# Superhuman Fairness

**Omid Memarrast**, Linh Vu, Brian D. Ziebart

**University of Illinois Chicago**

UIC COMPUTER SCIENCE

# Motivation

Defining desired fairness trade-offs precisely is difficult
- **Multiple fairness metrics** [dp, eqodds, eqopp, prp, ...]

**A new perspective: Multiple stakeholders**

- with **different notions of fairness** and **desired performance-fairness trade-offs**
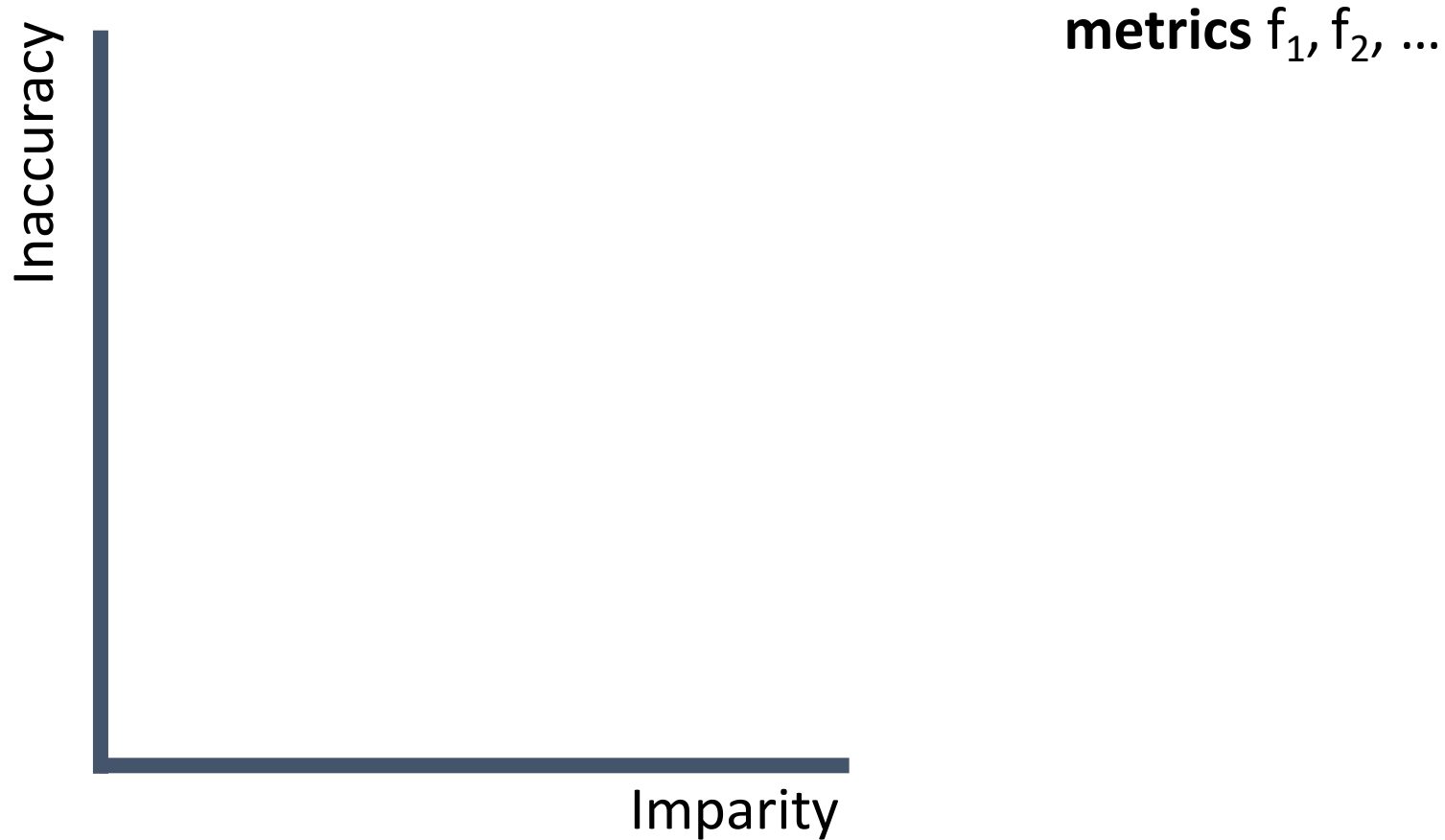
**Example**

- Admission: [CS department, Civil department, ..]

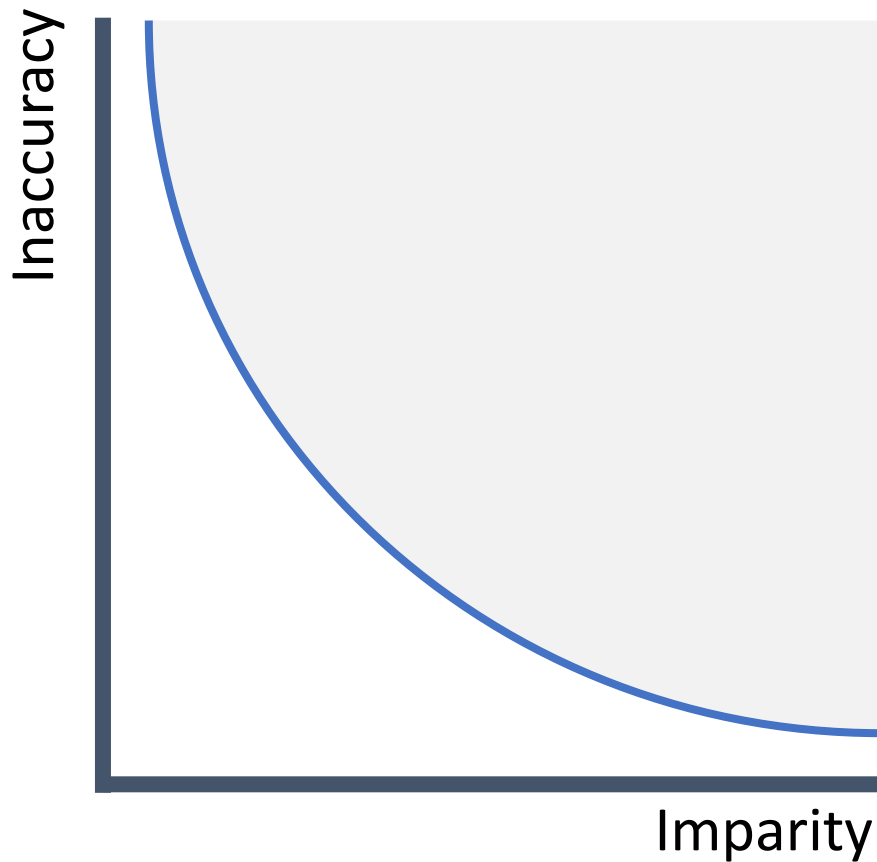- Each department: Their own perception of fairness

