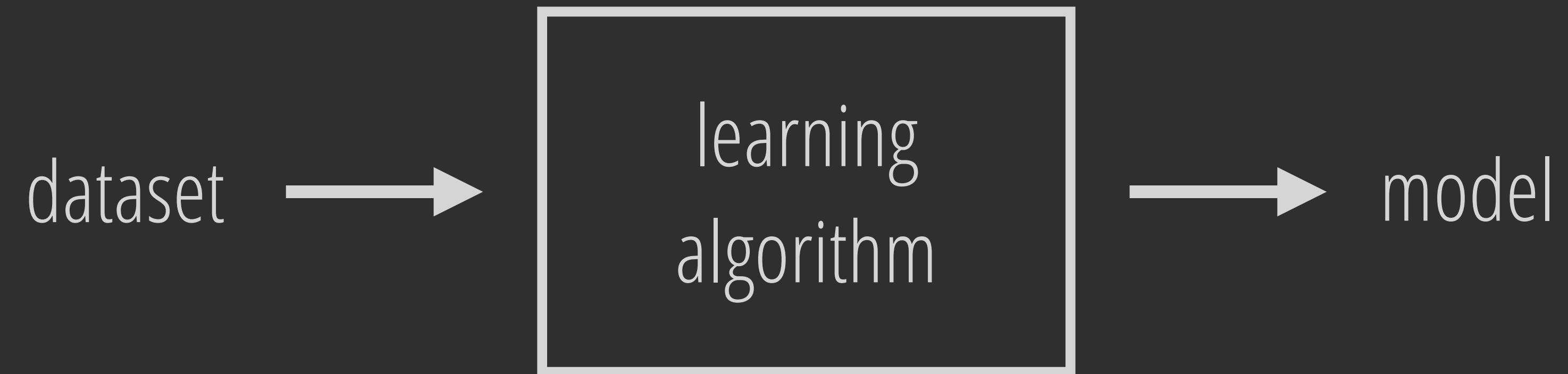


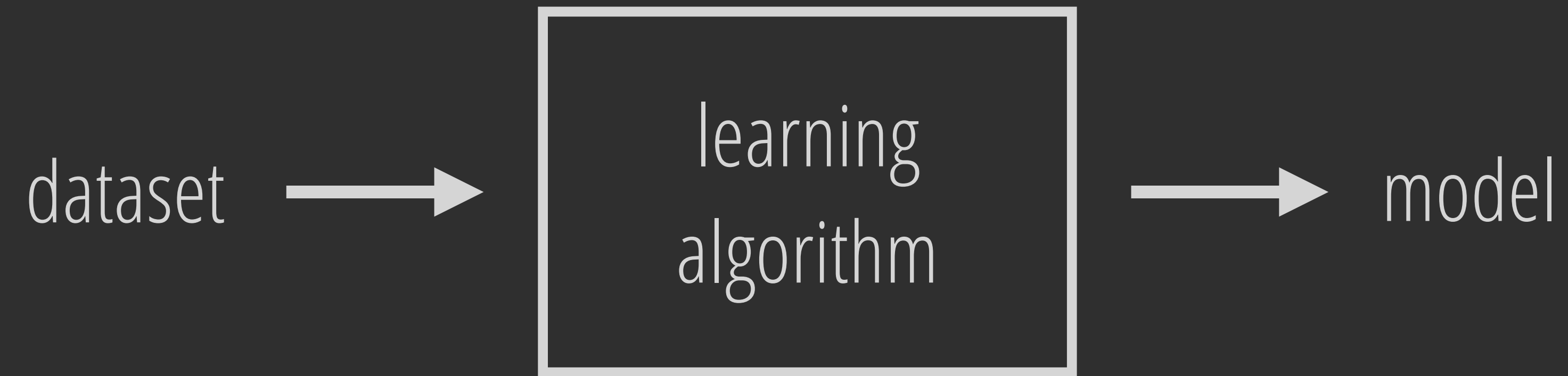
# Certifying Robustness to Programmable Data Bias in Decision Tree Learning

