



Through the Data Management Lens: Experimental Analysis and Evaluation of Fair Classification

Maliha Tashfia Islam
mtislam@cs.umass.edu
University of Massachusetts Amherst

Anna Fariha
annafariha@microsoft.com
Microsoft

Alexandra Meliou
ameli@cs.umass.edu
University of Massachusetts Amherst

Babak Salimi
bsalimi@ucsd.edu
University of California, San Diego



Classifiers can exhibit discriminatory behavior



Historical biases in training data often cause **discrimination!**

Two new challenges in fair classification

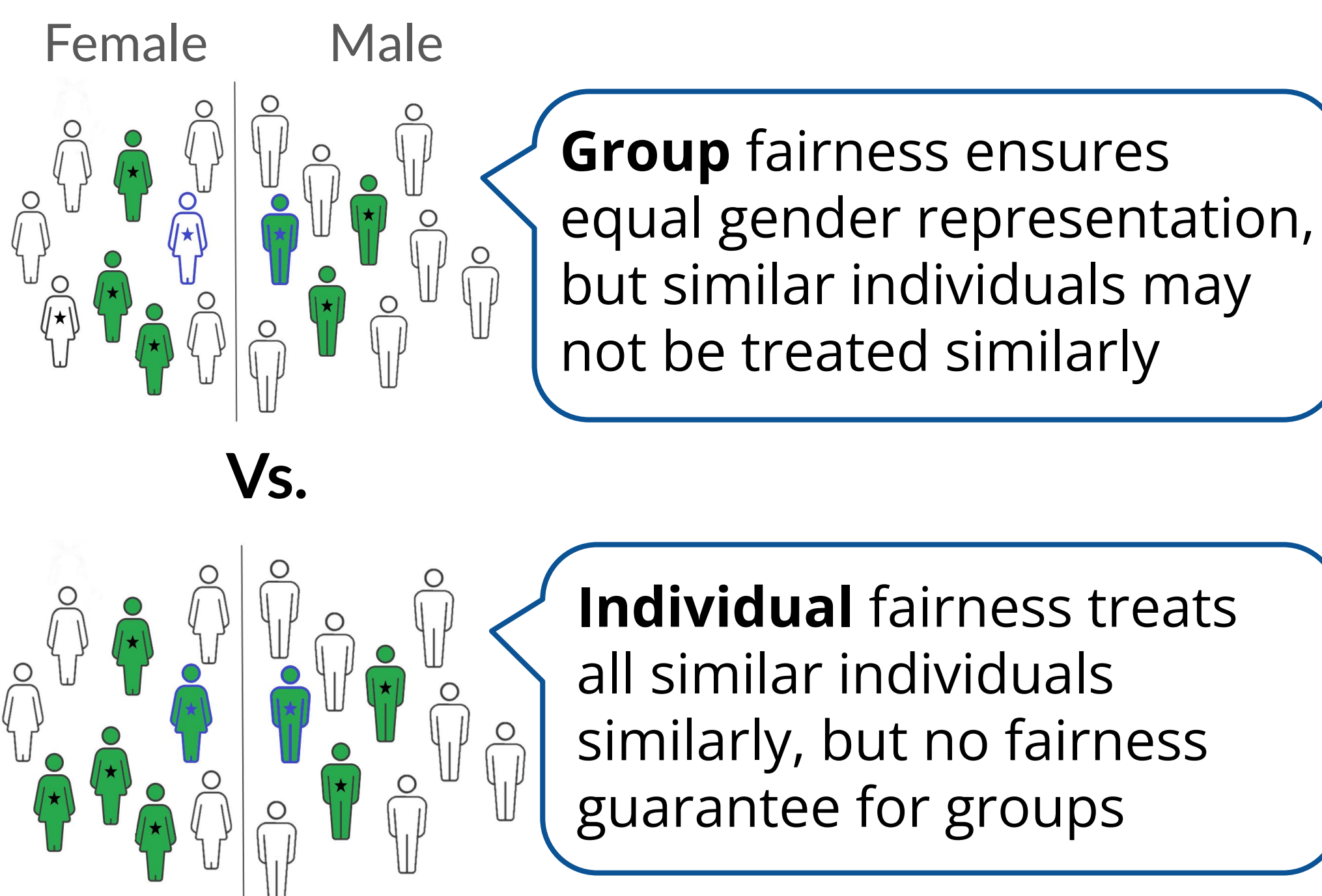
Classification algorithms can be modified to behave fairly, but present **two challenges**:

