

Background

When building classification systems with demographic fairness considerations, there are **two objectives to satisfy**:
1) maximizing utility for the specific task and **2) ensuring fairness** w.r.t. a known demographic attribute. These objectives often compete, so optimizing both can lead to a trade-off between utility and fairness.

Questions We Answer

- What** are the optimal trade-offs between utility and fairness?
- How** can we numerically quantify these trade-offs from data for a desired prediction task and demographic attribute of interest?

Trade-Offs Definitions

Definition 1. Data Space Trade-Off (DST)

$$f_{\lambda}^{DST} := \arg \inf_{f \in \mathcal{H}_X} \left\{ (1 - \lambda) \inf_{g_Y \in \mathcal{H}_Y} \mathbb{E}_{X,Y} [\mathcal{L}_Y(g_Y(f(X)), Y)] + \lambda \text{Dep}(f(X), S|Y=y) \right\}, \quad 0 \leq \lambda < 1$$

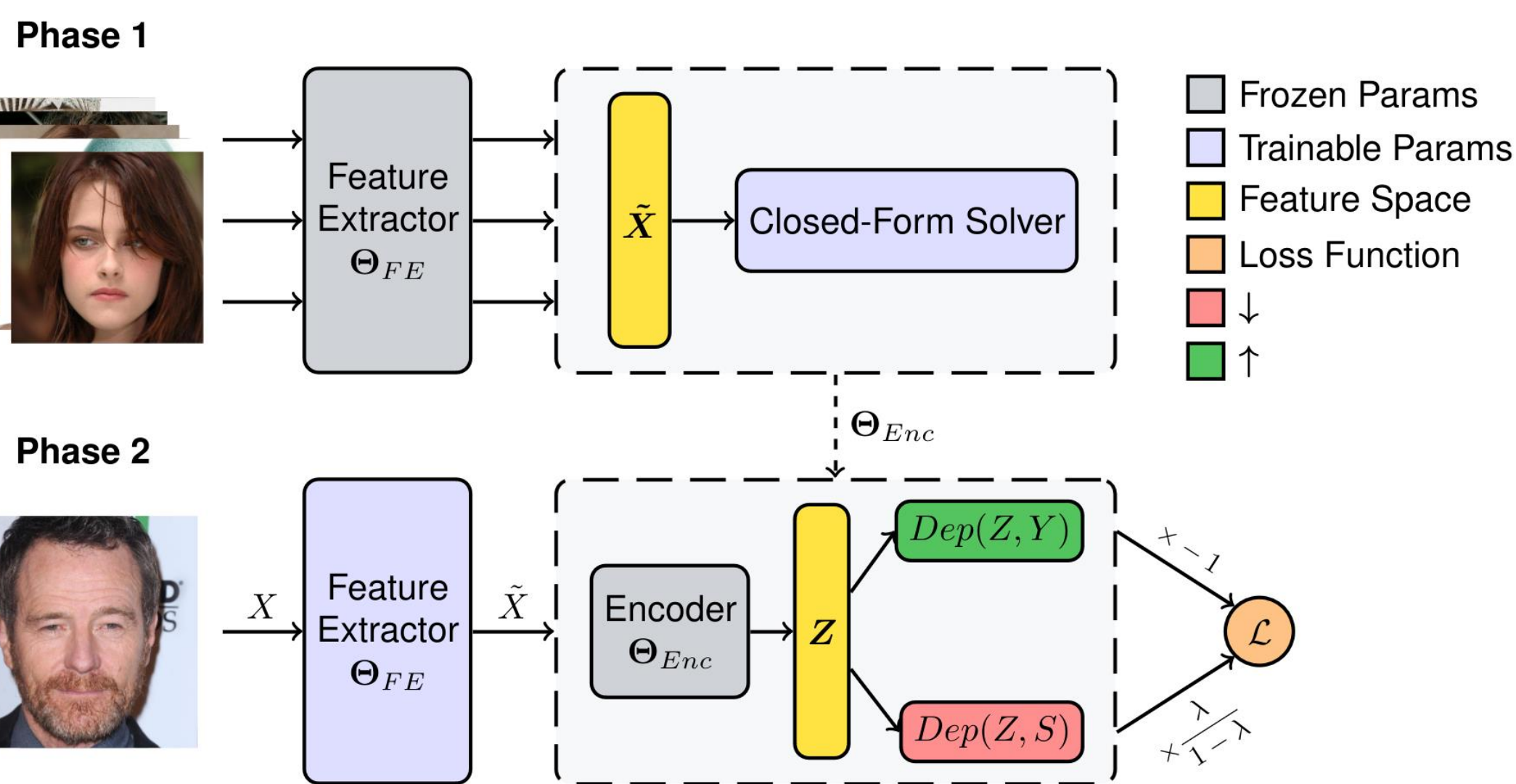
Definition 2. Label Space Trade-Off (LST)