Anna P. Meyer

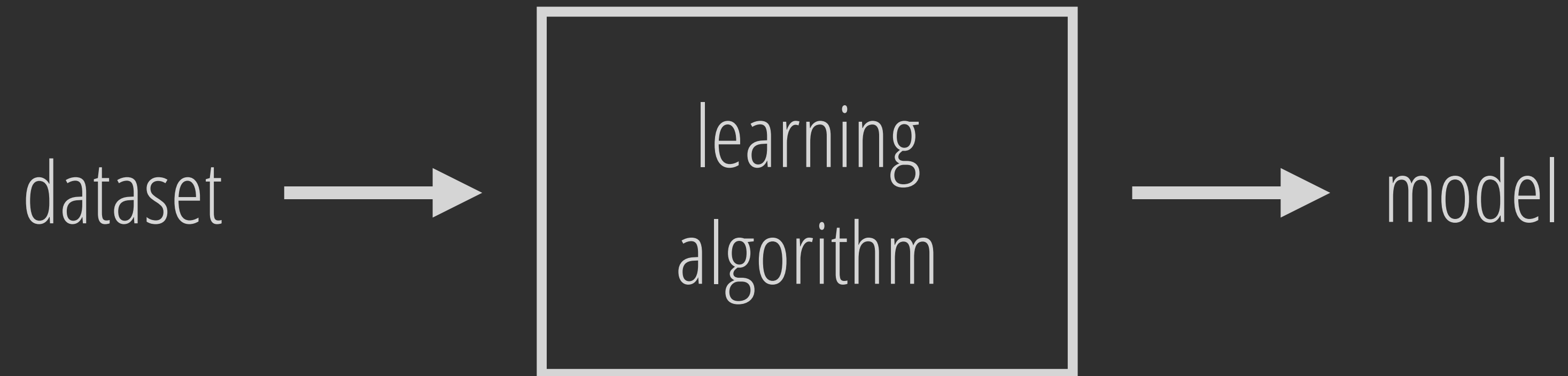
Joint work with Aws Albarghouthi and Loris D'Antoni







Is the dataset unbiased?  
complete?  
representative?



Is the dataset biased?

complete?

representative?

Probably not.

What is the impact on the  
model's predictions?

Goal: certify robustness to training-data bias

# Types of dataset bias

- Label-flipping
- Missing data
- Fake data (data-poisoning)

# Types of dataset bias

- Label-flipping
- Missing data
- Fake data (data-poisoning)

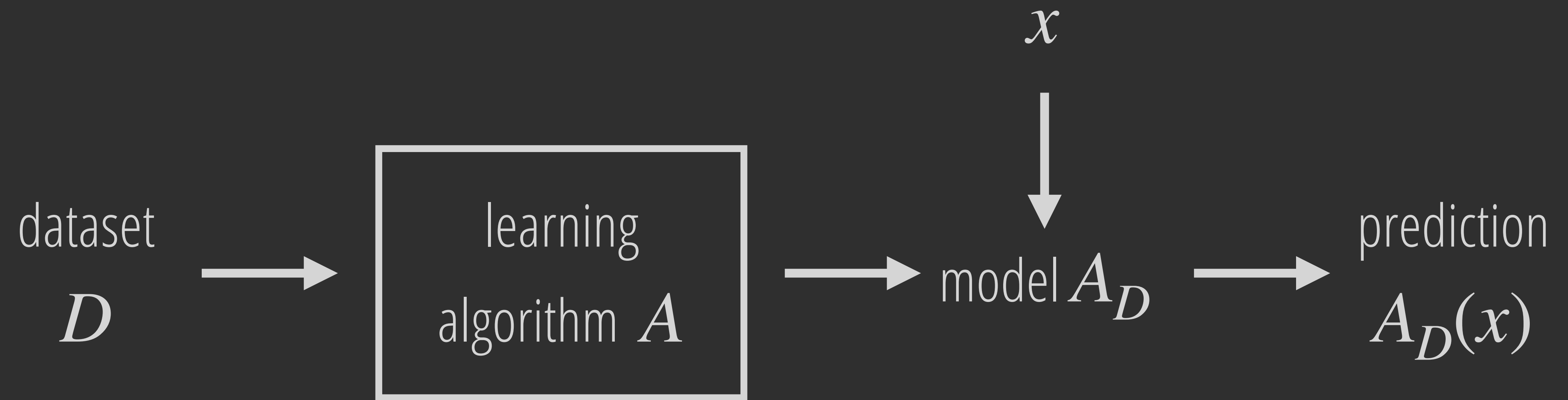
Assume fixed amount and type of data bias