- How to choose an appropriate **definition of fairness**?
- Which **technique** to use to incorporate fairness into a classifier?

Our work addresses both challenges!

Fairness is subjective and application dependent

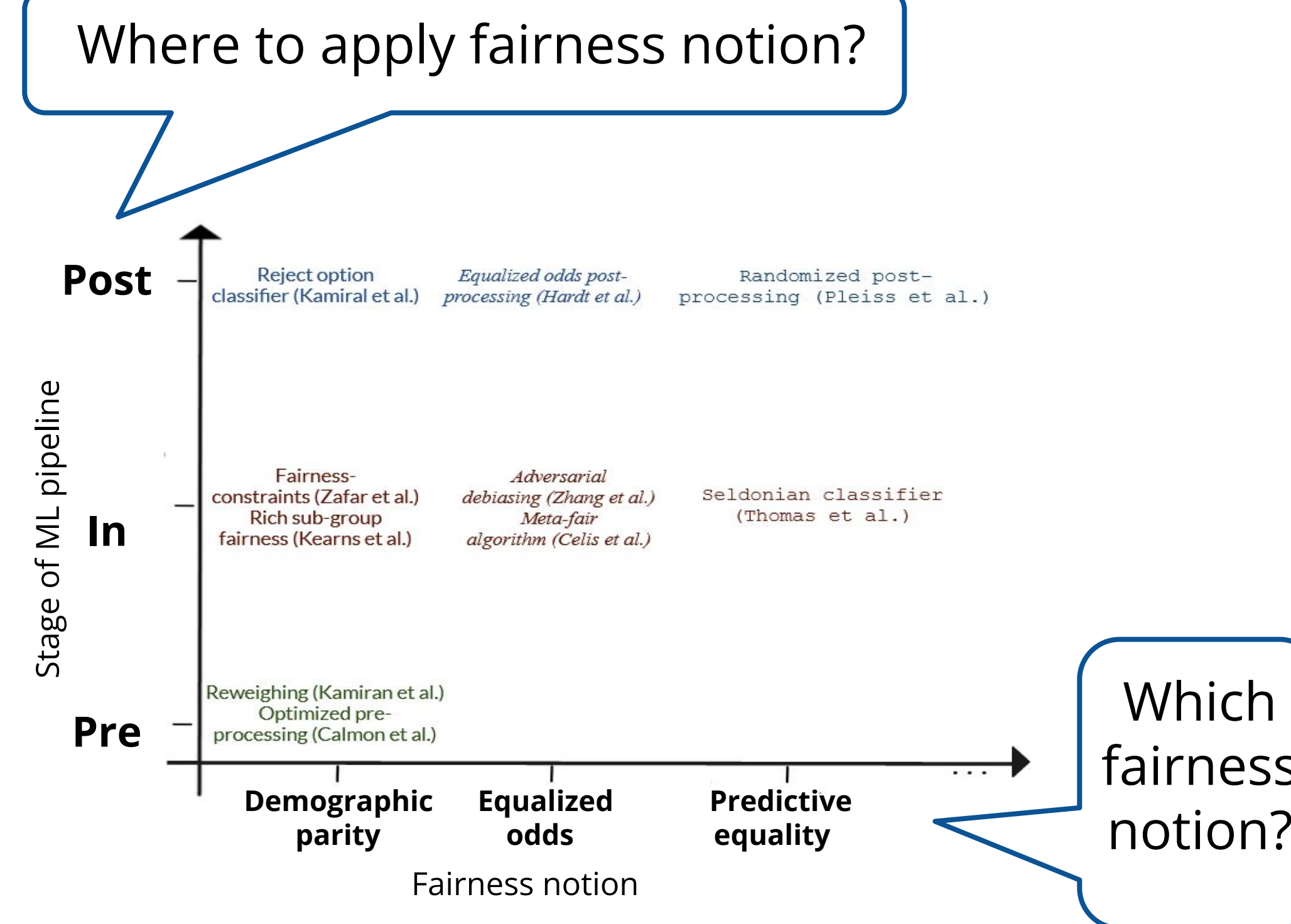
★: similar individuals
green: good outcome