$$Z_{\lambda}^{LST} := \arg \inf_{Z \in L^2} \left\{ (1 - \lambda) \inf_{g_Y \in \mathcal{H}_Y} \mathbb{E}_Y [\mathcal{L}_Y(g_Y(Z), Y)] + \lambda \text{Dep}(Z, S|Y=y) \right\}, \quad 0 \leq \lambda < 1$$

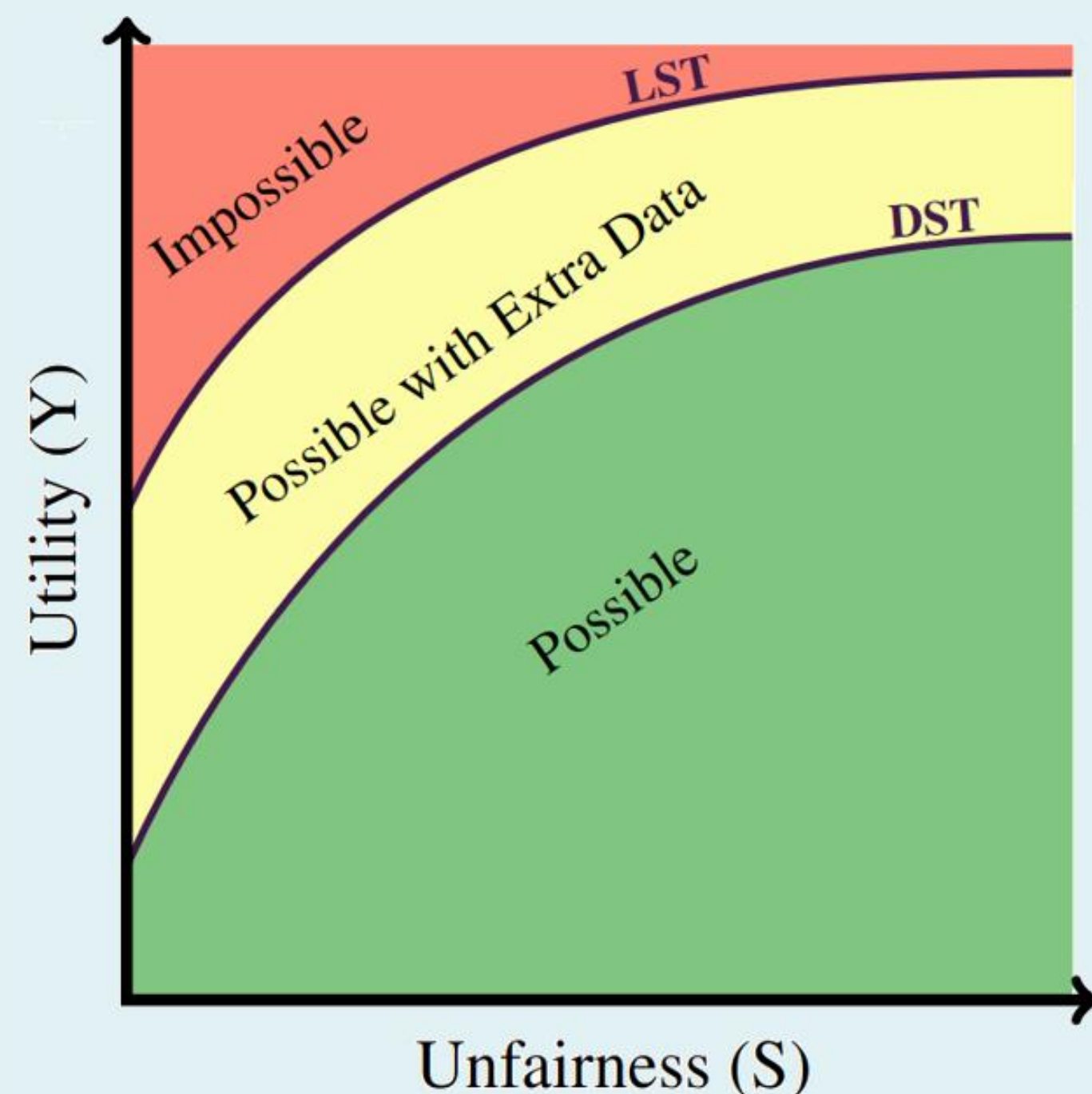
Fairness Criteria

- Demographic Parity (**DP**)
- Equalized Opportunity (**EO**)
- Equality Of Odds (**EOO**)