**Solution:**

Instead of optimal fairness, **outperform humans** across many metrics
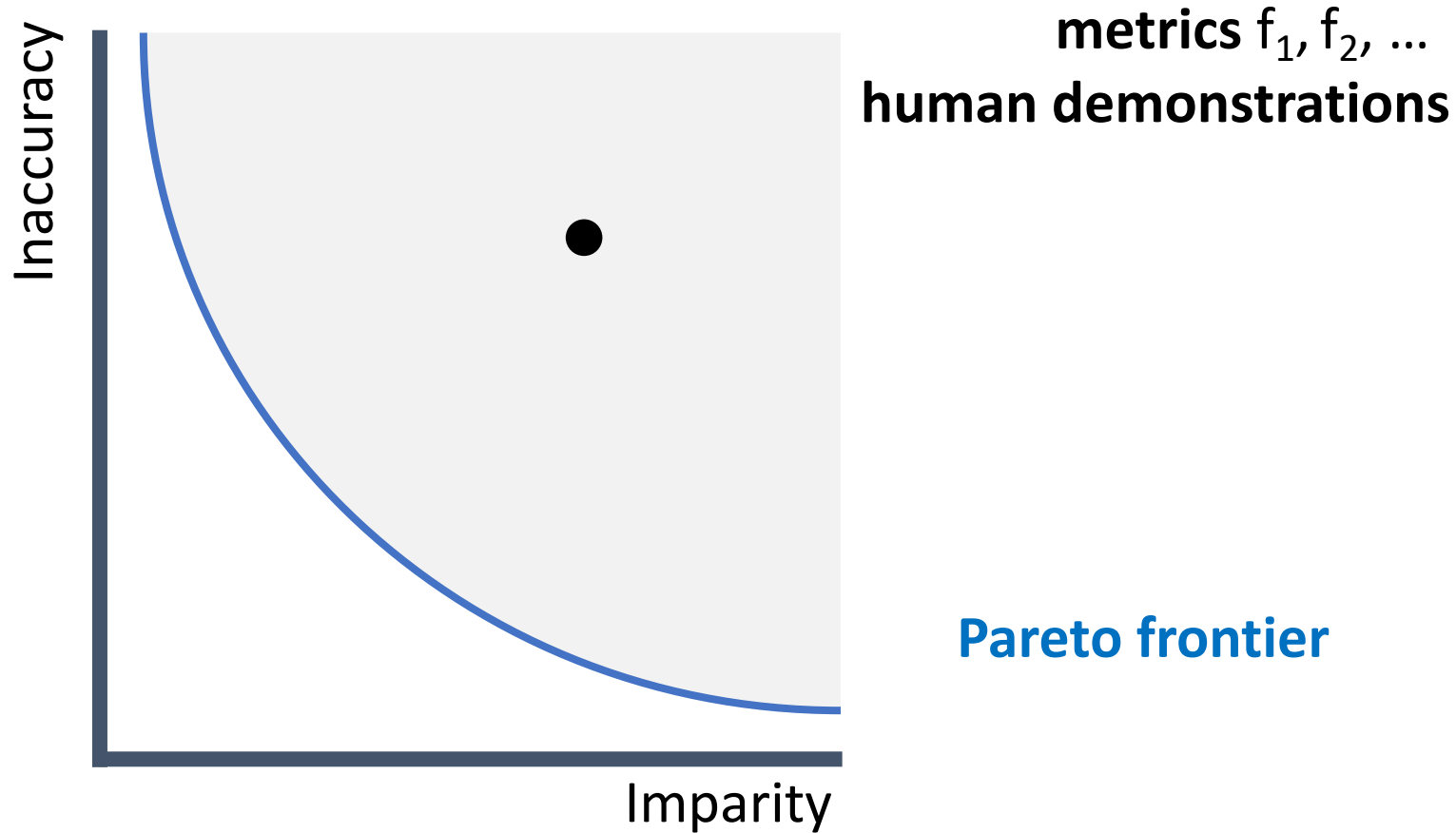
# Defining Superhuman Behavior



**metrics** $f_1$, $f_2$, …
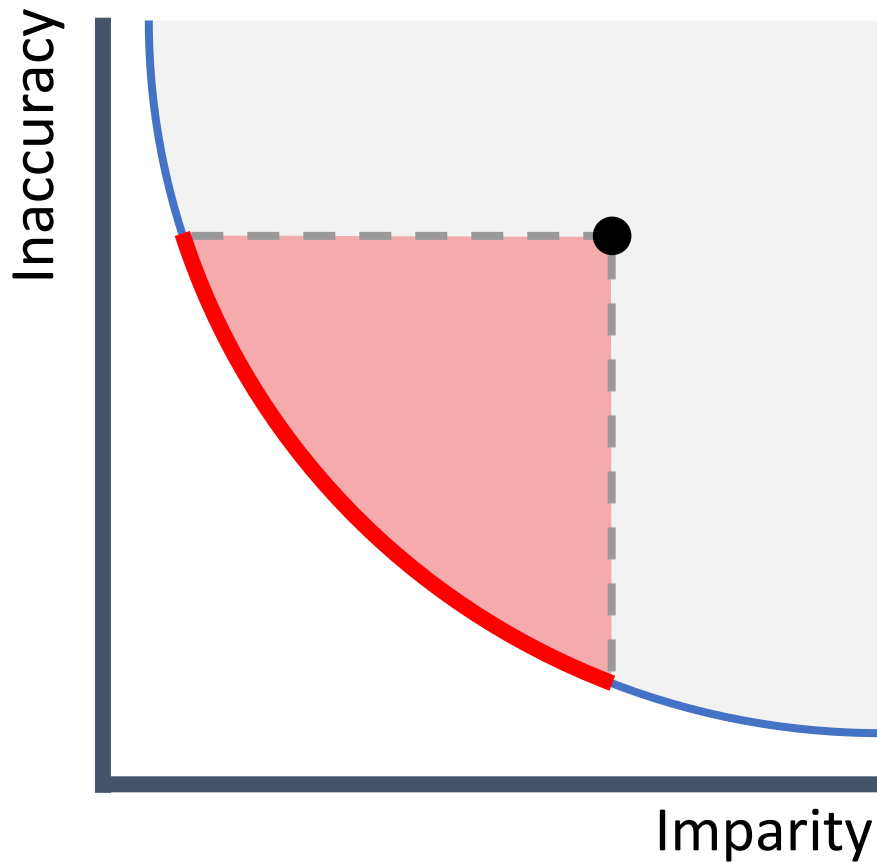
# Defining Superhuman Behavior



**metrics** $f_1$, $f_2$, ...

**Pareto frontier**

# Defining Superhuman Behavior



metrics $f_1$, $f_2$, …
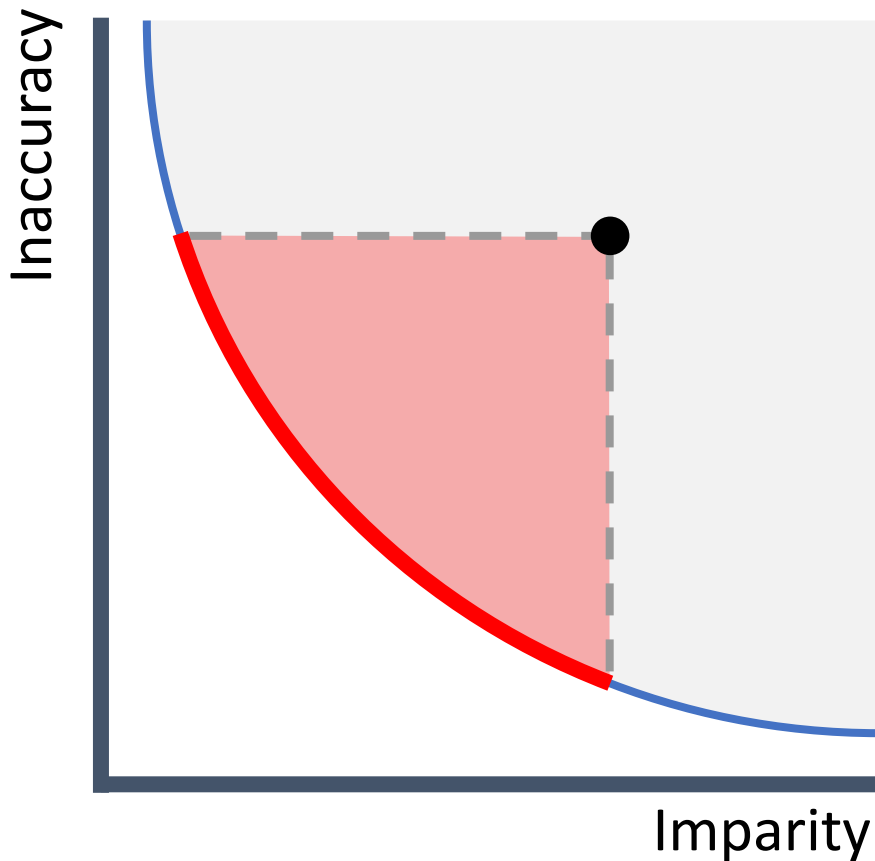human demonstrations