Informative categorization of 34 fairness notions

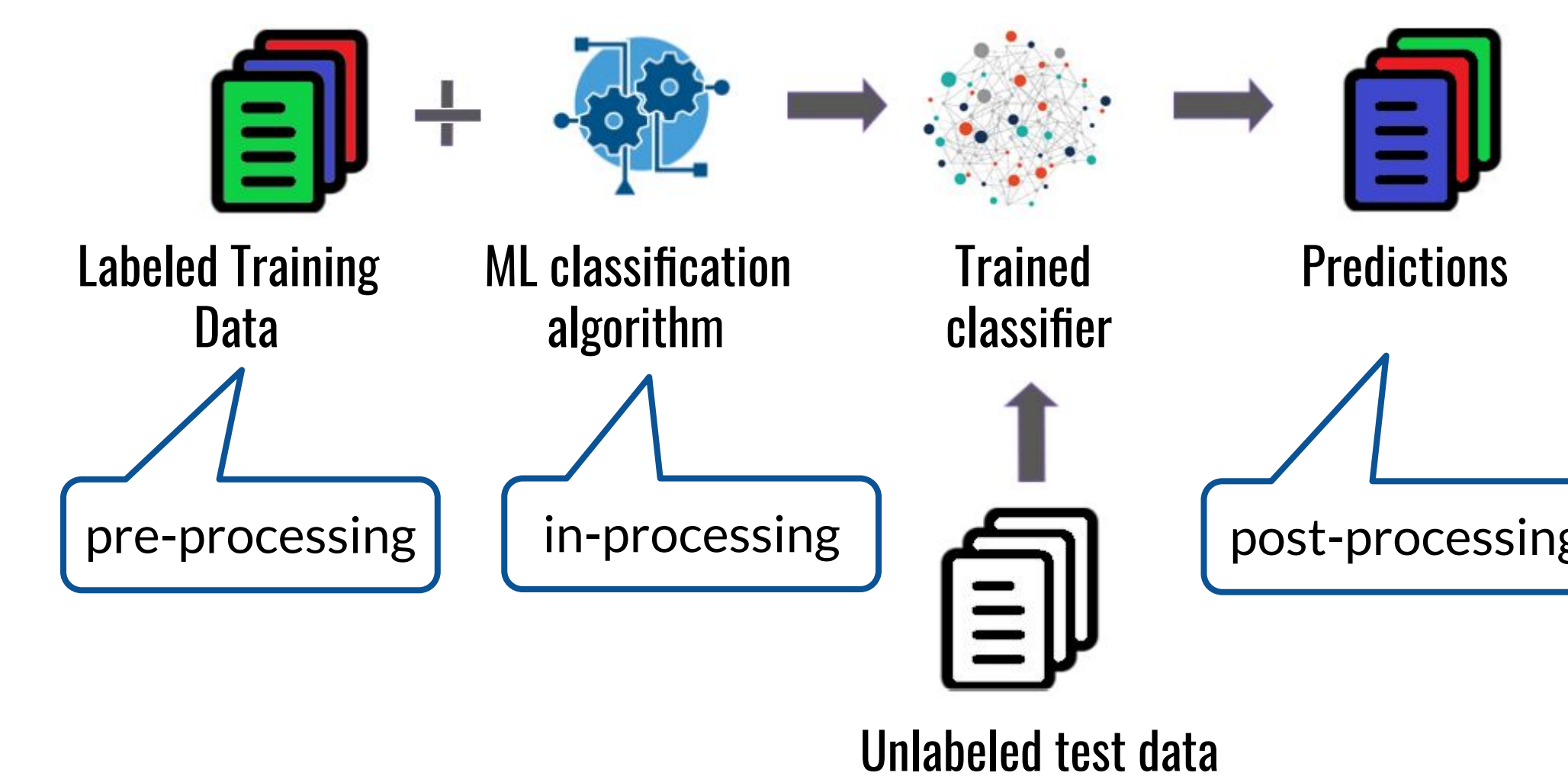
Fairness notion	Granularity		Causal hierarchy			Additional requirements		
	group	individual	observation	intervention	counterfactual	prediction	causality	resolving similarity
conditional statistical parity [24]	✓							
demographic parity [26]	✓							
intersectional fairness [32]	✓							
overall accuracy equality [10]		✓						
treatment equality [10]		✓						
equalized odds [41]		✓						
equal opportunity [41]		✓						
resilience to random bias [30]		✓						
preference-based fairness [96]		✓						
calibration [23]		✓						
calibration within groups [53]		✓						
positive class balance [53]		✓						
negative class balance [53]		✓						
individual discrimination [34]		✓						
metric multifairness [52]		✓						✓
proxy fairness [50]	✓							
total causal effect [72]	✓							
path-specific fairness [105]		✓						
unresolved discrimination [50]	✓							
interventional/justifiable fairness [80]	✓							
fair on average causal effect [49]	✓							
non-discrimination criterion [104]	✓							
equality of effort [42]	✓							
individual direct discrimination [103]		✓						