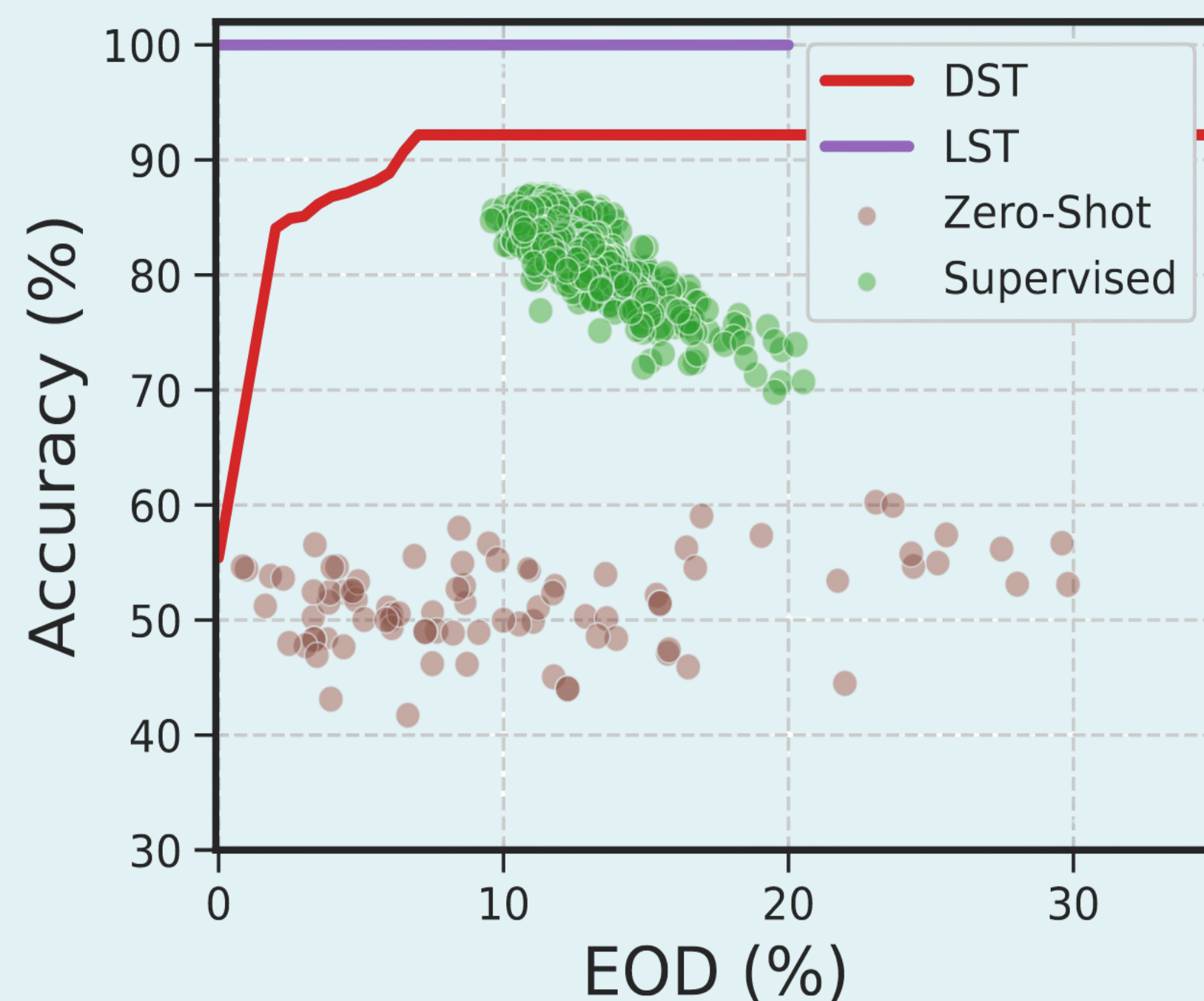
U-FaTE



Utility-Fairness Trade-Offs and How to Find Them



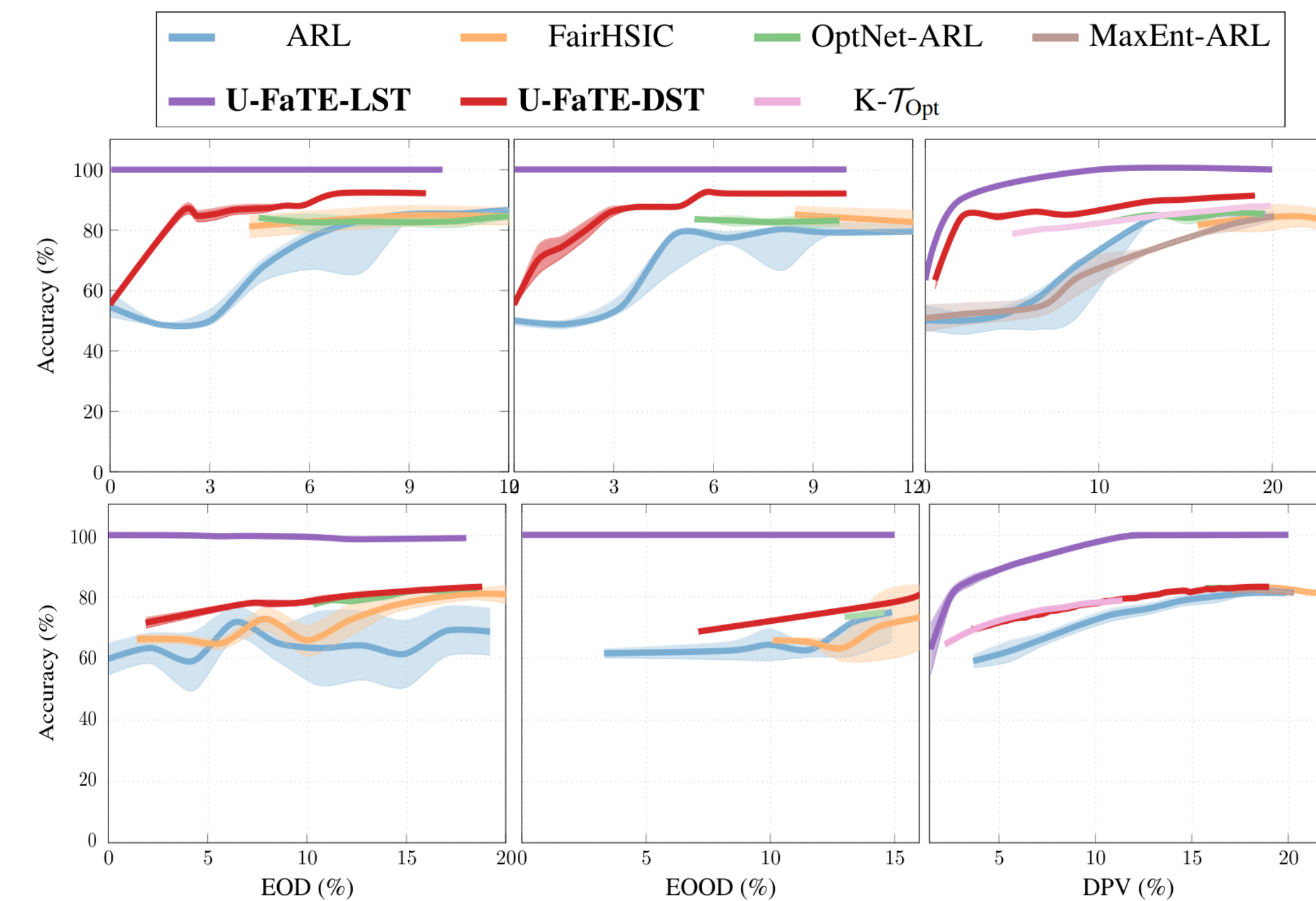
Data Space Trade-Off (**DST**) and Label Space Trade-Off (**LST**) divide the utility (e.g., accuracy) versus fairness space into **three regions**.



