

# Enhancing Group Fairness in Online Settings Using Oblique Decision Forests

Somnath Basu Roy Chowdhury<sup>1,3</sup>    Nicholas Monath<sup>2</sup>    Ahmad Beirami<sup>1</sup>    Rahul Kidambi<sup>1</sup>  
Kumar Avinava Dubey<sup>1</sup>    Amr Ahmed<sup>1</sup>    Snigdha Chaturvedi<sup>3</sup>

<sup>1</sup>  Google Research    <sup>2</sup>  Google DeepMind    <sup>3</sup>  UNC  
NLP

# Motivation



## Machine Bias

There's software used across the country to predict future criminals. And it's biased against blacks.

*by Julia Angwin, Jeff Larson, Surya Mattu and Lauren Kirchner, ProPublica*

May 23, 2016

- ML systems often produce unfair decisions against certain groups
- We study the challenging problem of achieving fairness in online settings

# Group Fairness

---

Group Fairness techniques focus on enhancing the fairness of ML algorithms by ensuring that different groups receive equal treatment.

# Batch-wise Group Fairness

---

- In batch-wise settings, a learning function  $f$  can be optimized as shown:

$$\min_f L(f(x), y), \text{ subject to } |\mathbb{E}[f(x) | a = 0)] - \mathbb{E}[f(x) | a = 1)]| < \epsilon .$$

$a$  is the sensitive attribute  
(e.g., gender)

# Batch-wise Group Fairness

---

- In batch-wise settings, a learning function  $f$  can be optimized as shown:

$$\min_f L(f(x), y), \text{ subject to } |\mathbb{E}[f(x | a = 0)] - \mathbb{E}[f(x | a = 1)]| < \epsilon.$$

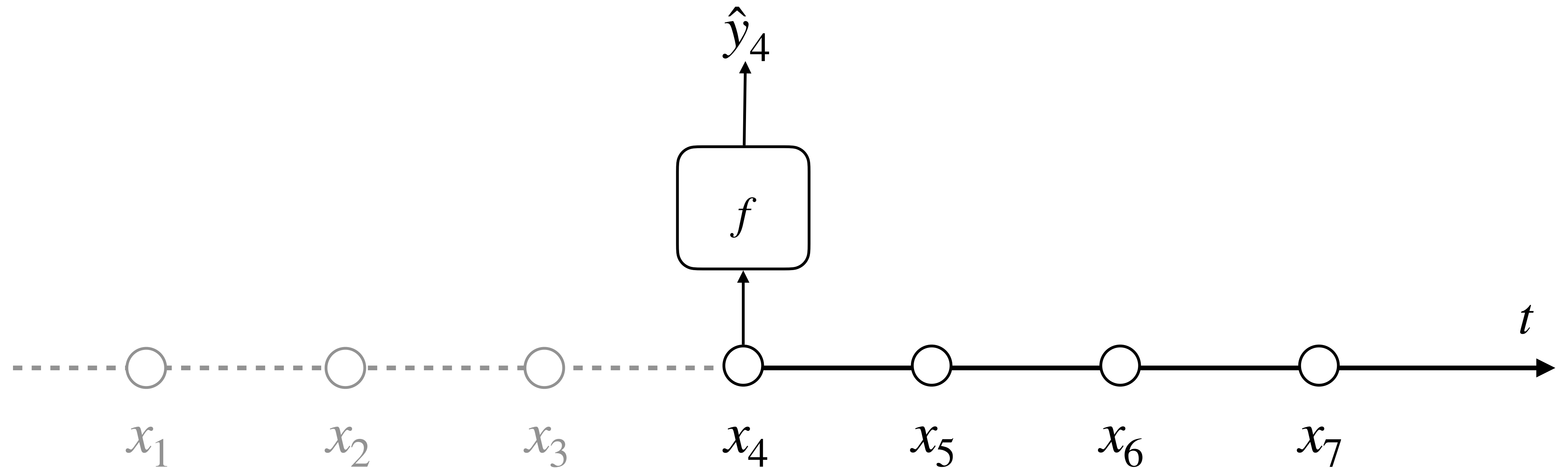
Difference between predictions of two groups