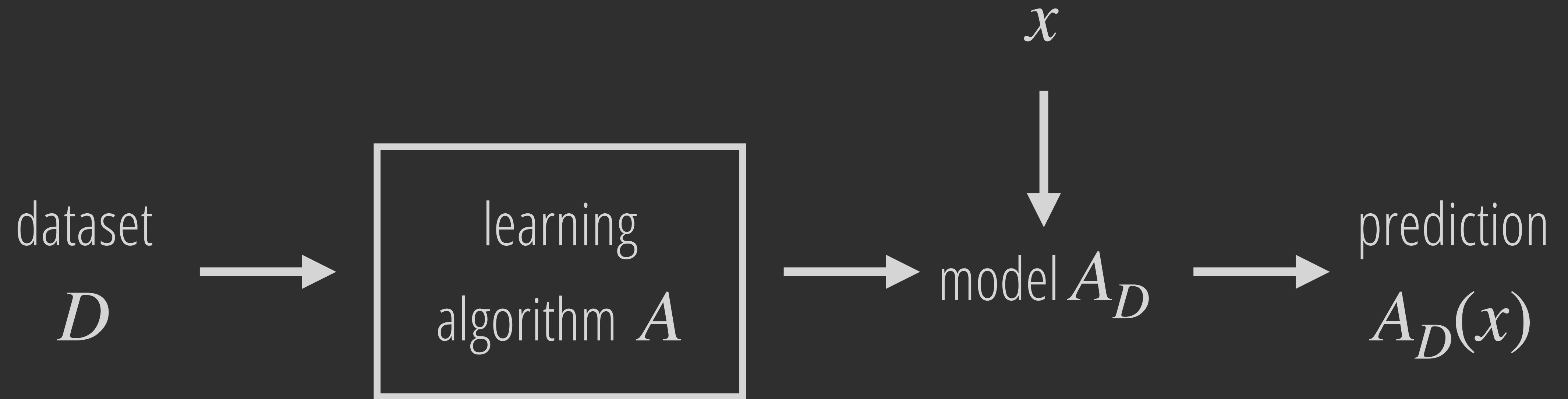
# Types of dataset bias

- Label-flipping
- Missing data
- Fake data (data-poisoning)