Fair classifiers vary across two dimensions



Three stages of ML pipeline to enforce fairness

Pre-processing modifies data **before** training
In-processing constrains classifier objective **during** training
Post-processing modifies prediction **after** training



Empirical investigation of fair classifiers

Approaches: 18 fair techniques

Datasets: 3 real-world datasets

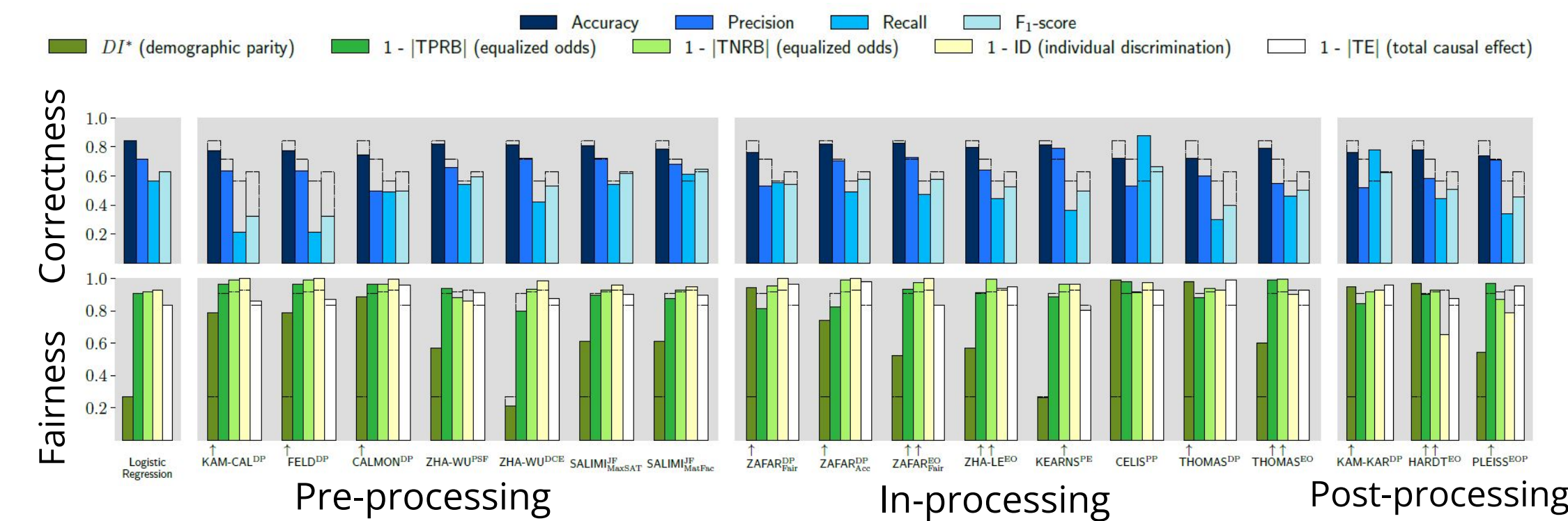
Biases: gender and racial

Evaluation axes: 4 areas

- Correctness- fairness **tradeoff**
- Scalability** issues
- Effect of training **data errors**
- Sensitivity to the choice of **ML model**

Correctness-fairness tradeoff

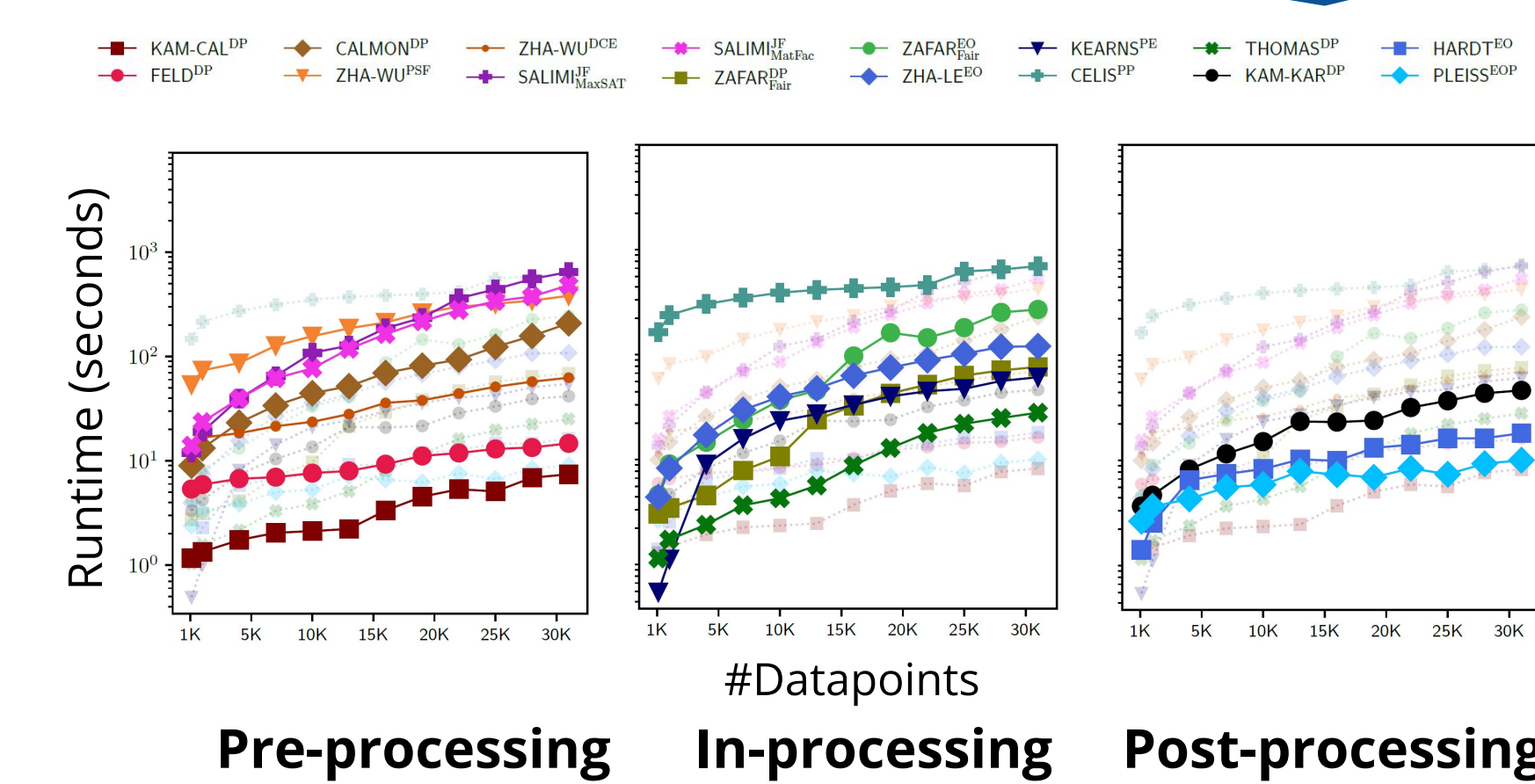
Fair approaches **trade** accuracy for fairness