# Online Setting

---

- In online setup, input points  $x_1, x_2, \dots$  arrive one at a time

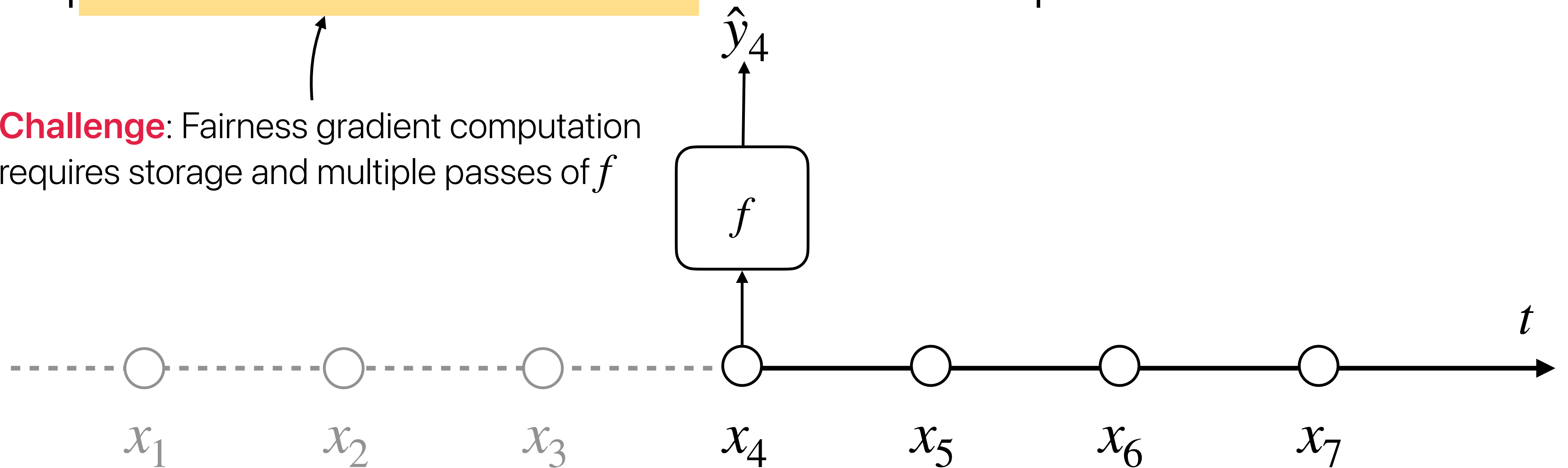


# Online Setting

---

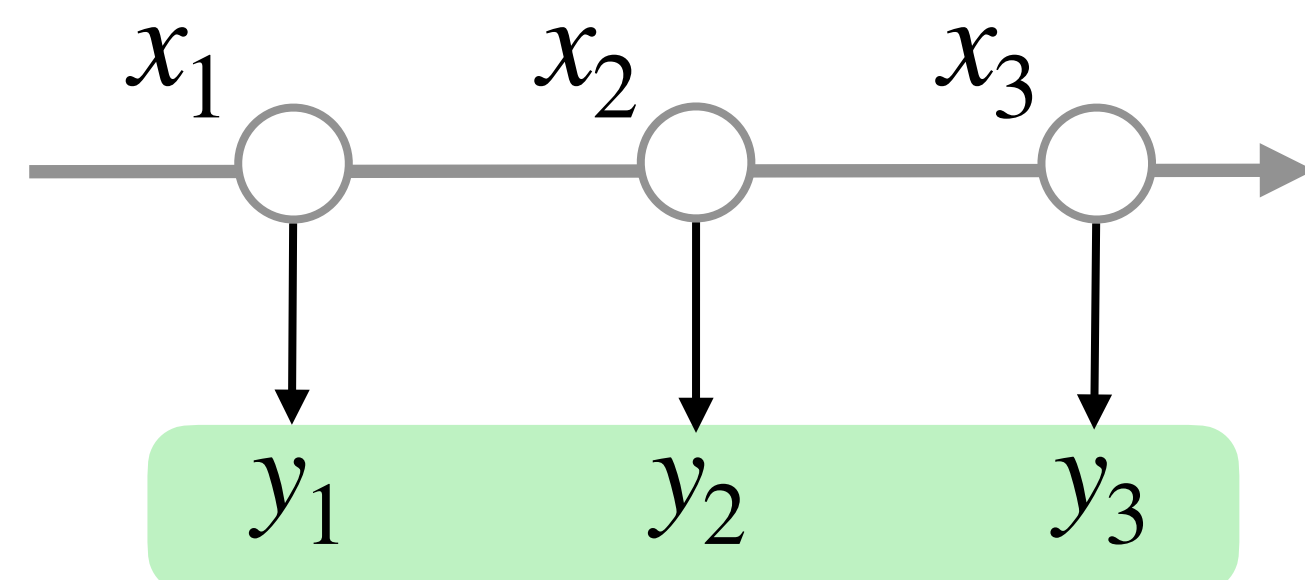
$$\left| \frac{f(x_1 | a = 0) + \dots + f(x_n | a = 0)}{n} - \mathbb{E}[f(x | a = 1)] \right| < \epsilon.$$