Pareto frontier

# Defining Superhuman Behavior



A **policy** is **superhuman** if it has smaller **metrics** $f_1$, $f_2$, ... for all **human demonstrations**

**Pareto frontier**
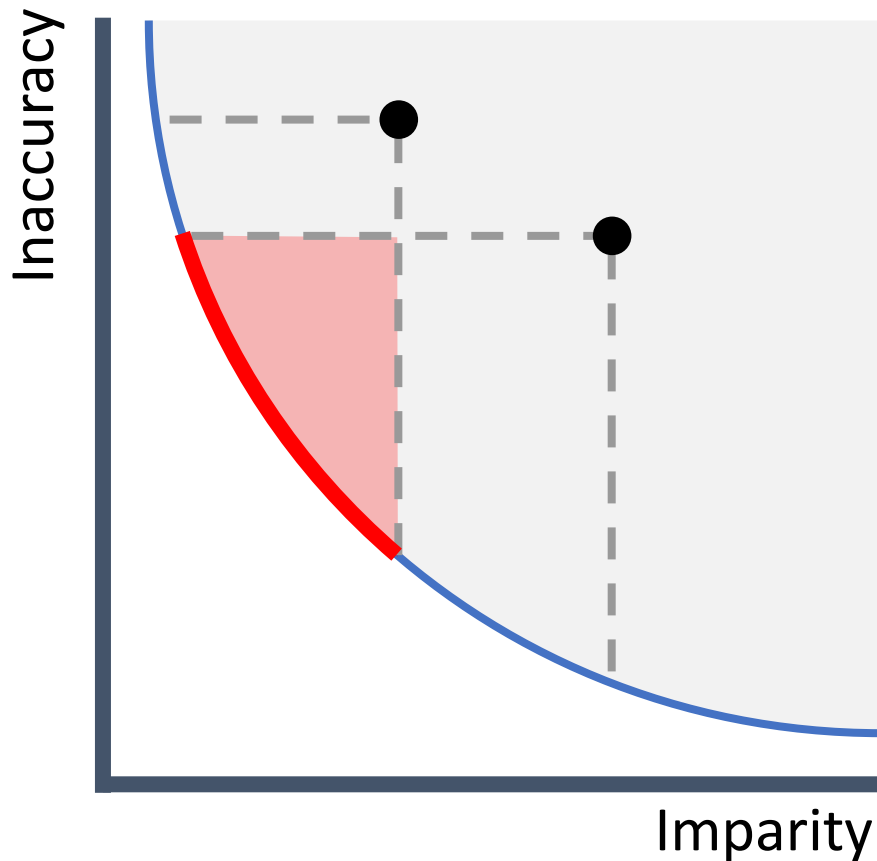
# Defining Superhuman Behavior



A **policy** is **superhuman** if it has smaller **metrics** $f_1$, $f_2$, … for all **human demonstrations**

Guarantees <u>lower cost</u> than demonstration costs for family of additive trade-offs

**Pareto frontier**
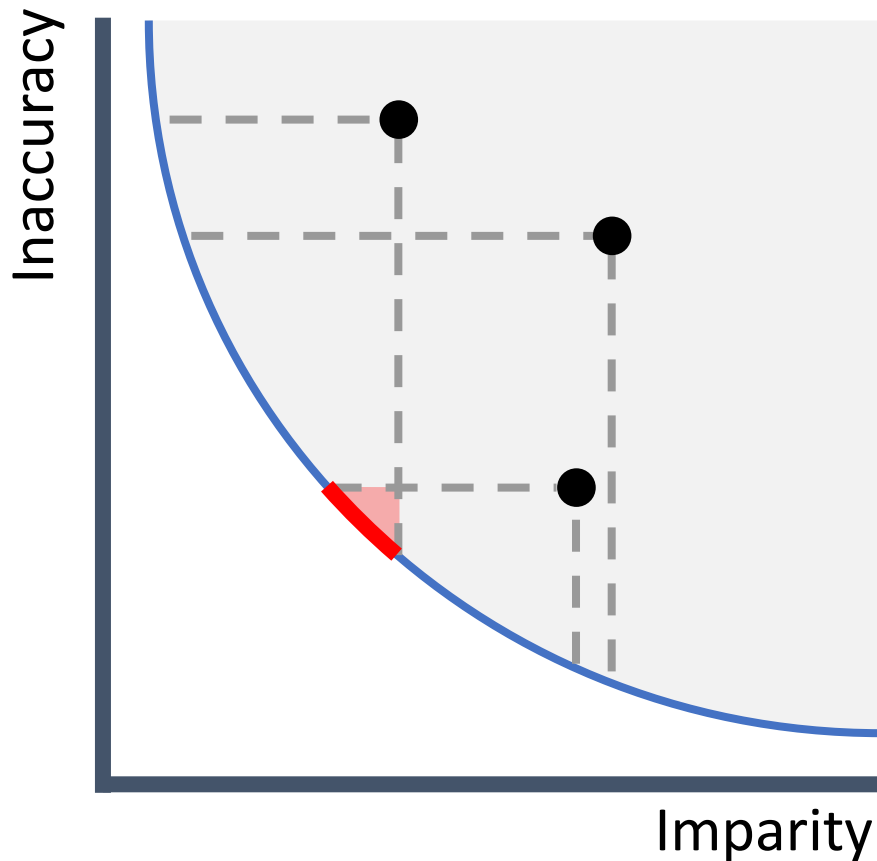
# Defining Superhuman Behavior



A **policy** is **superhuman** if it has smaller **metrics** $f_1$, $f_2$, ... for all **human demonstrations**

Guarantees <u>lower cost</u> than demonstration costs for family of additive trade-offs

Set of **superhuman policies** on the **Pareto frontier** shrinks as demonstrations grow

# Defining Superhuman Behavior

A **policy** is **superhuman** if it has smaller **metrics** $f_1$, $f_2$, … for all **human demonstrations**

Guarantees <u>lower cost</u> than demonstration costs for family of additive trade-offs