Each type can be general or targeted







bias robustness of  $x$

for all  $D'$  that disagree with  $D$  on  $\leq n$  labels

show that  $A_{D'}(x) = A_D(x)$

bias robustness of  $x$

for all  $D'$  that disagree with  $D$  on  $\leq n$  labels

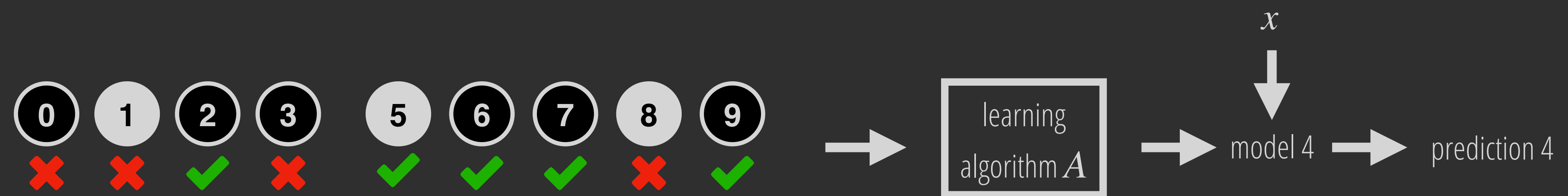
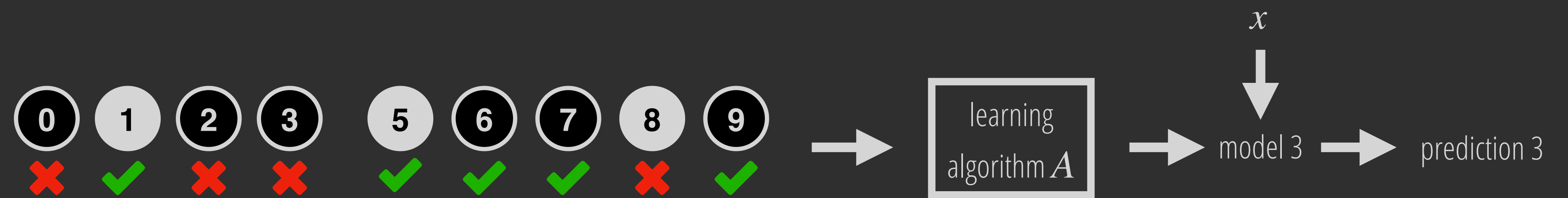
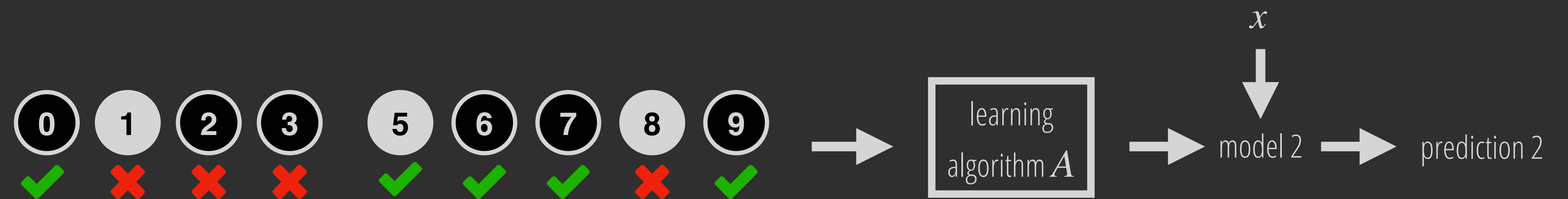
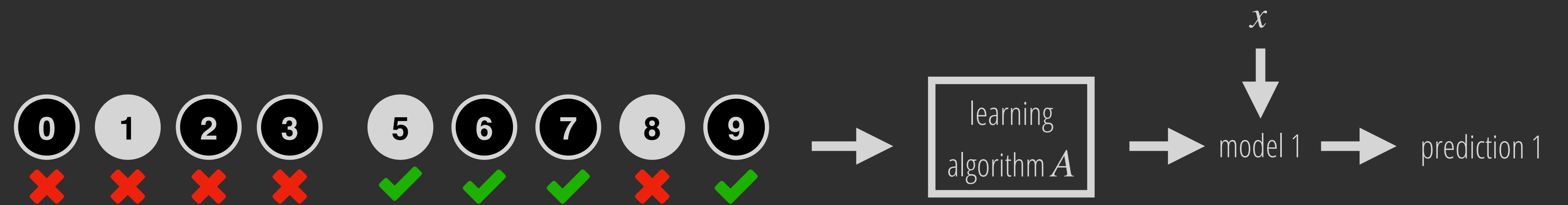
show that  $A_{D'}(x) = A_D(x)$

Dataset  $D$





etc.



etc.

$$|D| = 1000$$

$$n = 10$$

$\sim 10^{23}$  datasets!