**Challenge:** Fairness gradient computation requires storage and multiple passes of  $f$



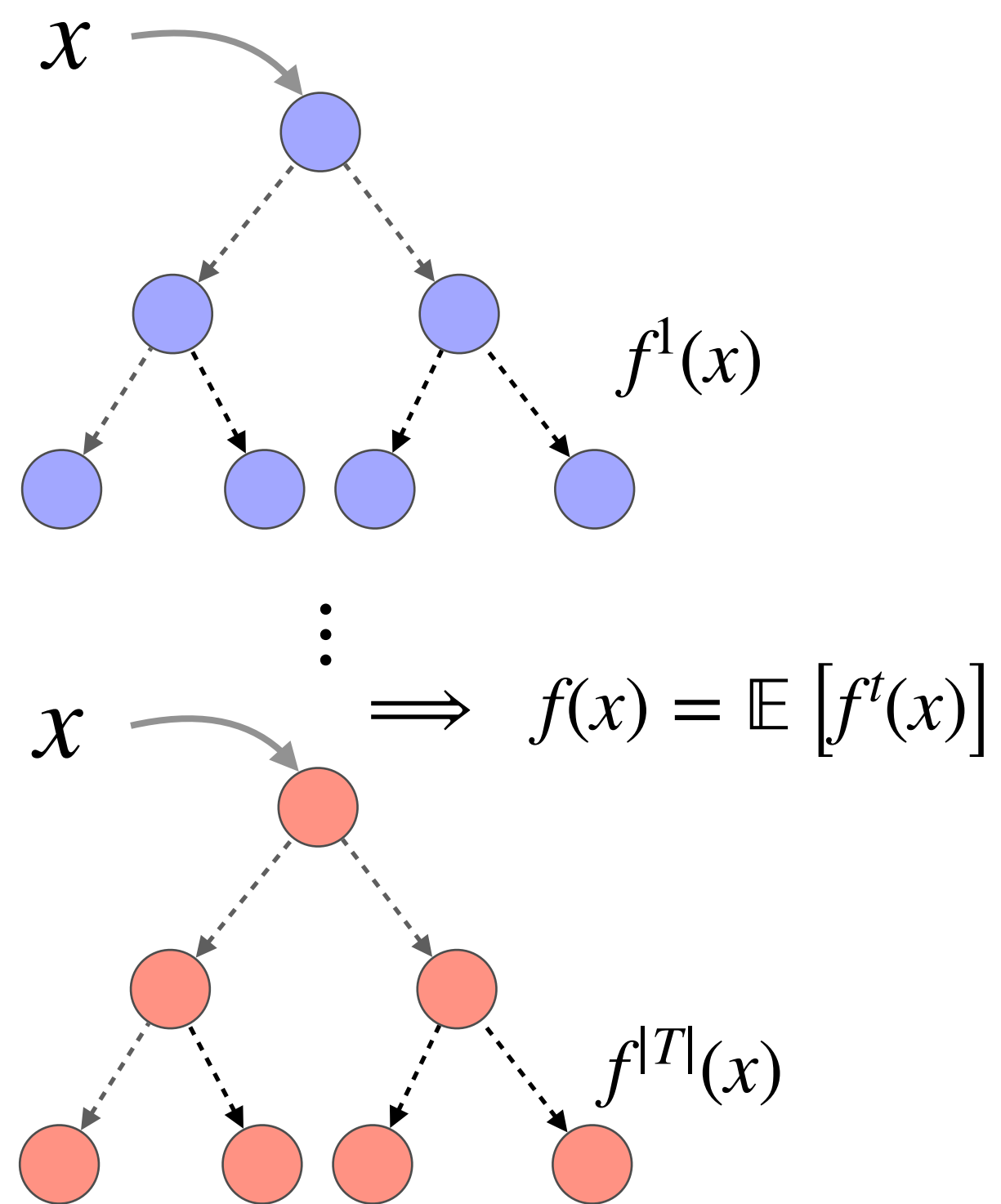
# Overview of Aranyani

## Online Learning For Group Fairness

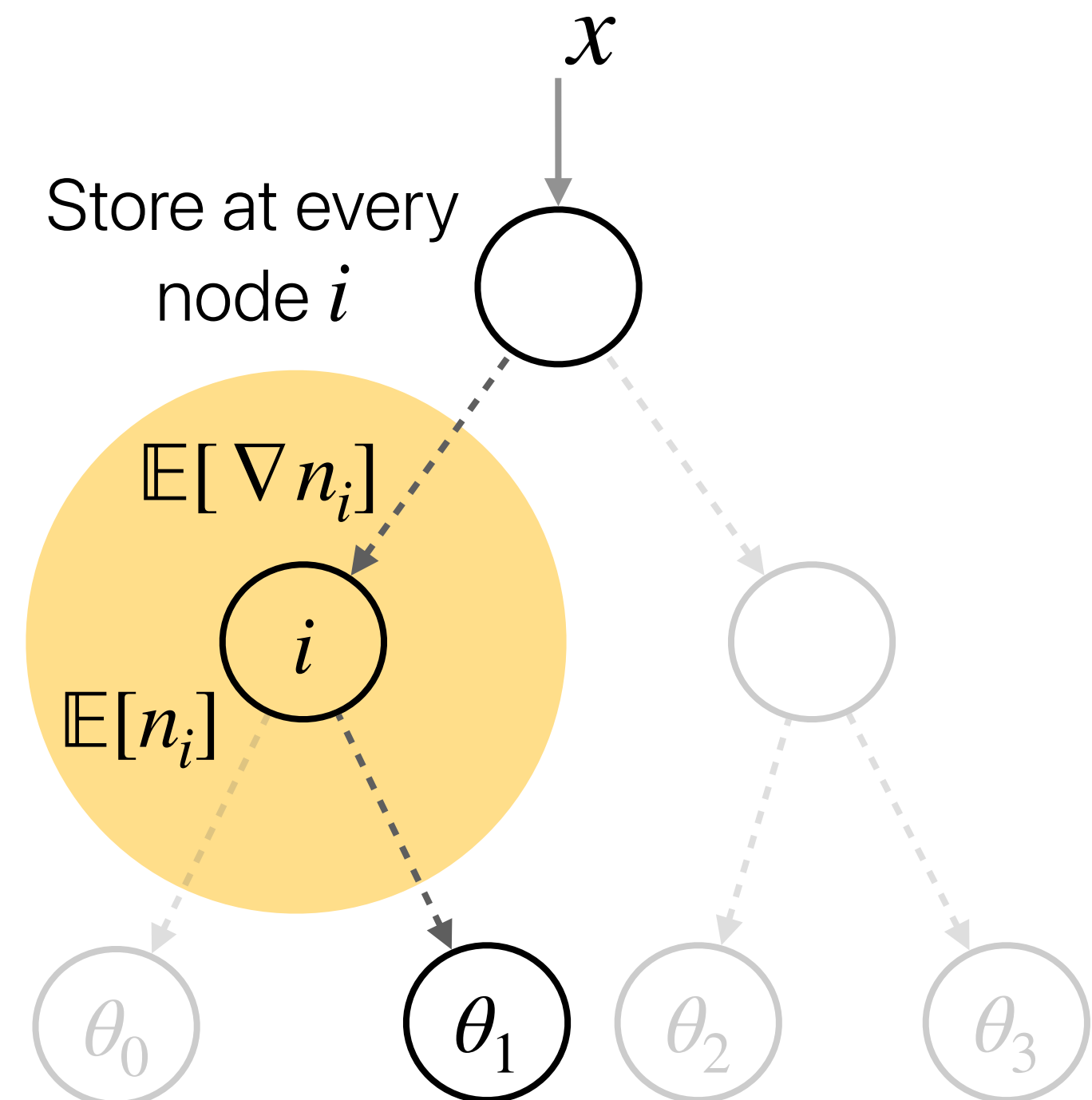


Discrimination  $< \epsilon$