Set of **superhuman policies** on the **Pareto frontier** shrinks as demonstrations grow
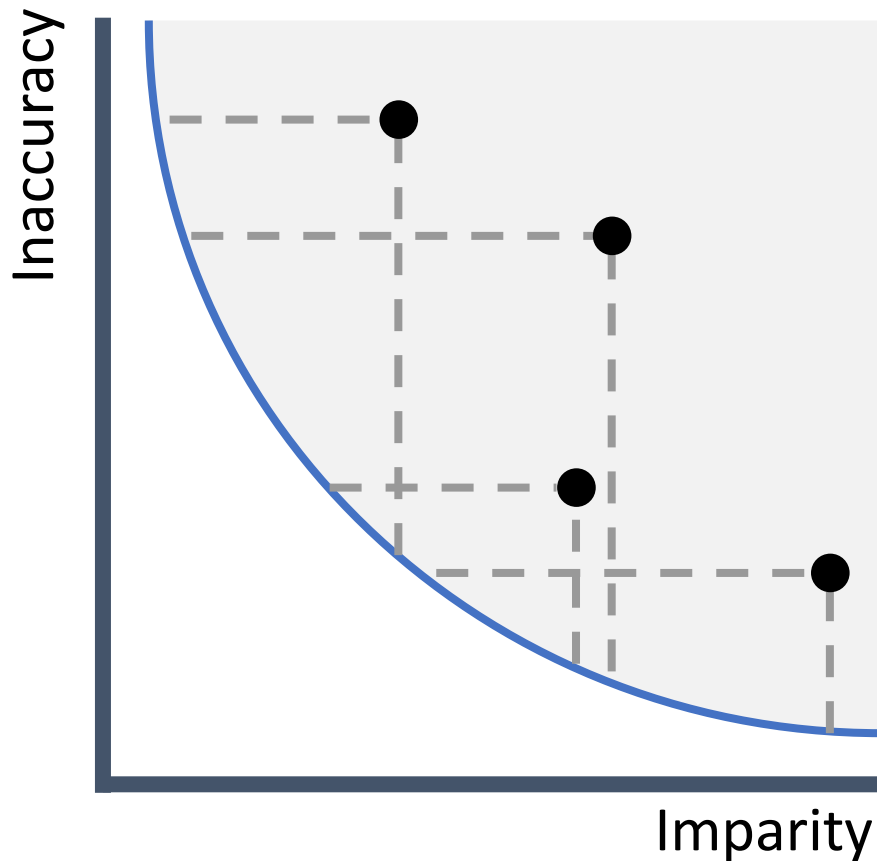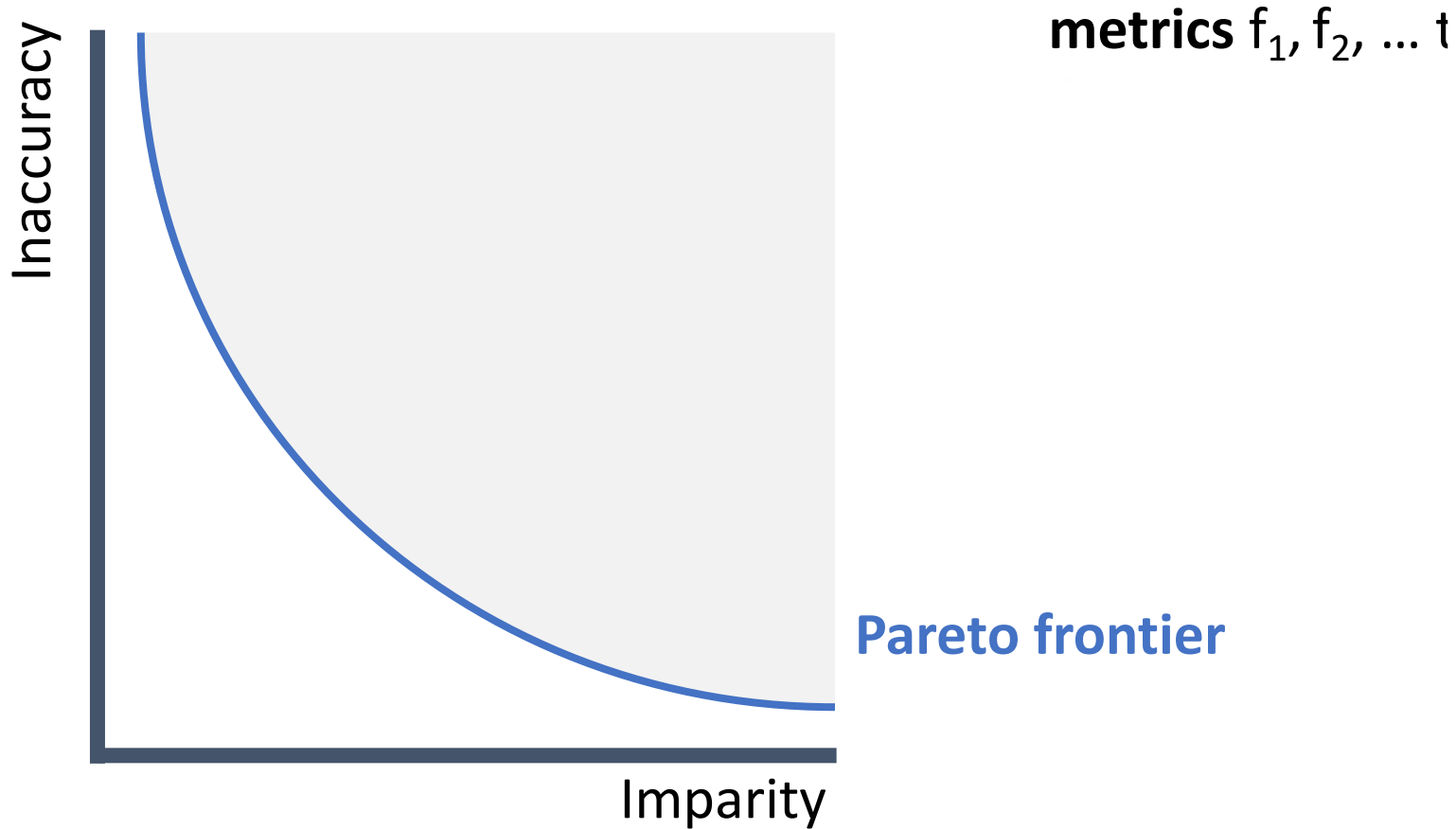
# Defining Superhuman Behavior



A **policy** is **superhuman** if it has smaller **metrics** $f_1$, $f_2$, … for all **human demonstrations**