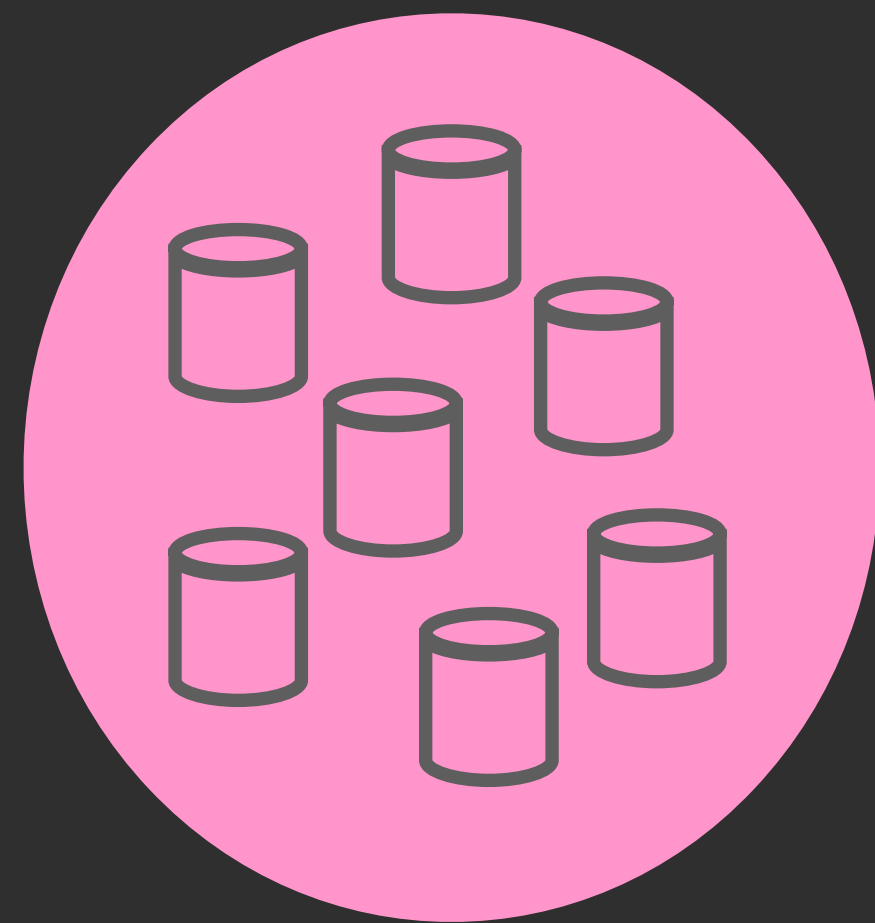
bias robustness of  $x$

for all  $D'$  that disagree with  $D$  on  $\leq n$  labels

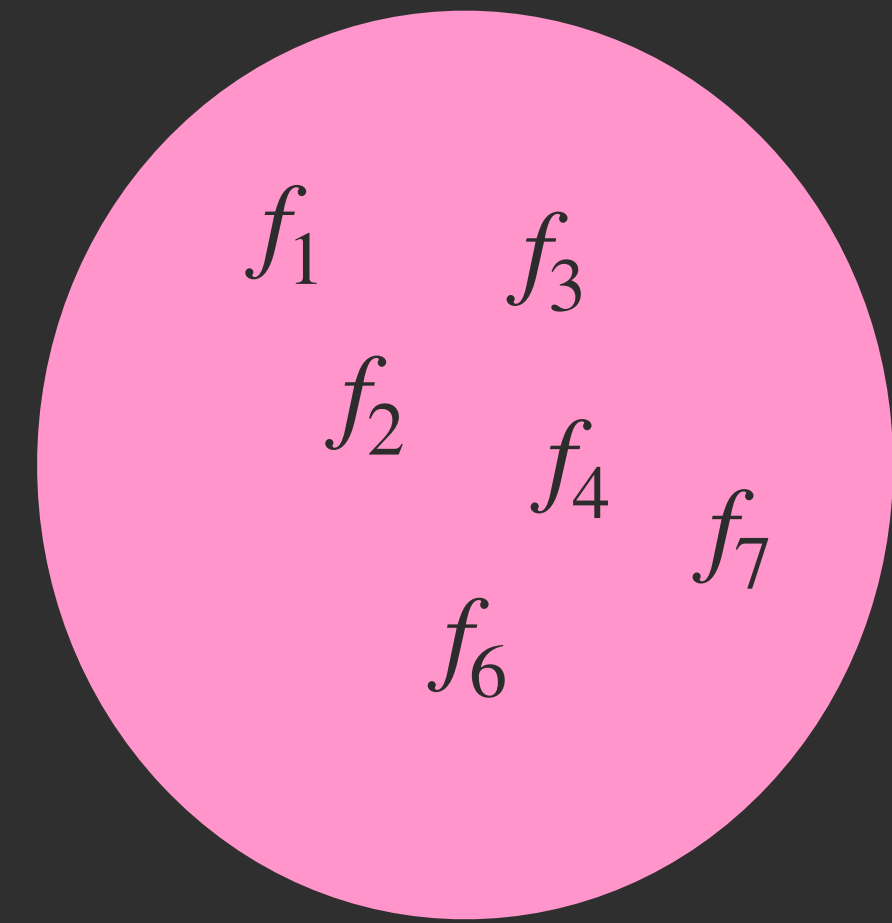
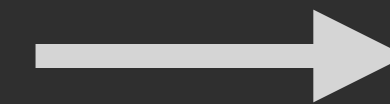
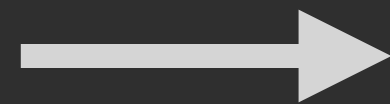
show that  $A_{D'}(x) = A_D(x)$