## Prediction Using Oblique Decision Forests



## Gradient Estimation Using Aggregate Statistics



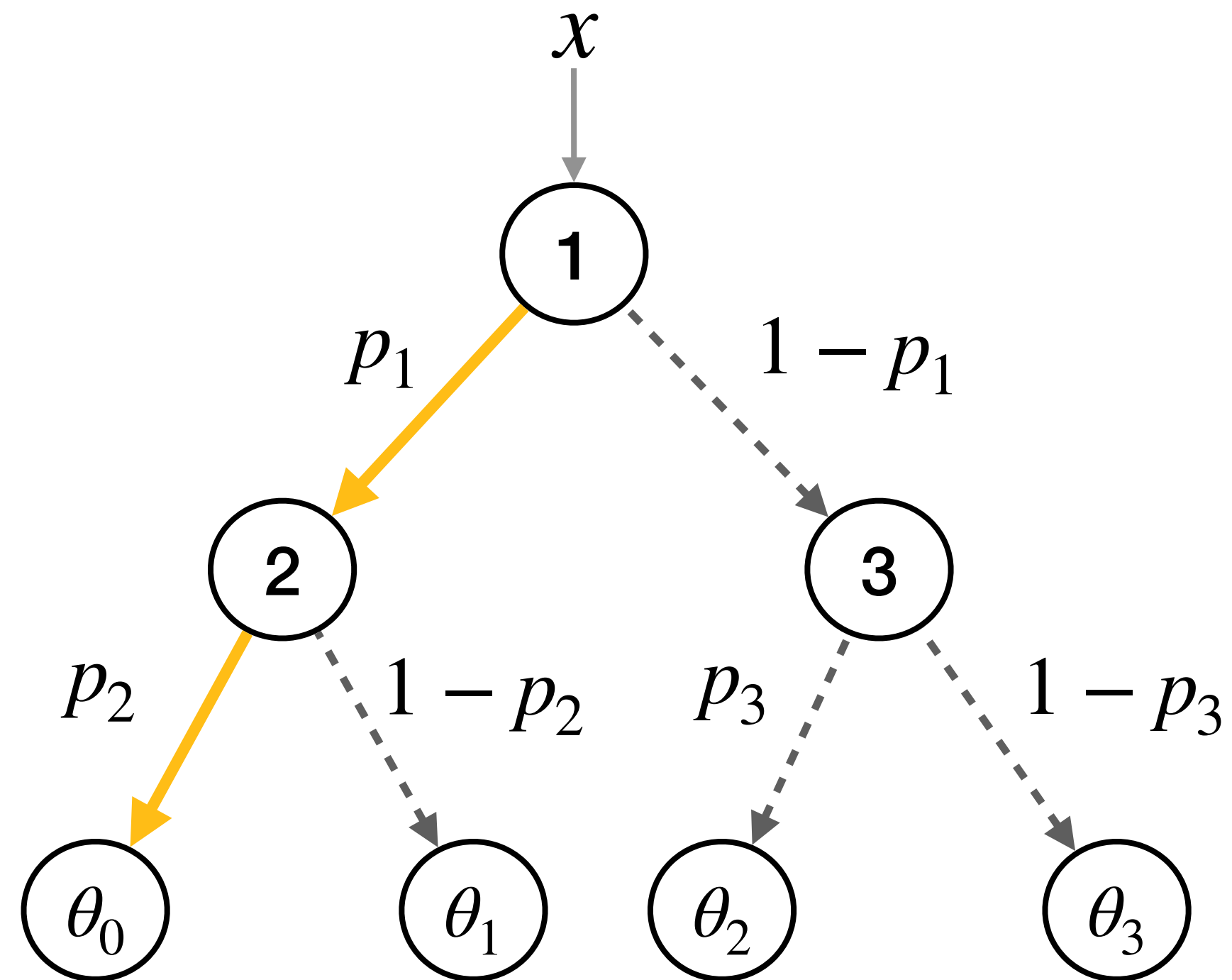


# Aranyani

---

# Aranyani

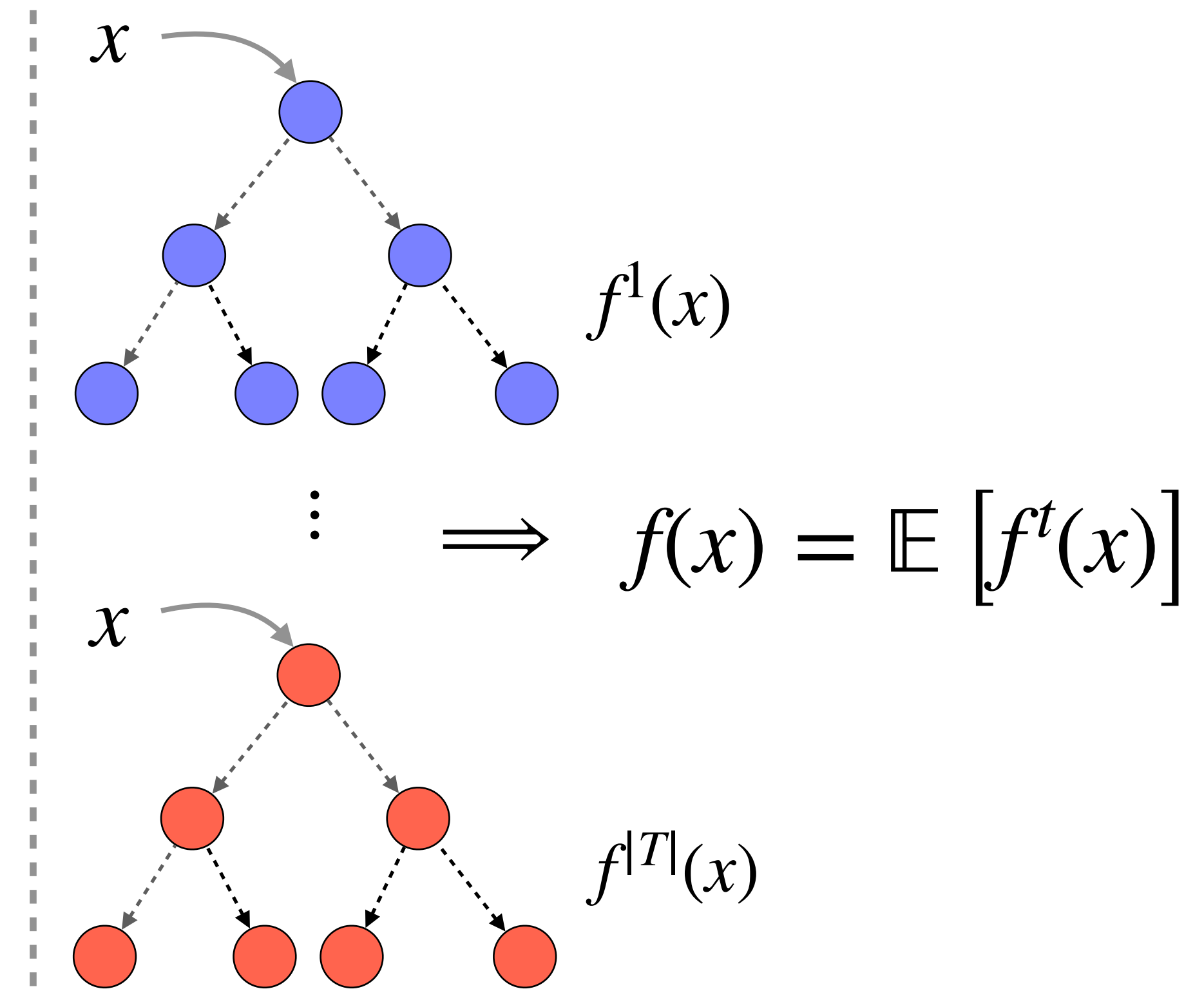
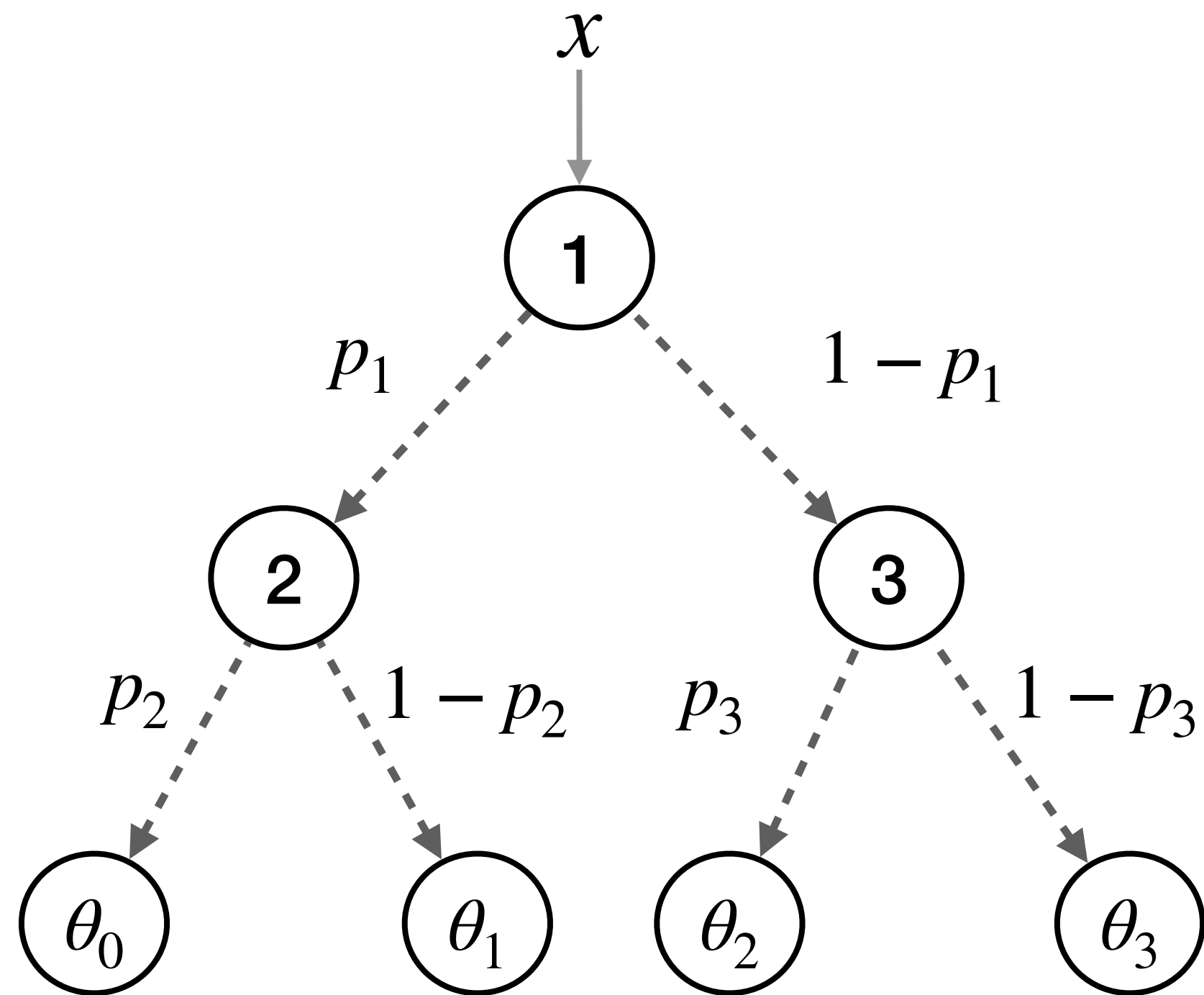
---



$$f(x) = p_1 p_2 \theta_0 + p_1 (1 - p_2) \theta_1 + (1 - p_1) p_3 \theta_2 + (1 - p_1) (1 - p_3) \theta_3$$

# Aranyani

---



# Fairness Gradient Estimation

---

- The fairness gradient estimation process is shown below:

$$G(\Theta) = \nabla_{\Theta} L(f(x), y) + \lambda \sum_{i,j} \nabla_{\Theta} H_{\delta}(F_{ij})$$