Bigger compromise in accuracy when target fairness is **low** in fairness-unaware setting

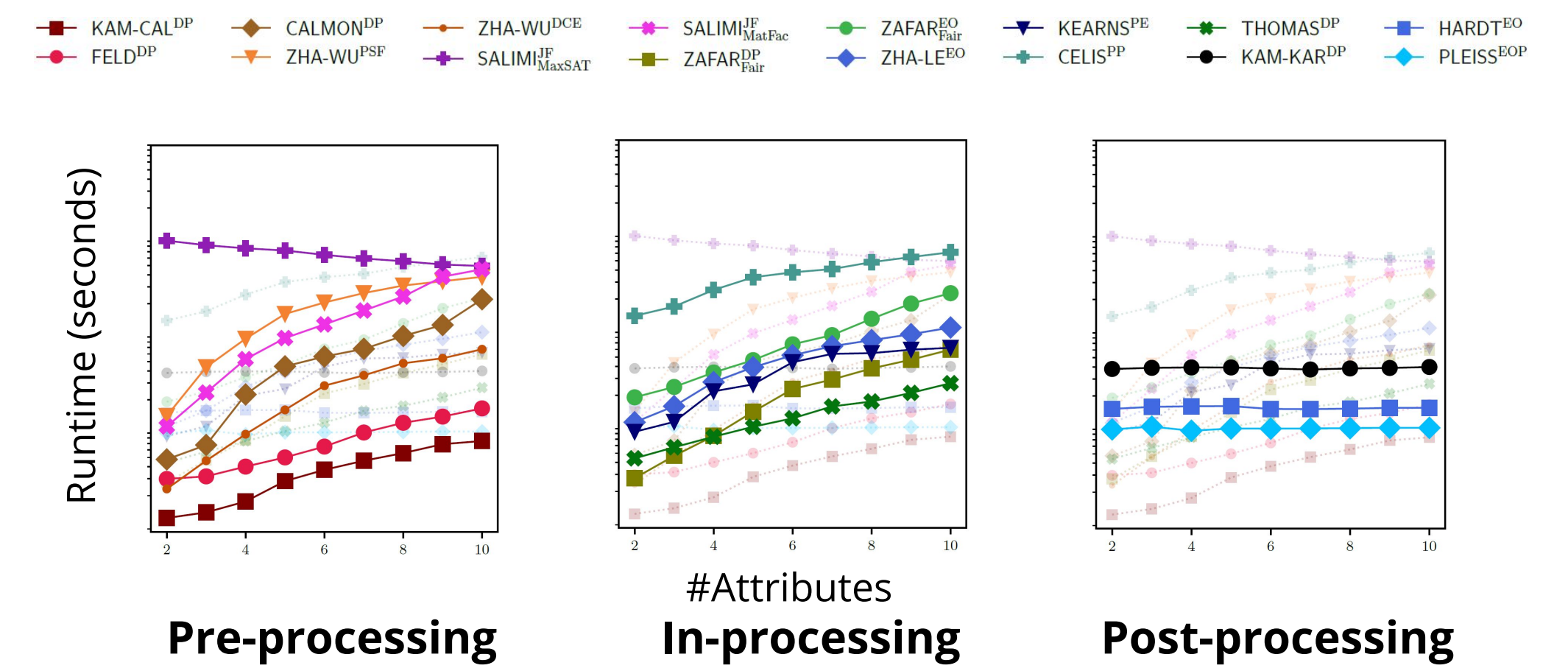
Efficiency gap

Post-processing is most **efficient** due to inherent simplicity