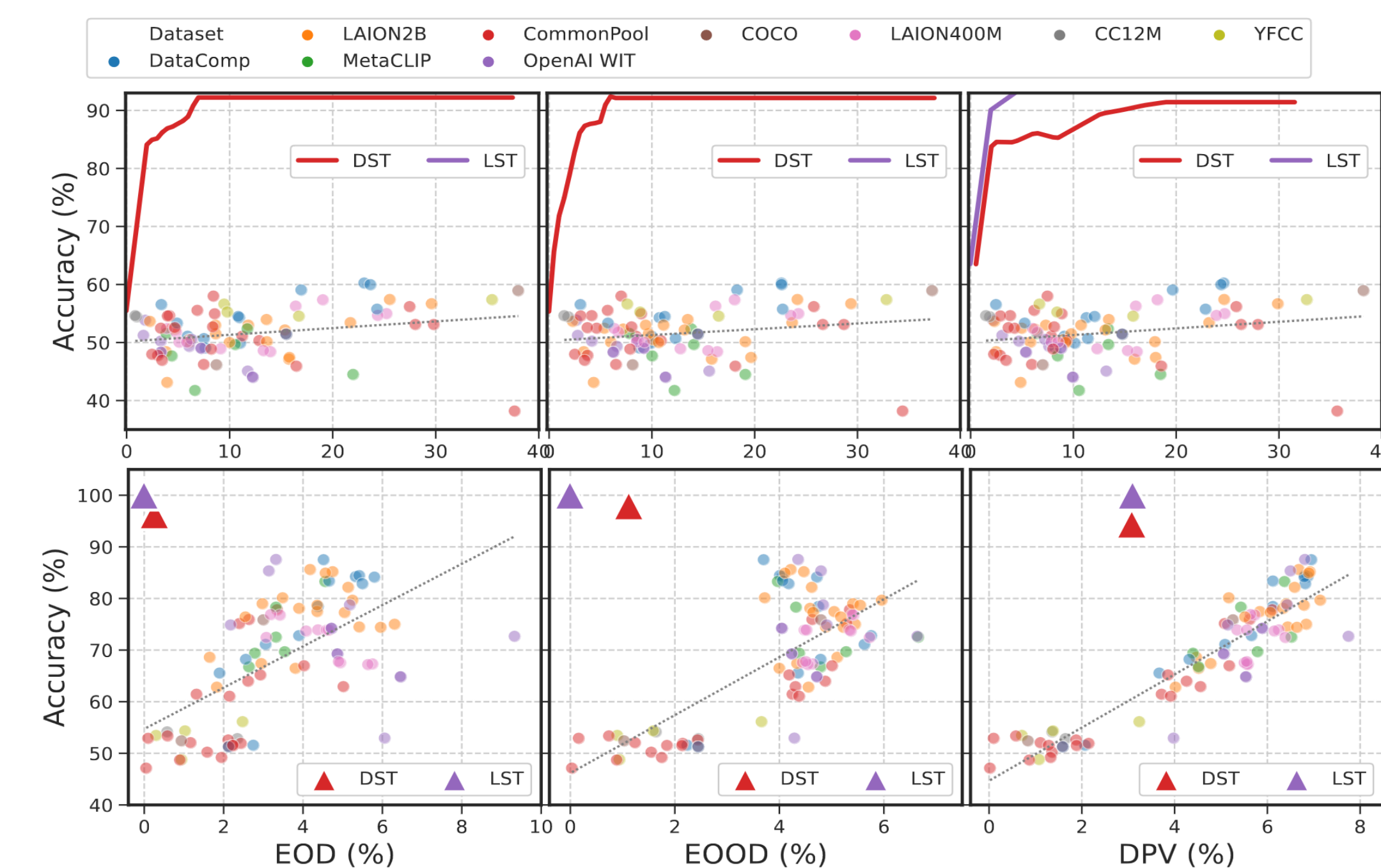
We empirically estimate **DST** and **LST** on **CelebA** and evaluate the utility (high cheekbones) and fairness (**gender & age**) of over **100 zero-shot** and **900 supervised** models.