Differentiable Huber loss for  
node-level decisions



# Theoretical Results

---

- Estimation error of fairness gradients is bounded:  $\delta B/2$
- The gradient norm  $\Phi_T$  is bounded by

$$\Phi_T \leq (\epsilon + \underline{2^{h-2} \lambda^2 \delta^2 B^2})$$

$h$ : tree height,  $\lambda$ : loss hyperparameter

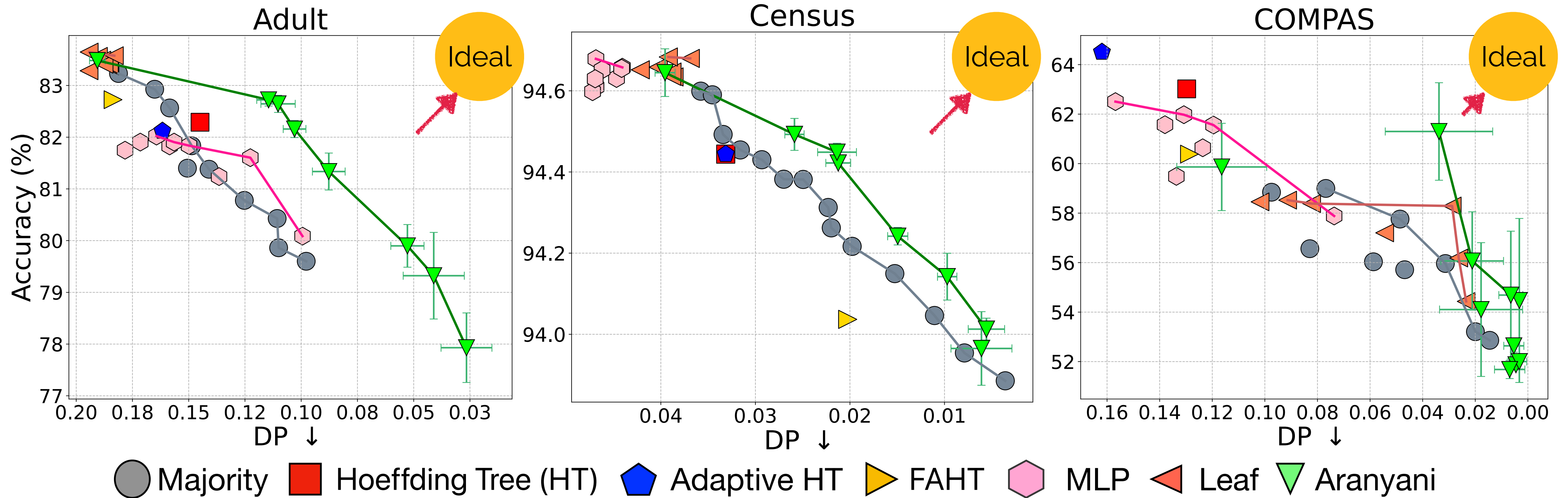
$\delta$ : Huber constant,  $B$ : input bound

# Experiments

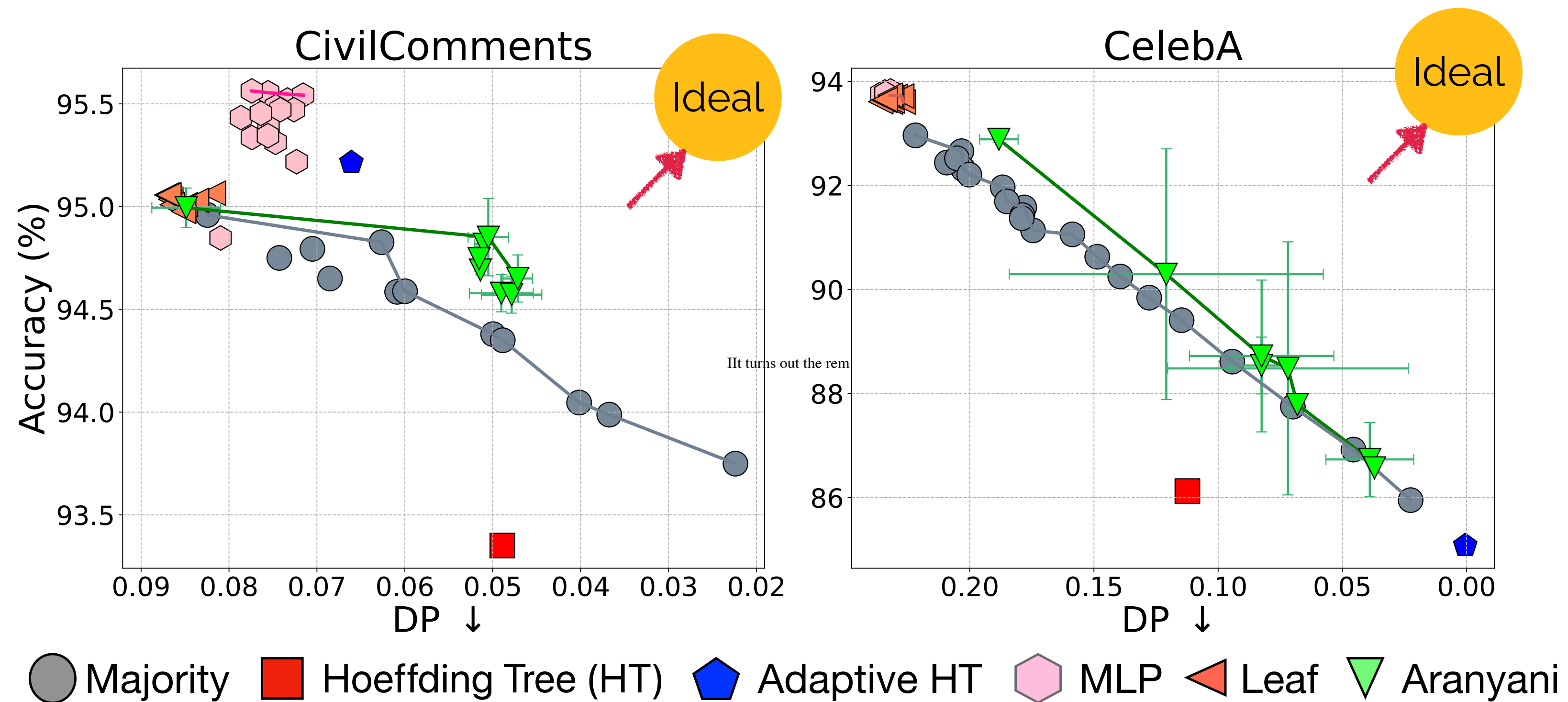
---

- Experiments show effectiveness in Tabular, Vision, and Language datasets
- During online learning, at each step we measure the task performance and fairness
- We report the average performances at the final step,  $T$

