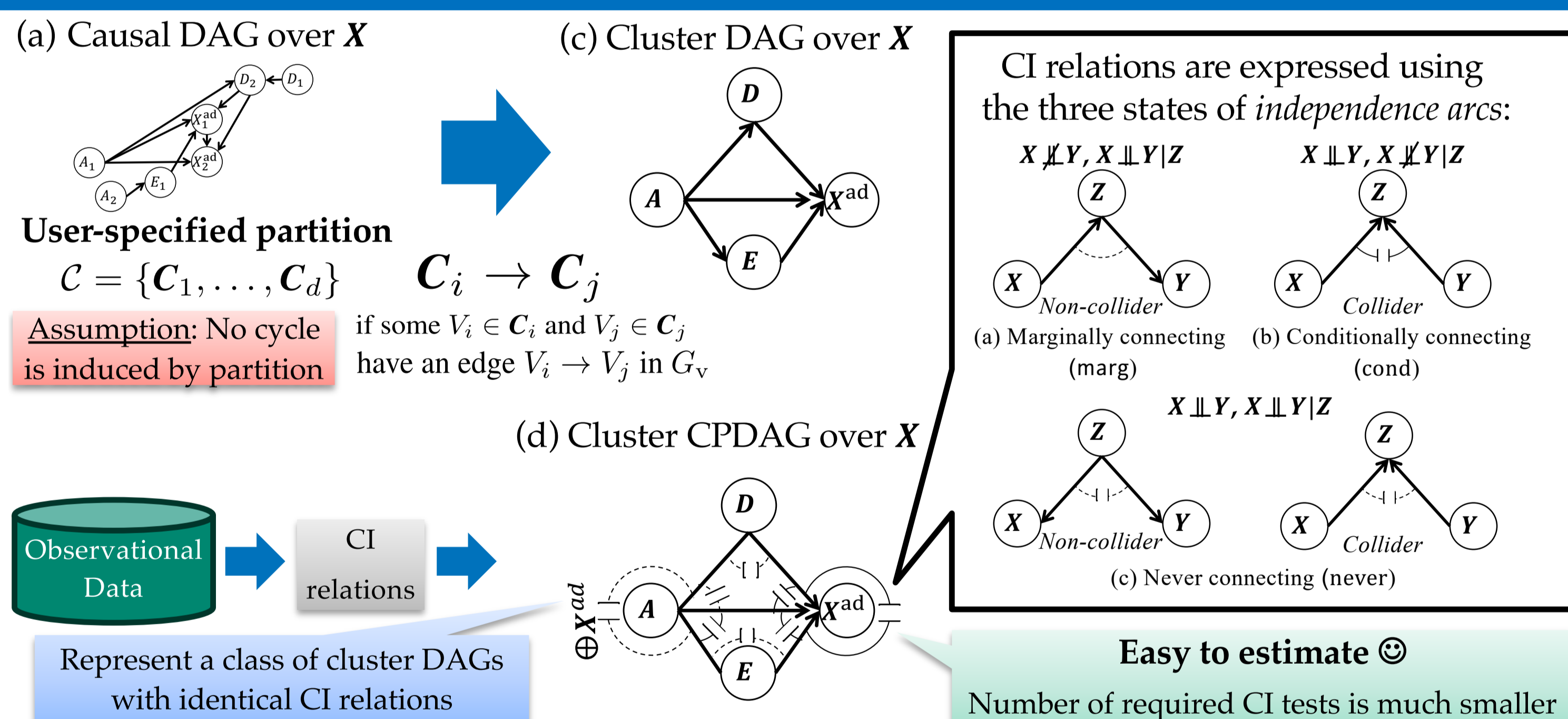
(a) Causal graph over features  $X$  (b) Causal graph over features  $X$  & prediction  $\hat{Y}$