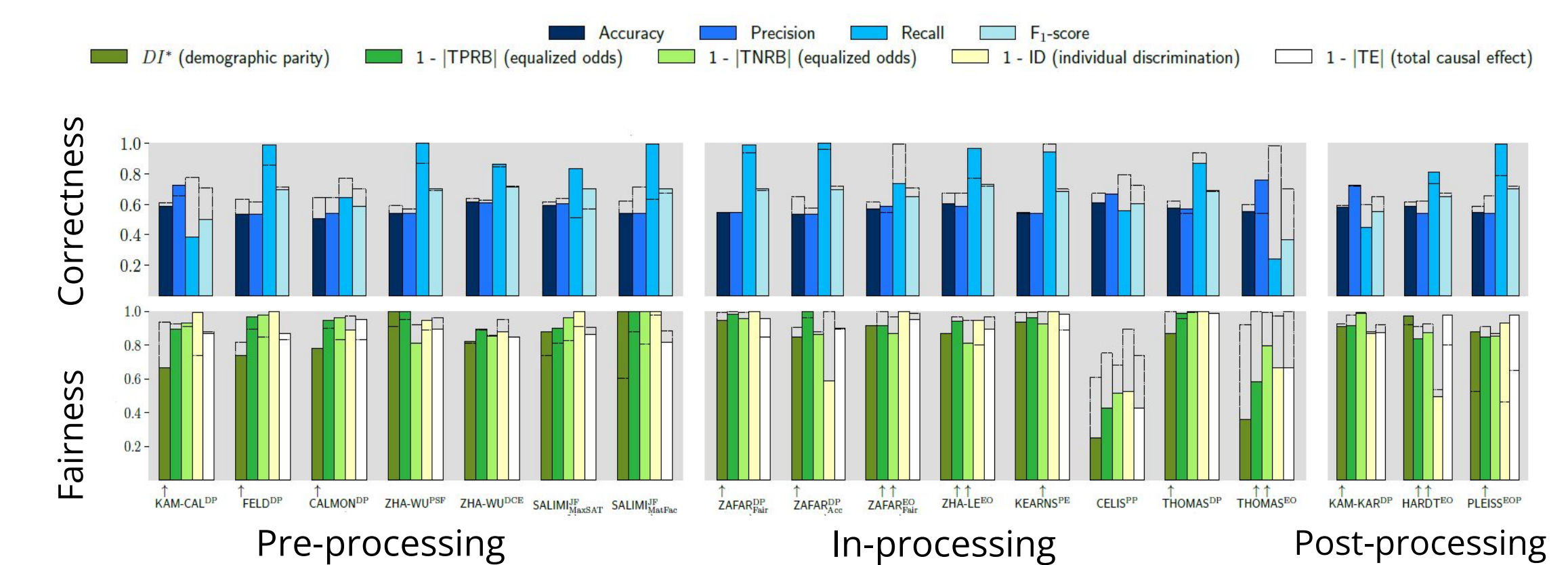
Complex mechanisms can **lower efficiency**: common in **pre-** and **in-processing**

Scalability issues



Pre-processing scales worse with increasing #attributes than with increasing #datapoints