# Tabular Datasets



# Vision & Language Datasets



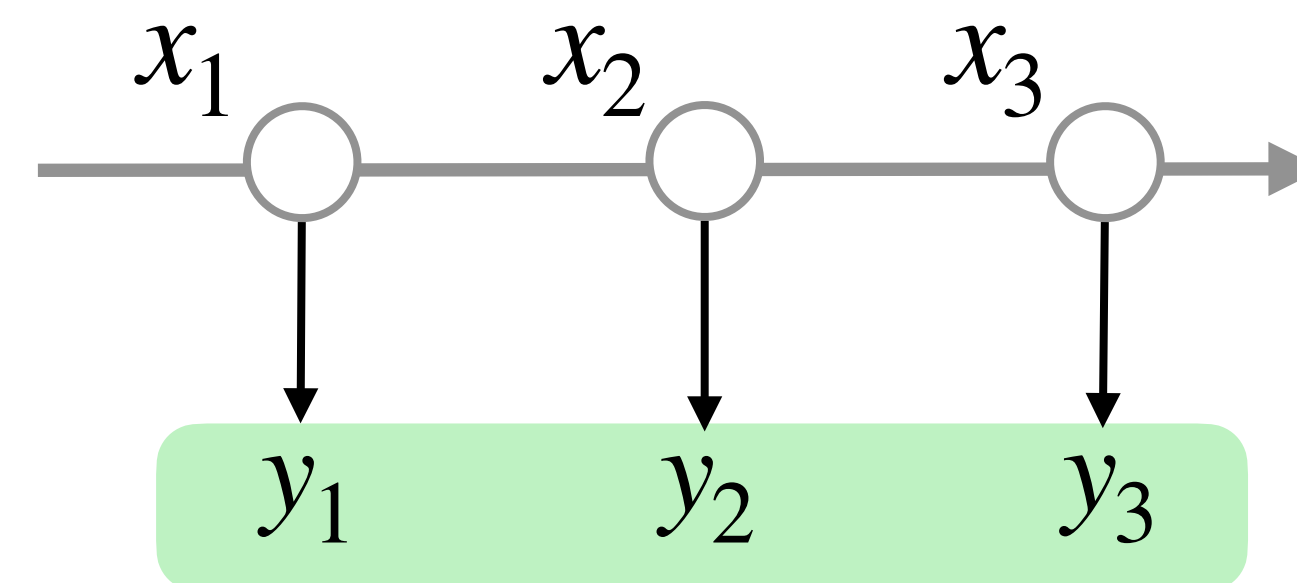


# Summary

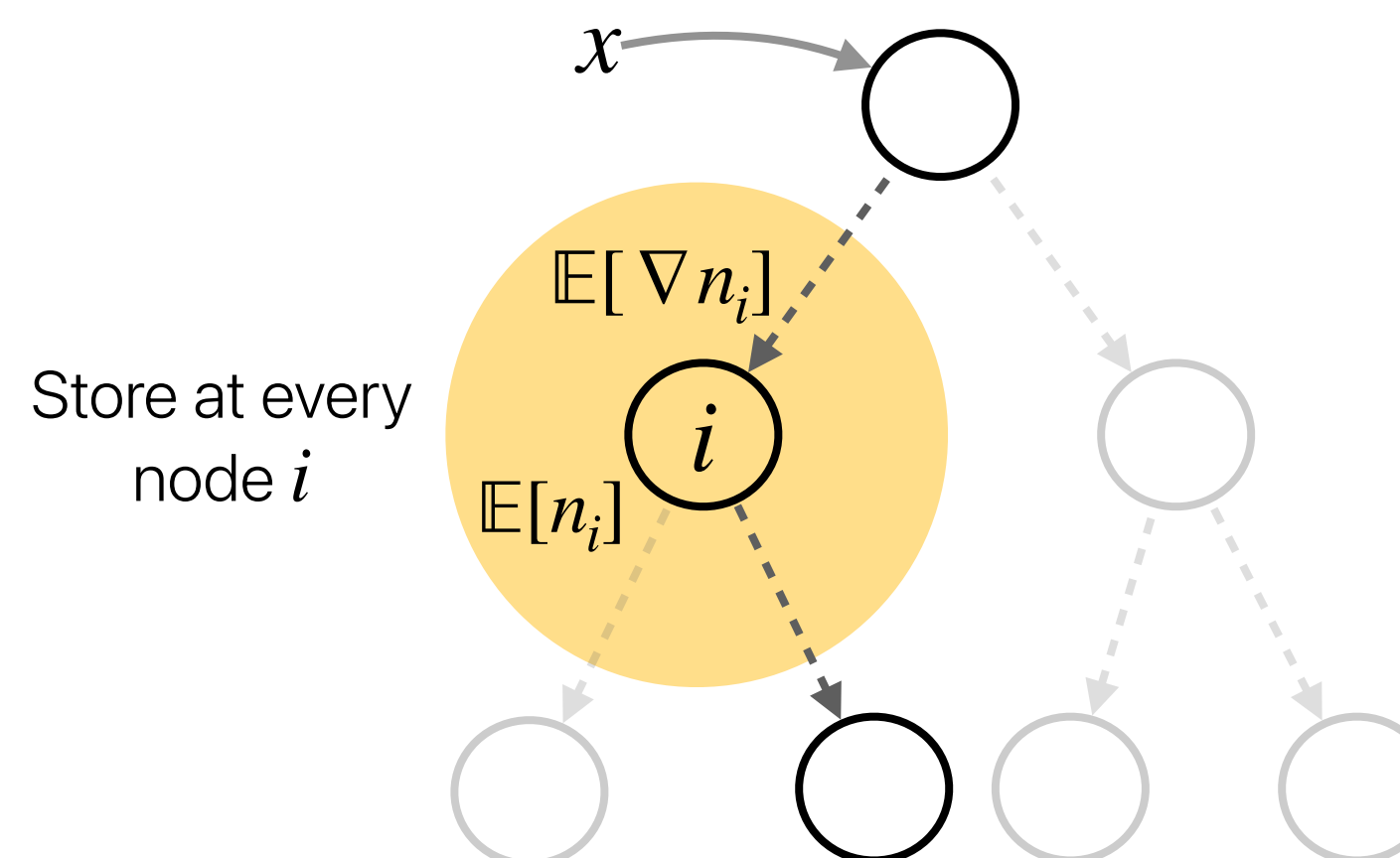
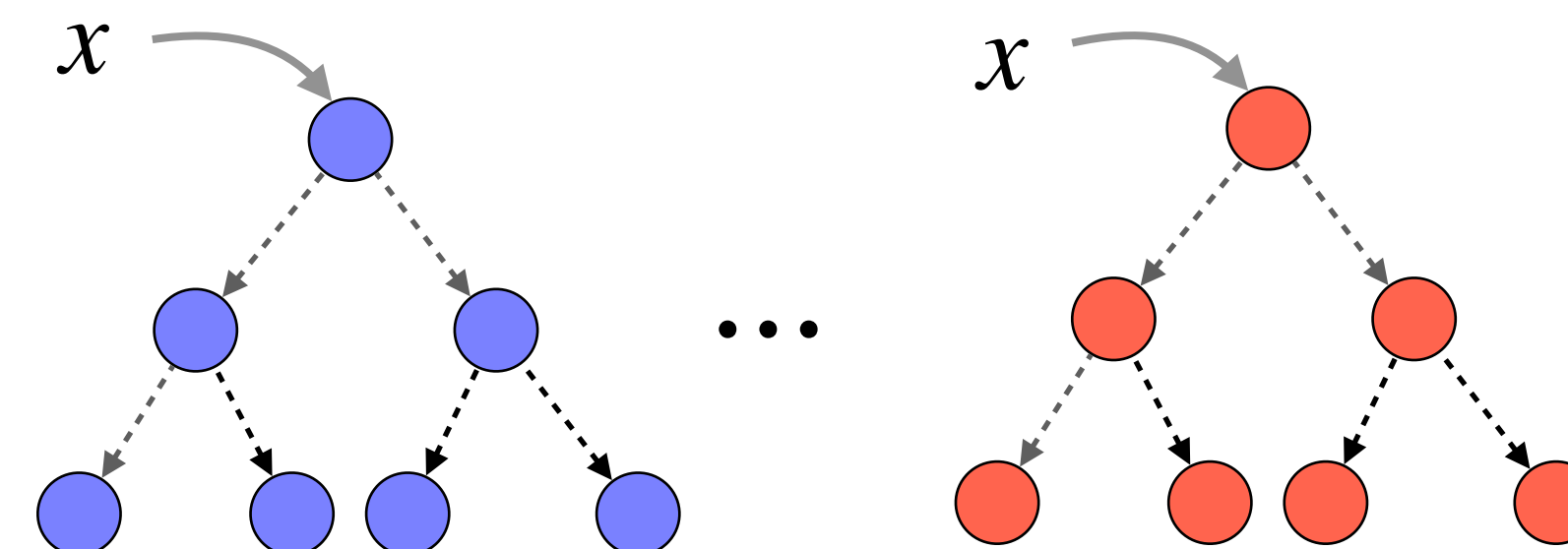
We propose **Aranyani** to achieve group fairness in online environments