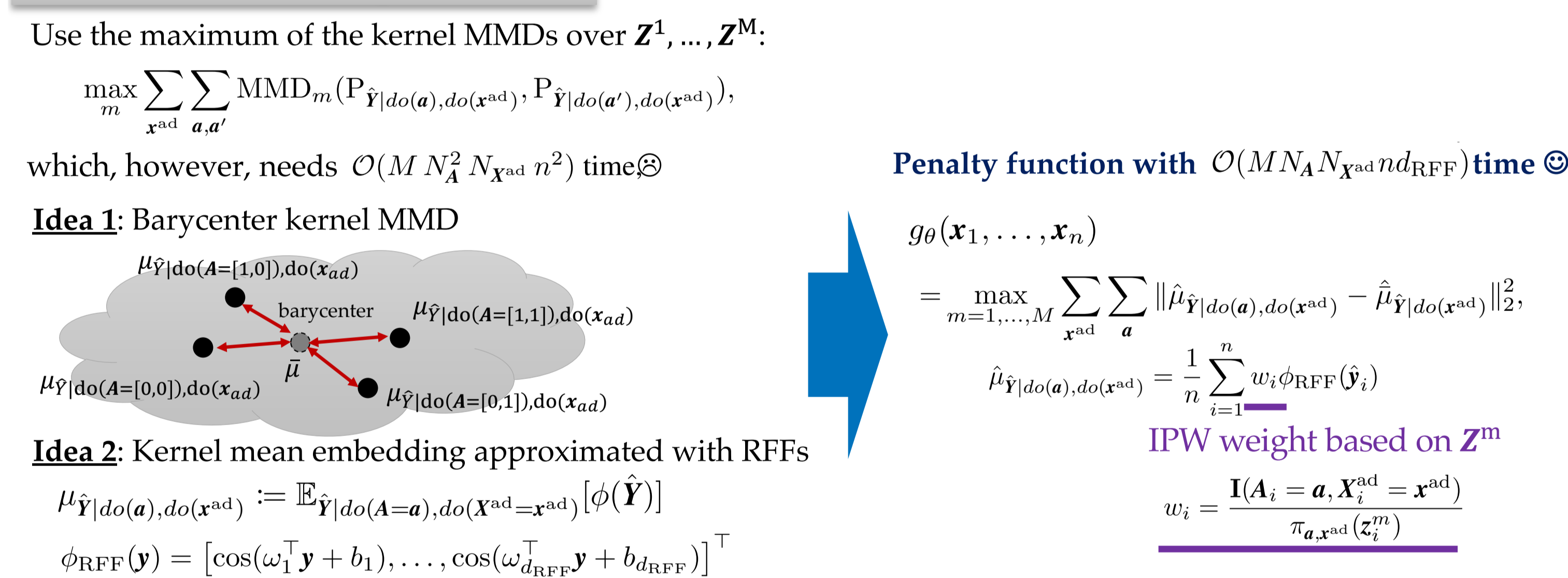
## Graphical Algorithm & Learning Framework



## Cluster DAGs & Cluster CPDAGs



## Worst-case unfairness penalization



## Difficulty & Our Proposed Strategies

**Q: How can we infer interventional distributions using a cluster CPDAG?**

If we have access to the variable-level DAG:

We can identify back-door adjustment variables  $Z$  from the DAG and compute

$$P(\hat{Y} = \hat{y} \mid do(A = a), do(X^{ad} = x^{ad})) = \mathbb{E}_{Z, X^{re} \mid a, x^{ad}} [P(\hat{Y} = \hat{y} \mid A = a, X^{ad} = x^{ad}, Z, X^{re})]$$

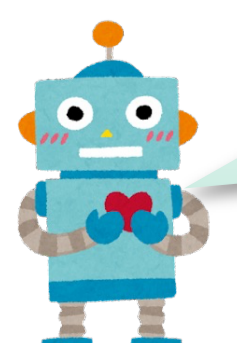
**Difficulty:** Cluster CPDAG represents multiple DAGs (in the cluster MEC)

**Our idea:** We obtain adjustment candidate sets  $Z^1, \dots, Z^M$ , at least one of which d-separates all back-door paths for the true DAG. We take 2 steps:

1. **Parent Enumeration:** From possible DAGs, enumerate the parents of  $A$  as

$$Z^m := \text{pa}(A) \cup S^m \quad \text{for } m = 1, \dots, M$$

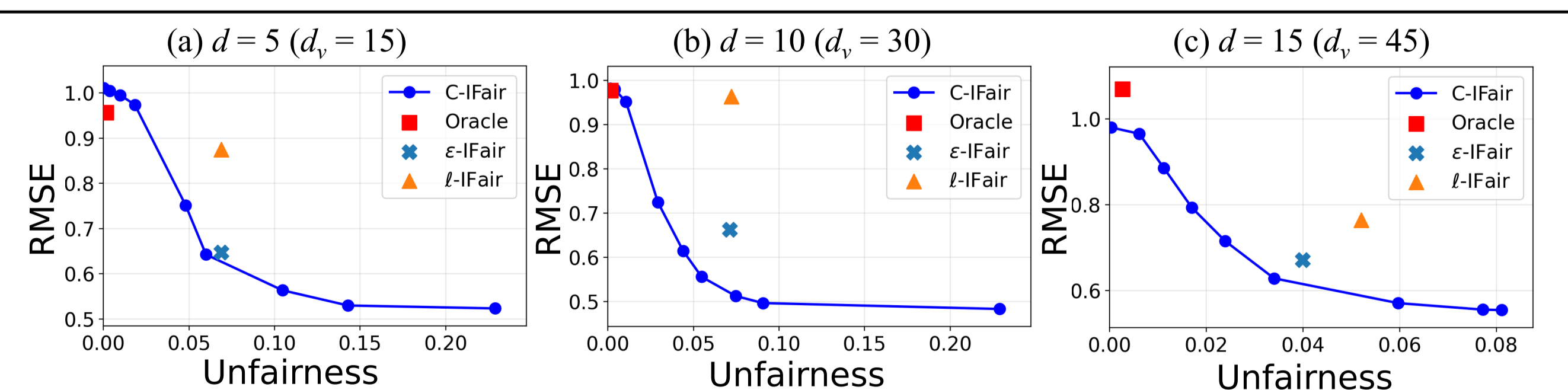
2. **Adjustment Set Completion:** Add clusters that are needed to achieve cluster d-separation [Anand+; NeurIPS2025]