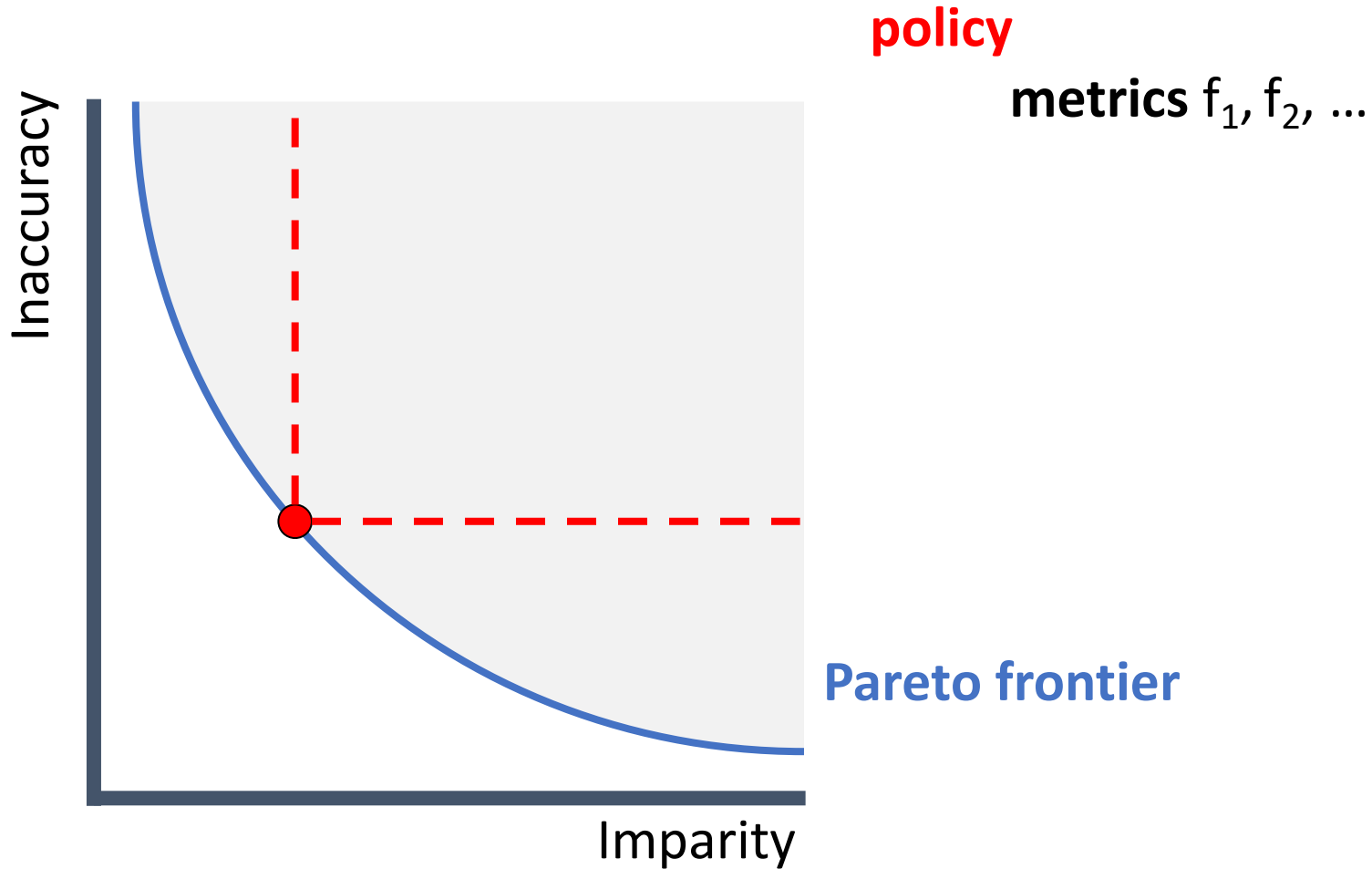Guarantees <u>lower cost</u> than demonstration costs for family of additive trade-offs

Set of **superhuman policies** on the **Pareto frontier** shrinks as demonstrations grow
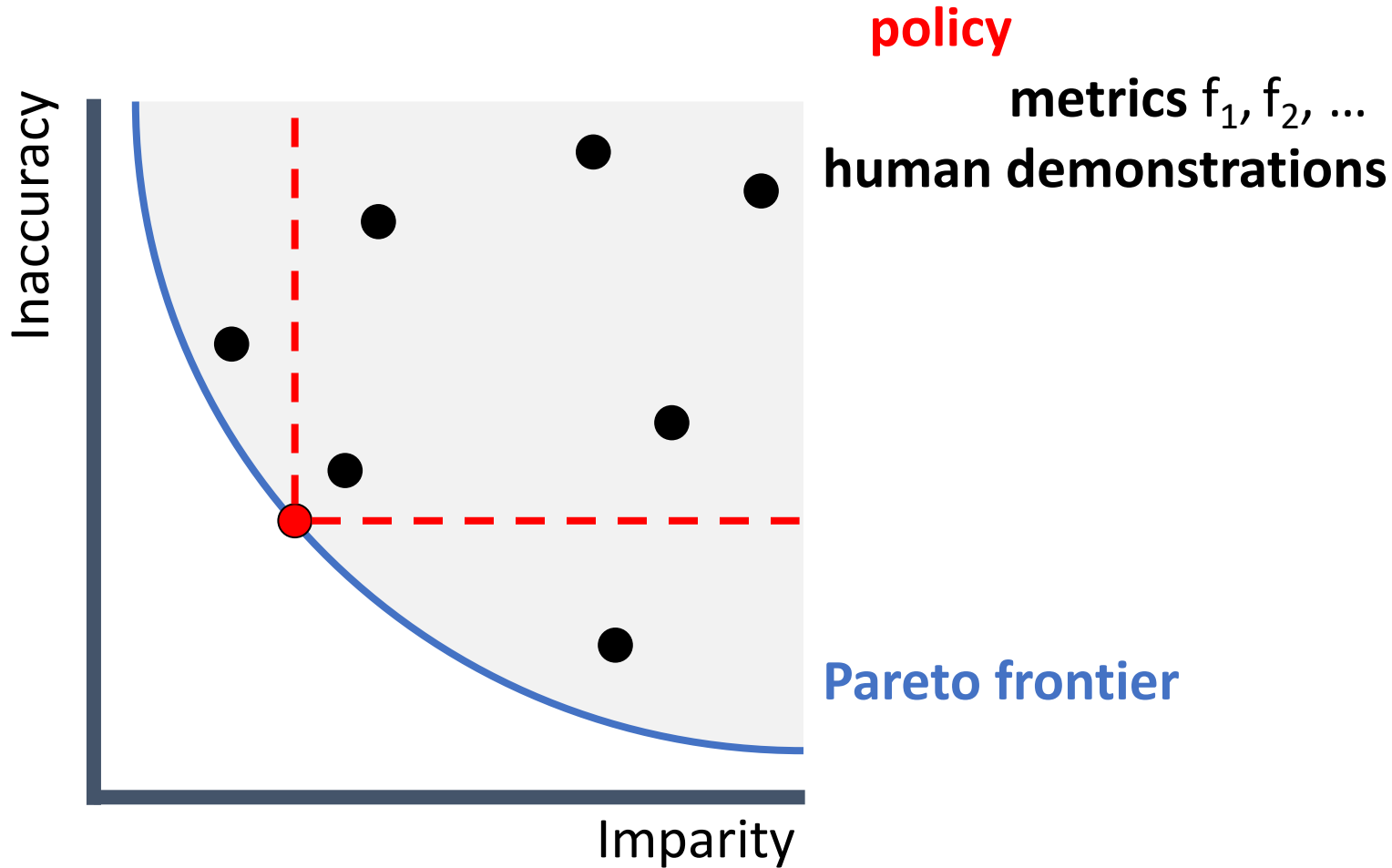
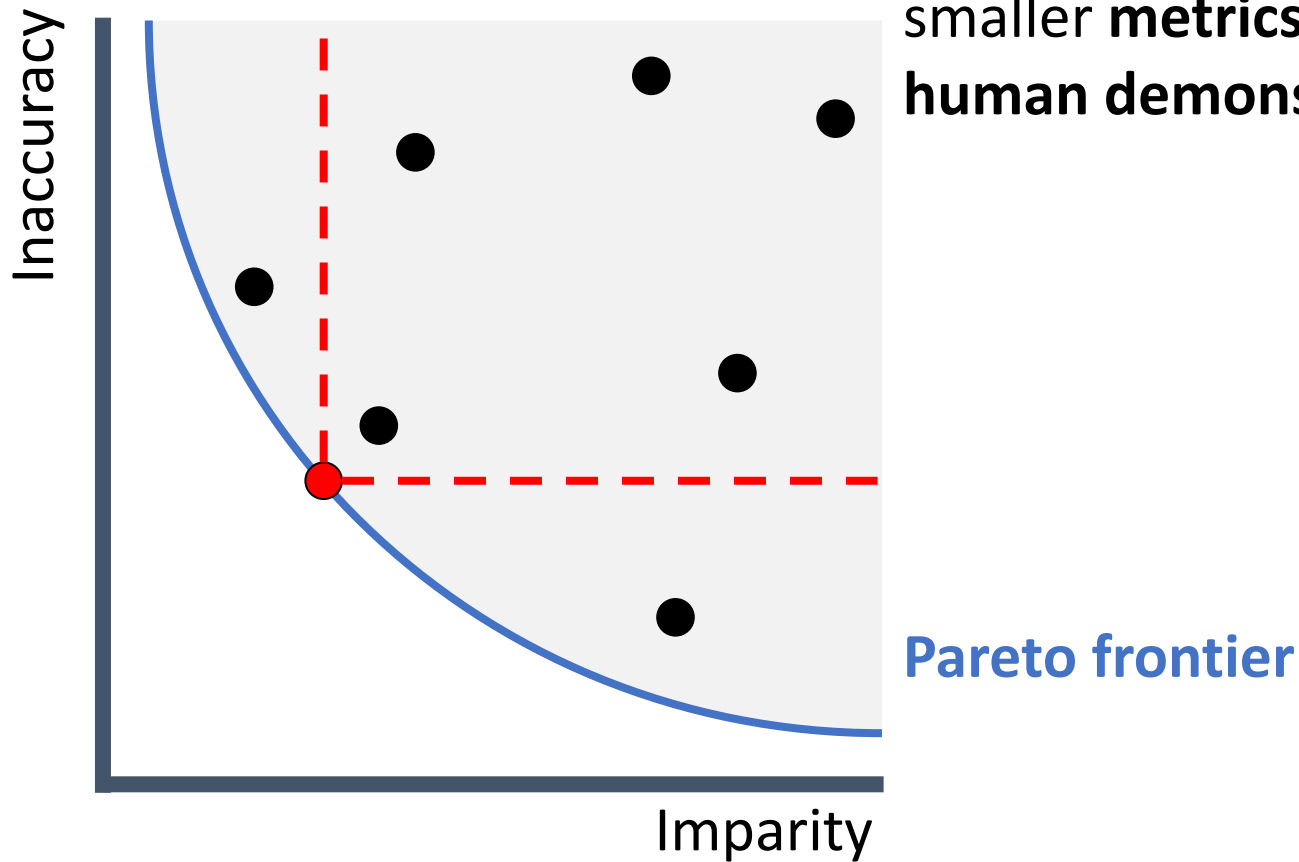Can become empty!