**Aranyani** leverages oblique decision forests for efficient online gradient computation

**Aranyani** achieves impressive performance in real-world scenarios



Discrimination  $< \epsilon$



# Thank You!

---

Contact Info:

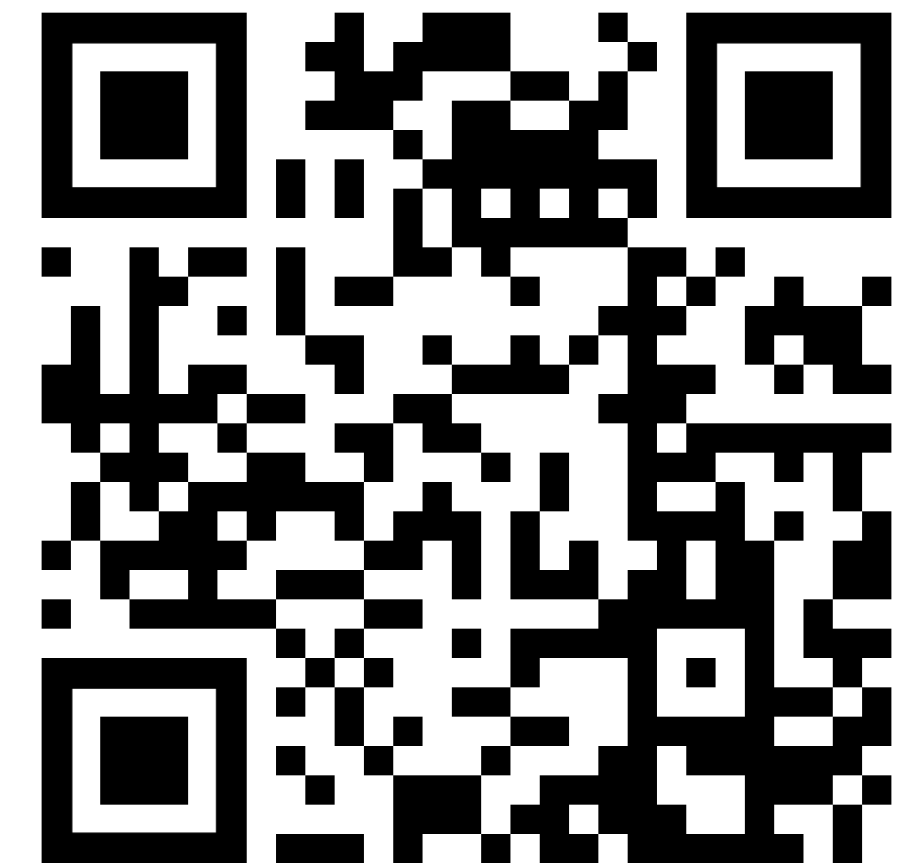
Somnath Basu Roy Chowdhury

UNC Chapel Hill

[somnath@cs.unc.edu](mailto:somnath@cs.unc.edu)



Paper



Code