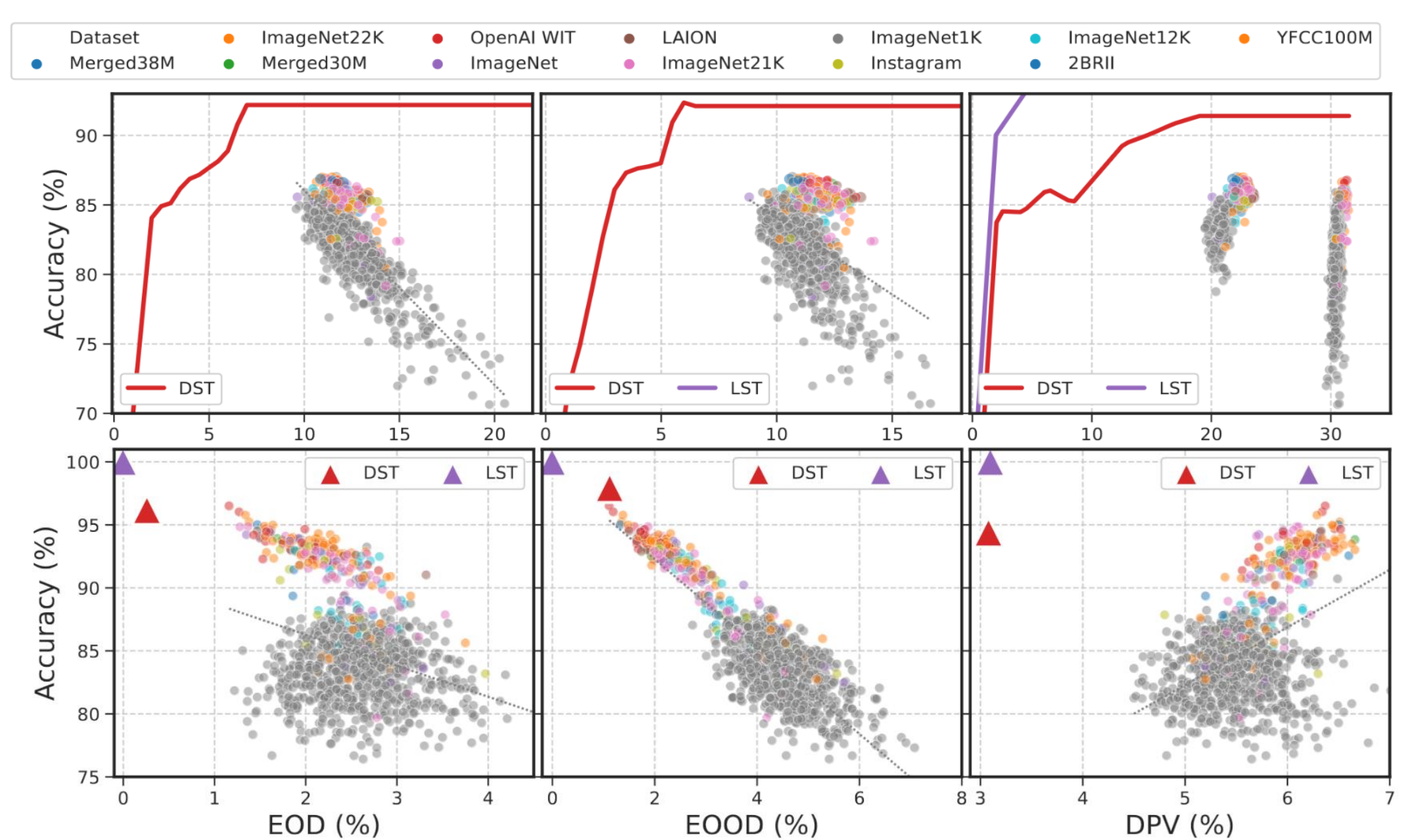
How do FRL Methods Compare?



How fair are CLIP Models?



How fair are pre-trained image models?



Sepehr Dehdashtian,
Bashir Sadeghi,
Vishnu Naresh Boddeti