# Superhuman Percentile & Subdominance



**metrics** $f_1, f_2, \ldots t$

# Superhuman Percentile & Subdominance



policy

**metrics** $f_1$, $f_2$, …

Inaccuracy

Imparity

Pareto frontier

# Superhuman Percentile & Subdominance

# Superhuman Percentile & Subdominance



A **policy** is **γ-superhuman** if it has smaller **metrics** $f_1$, $f_2$, … than γ% of **human demonstrations**

# Superhuman Percentile & Subdominance



A **policy** is **γ-superhuman** if it has smaller **metrics** $f_1$, $f_2$, ... than γ% of **human demonstrations**
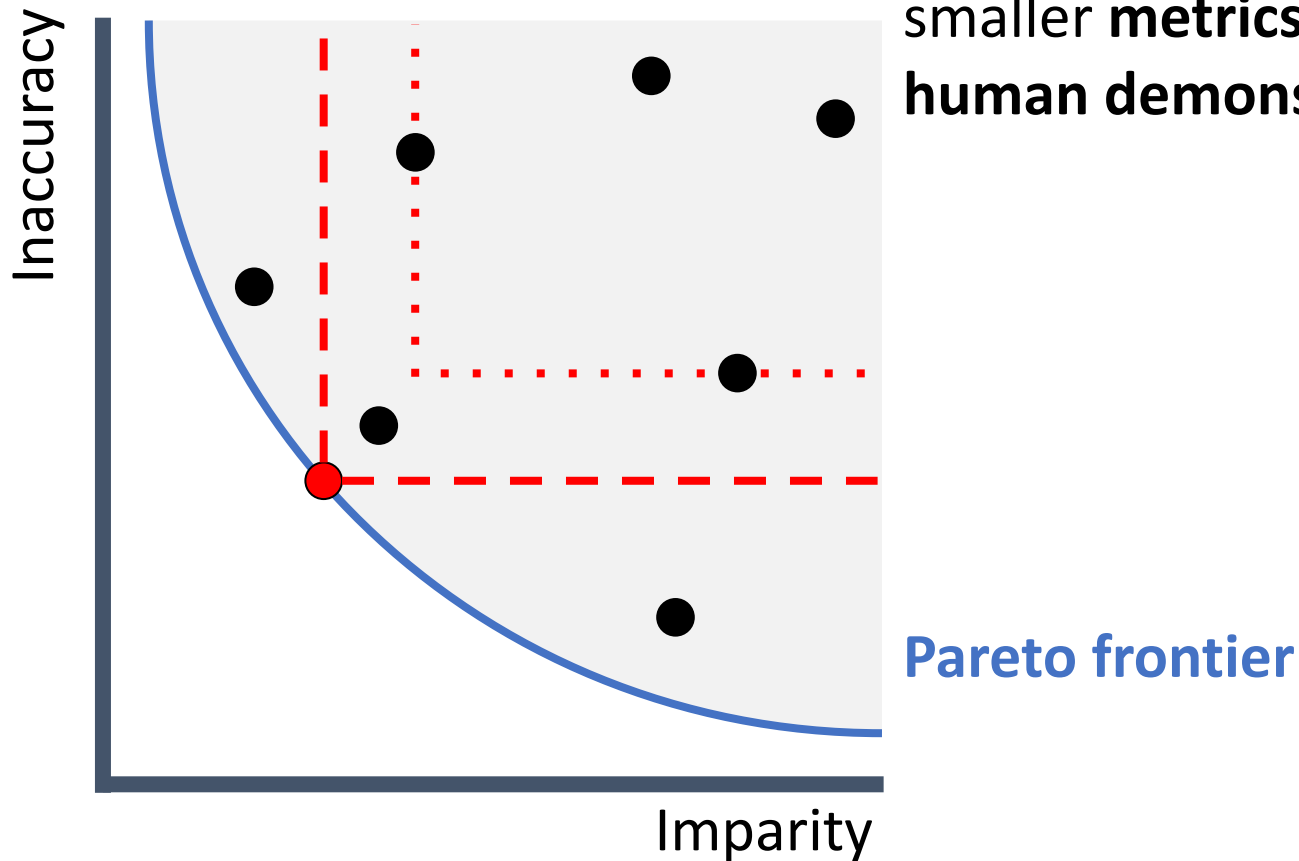
# Superhuman Percentile & Subdominance



A **policy** is **γ-superhuman** if it has smaller **metrics** $f_1$, $f_2$, … than γ% of **human demonstrations**