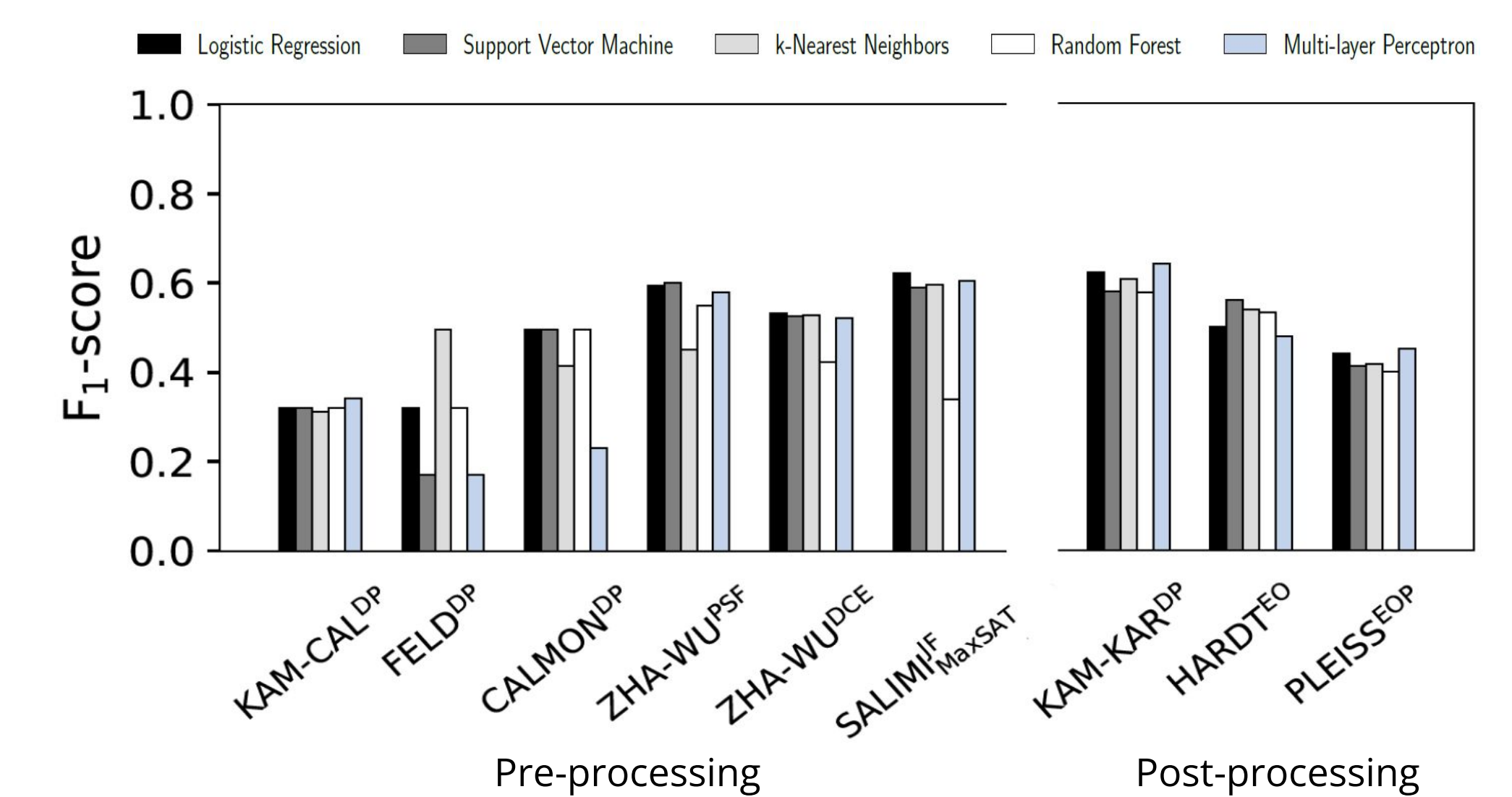
What happens if training data has errors?



Pre- and in-processing can **fail** to build fair models when training data contains **errors**

Sensitivity to the choice of ML model

Post-processing is more **stable**



Pre-processing can **fluctuate** depending on classifier model