After taking these 2 steps, we formulate unfairness penalty  $g_{\theta}$  as the maximum of the distributional distances over  $Z^1, \dots, Z^M$

## Experimental Results

	$d=5$ ( $d_v=15$ )		$d=10$ ( $d_v=30$ )		$d=15$ ( $d_v=45$ )	
	RMSE $\downarrow$	UNFAIRNESS $\downarrow$	RMSE $\downarrow$	UNFAIRNESS $\downarrow$	RMSE $\downarrow$	UNFAIRNESS $\downarrow$
LINEAR						
ORACLE	0.957 $\pm$ 0.037	0.000 $\pm$ 0.000	0.977 $\pm$ 0.044	0.000 $\pm$ 0.000	1.069 $\pm$ 0.143	0.000 $\pm$ 0.000
FULL	0.523 $\pm$ 0.195	0.259 $\pm$ 0.129	0.483 $\pm$ 0.195	0.229 $\pm$ 0.136	0.554 $\pm$ 0.210	0.081 $\pm$ 0.117
UNAWARE	0.774 $\pm$ 0.154	0.071 $\pm$ 0.059	0.774 $\pm$ 0.154	0.062 $\pm$ 0.079	0.793 $\pm$ 0.120	0.046 $\pm$ 0.104
NO-DESCs	0.739 $\pm$ 0.140	0.064 $\pm$ 0.089	0.739 $\pm$ 0.140	0.065 $\pm$ 0.100	0.699 $\pm$ 0.135	0.033 $\pm$ 0.028
$\epsilon$ -IFAIR	0.647 $\pm$ 0.030	0.069 $\pm$ 0.031	0.663 $\pm$ 0.062	0.071 $\pm$ 0.011	0.671 $\pm$ 0.044	0.040 $\pm$ 0.011
$\ell$ -IFAIR	0.875 $\pm$ 0.028	0.069 $\pm$ 0.081	0.964 $\pm$ 0.049	0.072 $\pm$ 0.052	0.764 $\pm$ 0.049	0.052 $\pm$ 0.052
<b>C-IFAIR</b>	0.643 $\pm$ 0.127	0.060 $\pm$ 0.054	0.660 $\pm$ 0.123	0.056 $\pm$ 0.052	0.669 $\pm$ 0.171	0.020 $\pm$ 0.036



	ADULT		GERMAN		OULAD	
	AUC $\uparrow$	UNFAIRNESS $\downarrow$	AUC $\uparrow$	UNFAIRNESS $\downarrow$	AUC $\uparrow$	UNFAIRNESS $\downarrow$
ORACLE	0.709 $\pm$ 0.009	0.000 $\pm$ 0.000	0.582 $\pm$ 0.004	0.000 $\pm$ 0.000	0.628 $\pm$ 0.015	0.000 $\pm$ 0.000
FULL	0.874 $\pm$ 0.004	0.030 $\pm$ 0.024	0.762 $\pm$ 0.060	0.173 $\pm$ 0.011	0.697 $\pm$ 0.024	0.061 $\pm$ 0.001
UNAWARE	0.805 $\pm$ 0.049	0.018 $\pm$ 0.014	0.695 $\pm$ 0.061	0.101 $\pm$ 0.013	0.654 $\pm$ 0.019	0.004 $\pm$ 0.000
NO-DESCs	0.815 $\pm$ 0.007	0.026 $\pm$ 0.022	0.726 $\pm$ 0.059	0.085 $\pm$ 0.008	0.654 $\pm$ 0.020	0.004 $\pm$ 0.000
$\epsilon$ -IFAIR	0.812 $\pm$ 0.006	0.015 $\pm$ 0.012	0.756 $\pm$ 0.008	0.082 $\pm$ 0.000	0.641 $\pm$ 0.005	0.005 $\pm$ 0.002
$\ell$ -IFAIR	0.764 $\pm$ 0.002	0.015 $\pm$ 0.008	0.750 $\pm$ 0.035	0.151 $\pm$ 0.011	0.639 $\pm$ 0.011	0.003 $\pm$ 0.001
<b>C-IFAIR</b>	0.829 $\pm$ 0.021	0.014 $\pm$ 0.010	0.760 $\pm$ 0.061	0.065 $\pm$ 0.008	0.660 $\pm$ 0.018	0.001 $\pm$ 0.000