**Subdominance** measures how far a policy is from superhuman by some **margins**

# Superhuman Percentile & Subdominance



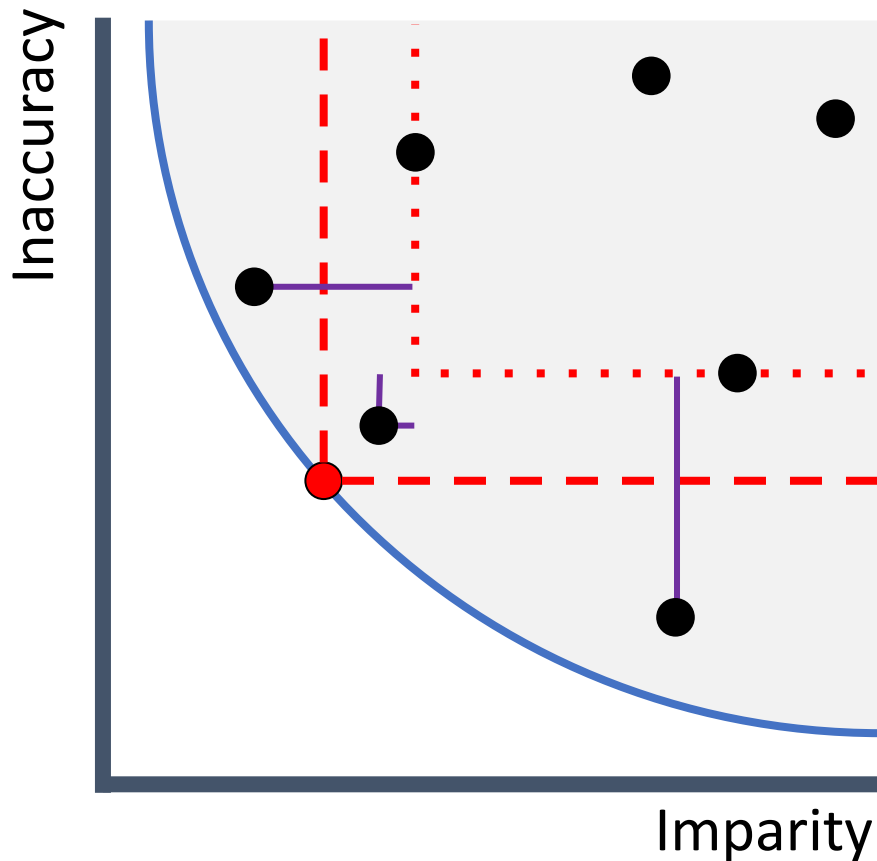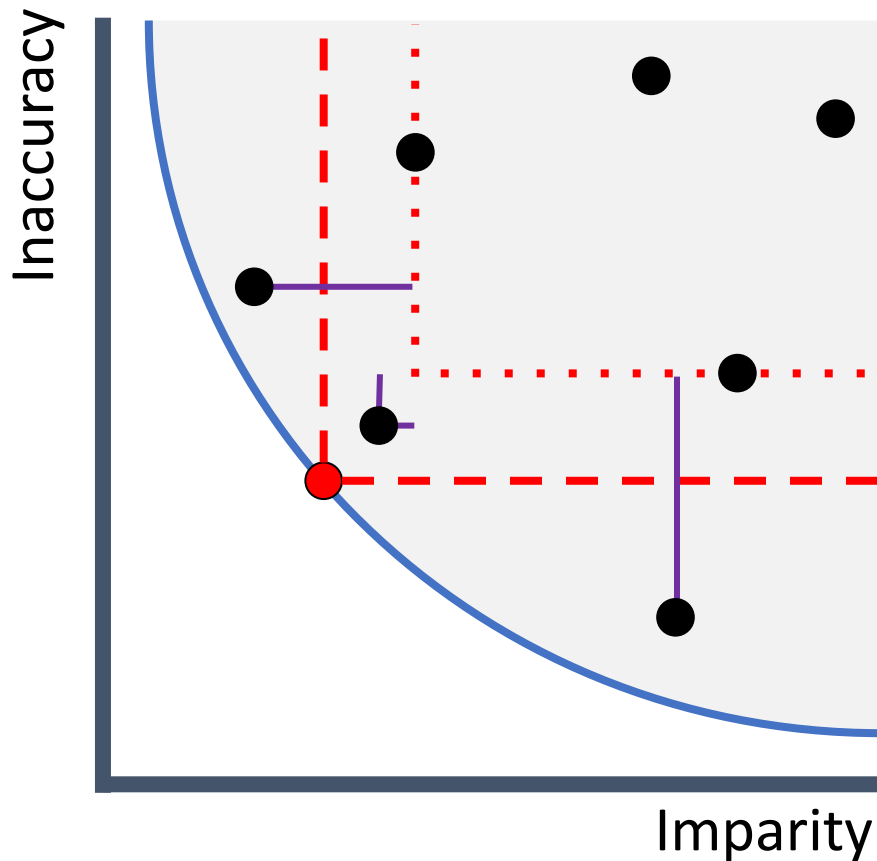A **policy** is **γ-superhuman** if it has smaller **metrics** $f_1, f_2, \ldots$ than γ% of **human demonstrations**

**Subdominance** measures how far a policy is from superhuman by some **margins**

**Minimally subdominant policy** tends to reside close to the **Pareto frontier**
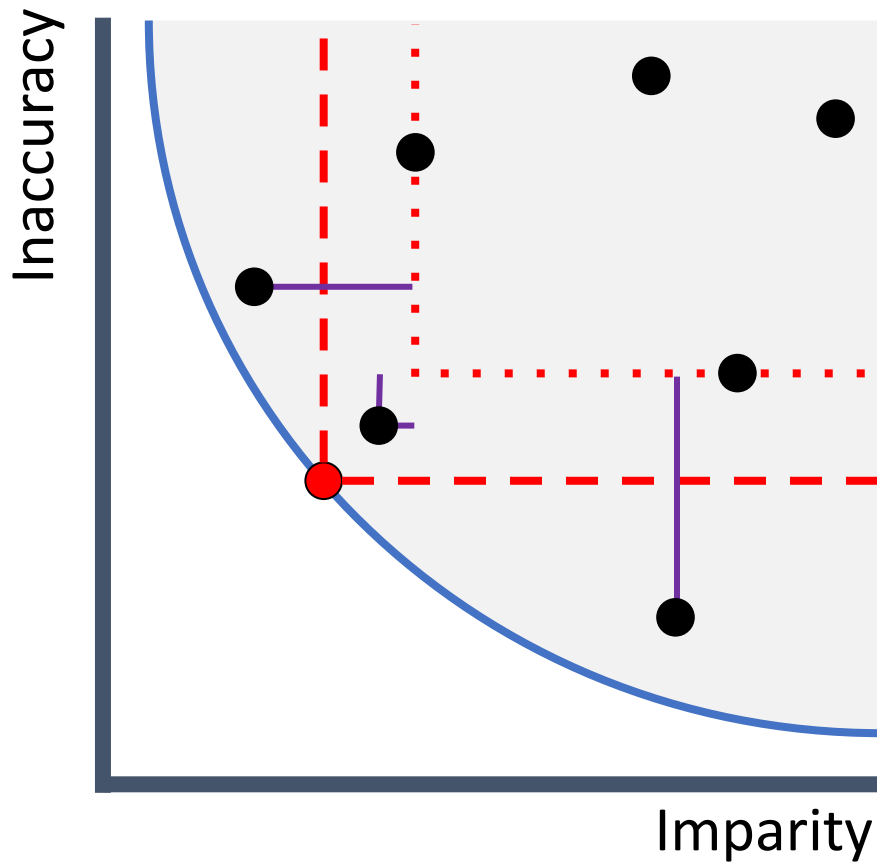
# Superhuman Percentile & Subdominance



A **policy** is **γ-superhuman** if it has smaller **metrics** $f_1, f_2, \ldots$ than γ% of **human demonstrations**

**Subdominance** measures how far a policy is from superhuman by some **margins**

**Minimally subdominant policy** tends to reside close to the **Pareto frontier**

**Subdominance** bounds the **superhuman percentile**

# Subdominance

$$\hat{\mathbf{y}} = \{\hat{y}_j\}_{j=1}^{\mathrm{M}}$$

Model Predictions

$$\tilde{\mathbf{y}} = \{\tilde{y}_j\}_{j=1}^{\mathrm{M}}$$

demonstrations

**The minimally subdominant policy:**

$$\operatorname*{argmin}_{\boldsymbol{\theta}} \min_{\boldsymbol{\alpha} \succeq 0} \mathbb{E}_{\hat{\mathbf{y}}|\mathbf{X} \sim P_{\boldsymbol{\theta}}} \left[ \operatorname{subdom}_{\boldsymbol{\alpha}} \left( \hat{\mathbf{y}}, \tilde{\boldsymbol{\mathcal{Y}}}, \mathbf{y}, \mathbf{a} \right) \right] + \lambda \|\boldsymbol{\alpha}\|_1$$

$\boldsymbol{\theta}$: **Model parameter**

$\alpha$: **Sensitivity to underperform demonstrations**

**If we have metrics inacc, dp, eqodds:**

$\operatorname{subdom}_{\alpha} = \alpha_{\text{inacc}} \operatorname{subdom}_{\text{inacc}} + \alpha_{\text{dp}} \operatorname{subdom}_{\text{dp}} + \alpha_{\text{eqodds}} \operatorname{subdom}_{\text{eqodds}}$
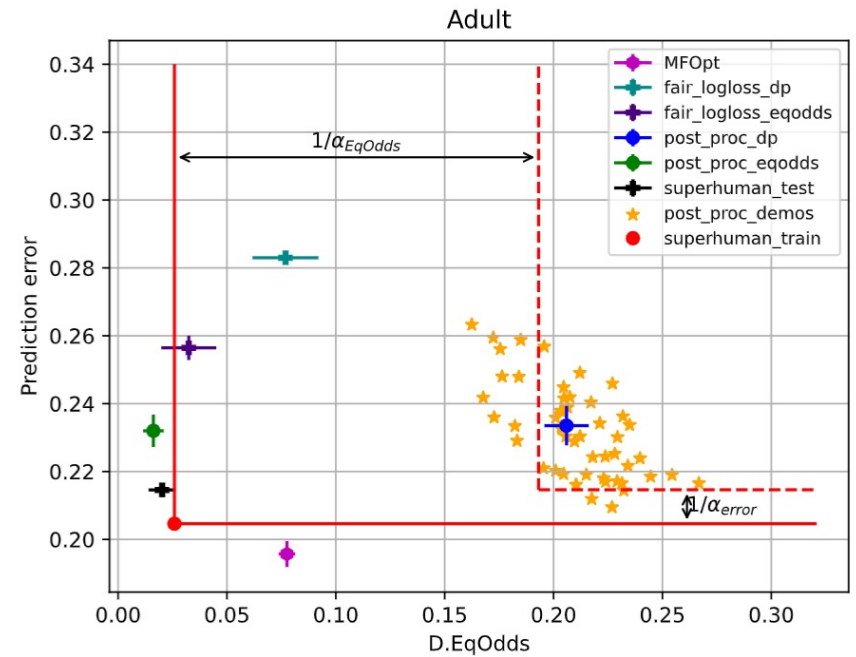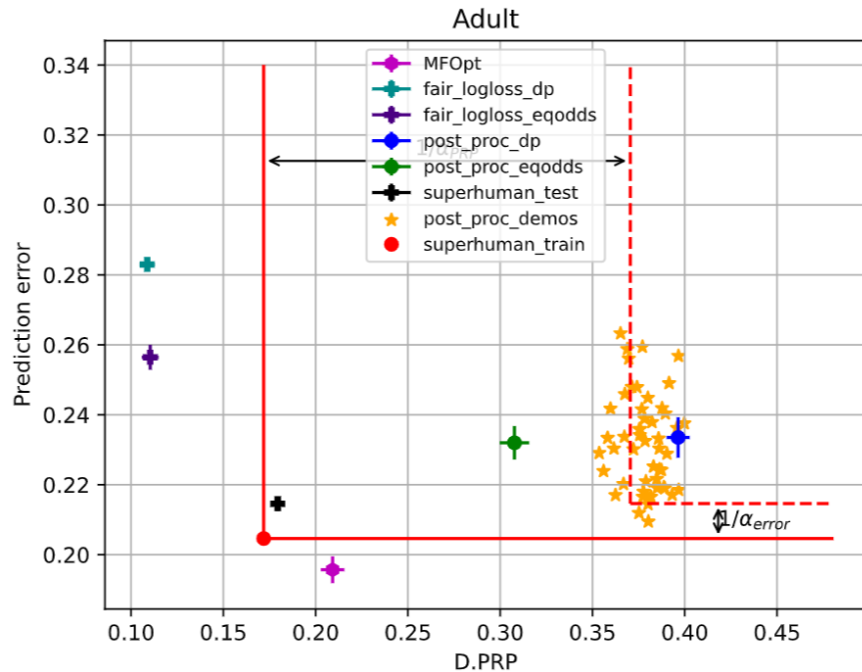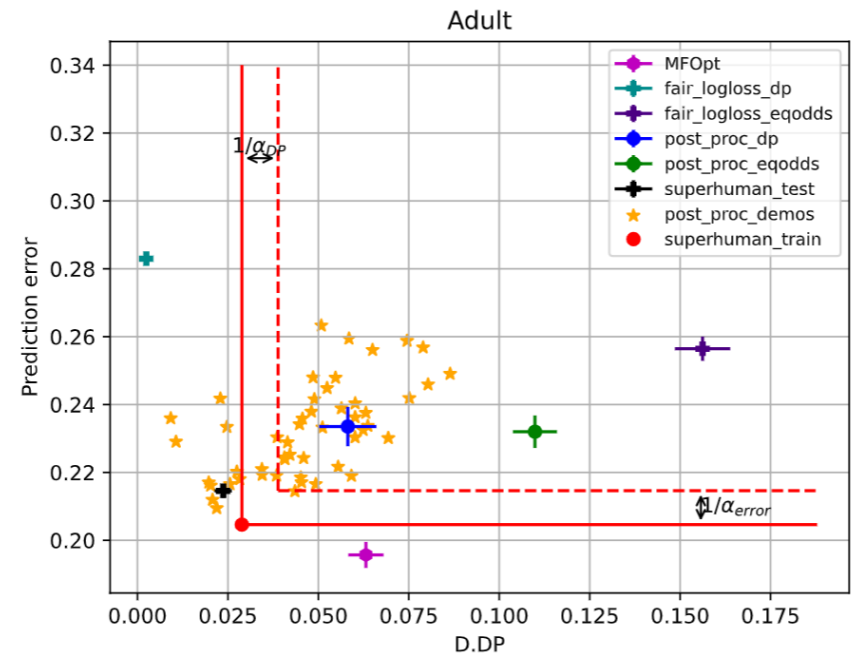
# Experiments

**Metrics:**

(In)Accuracy (Prediction error)

VS

[DP, EqOdds, PRP]

# Thank you!