Key challenge

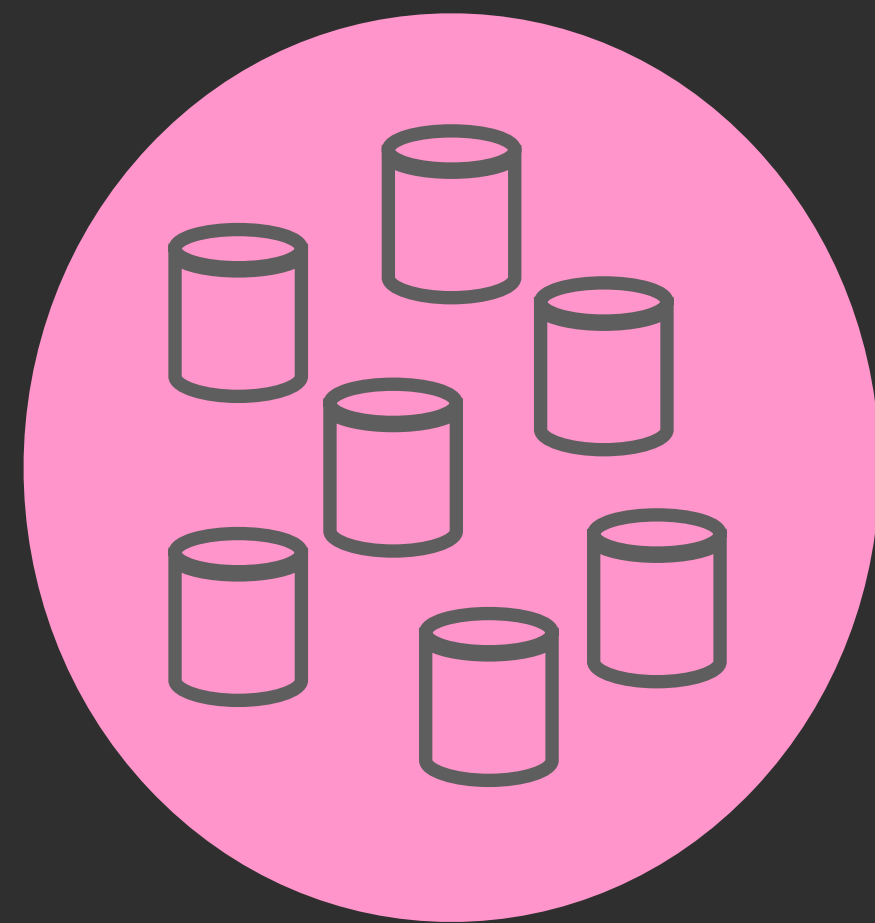
Combinatorial explosion in the number of datasets



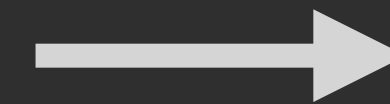
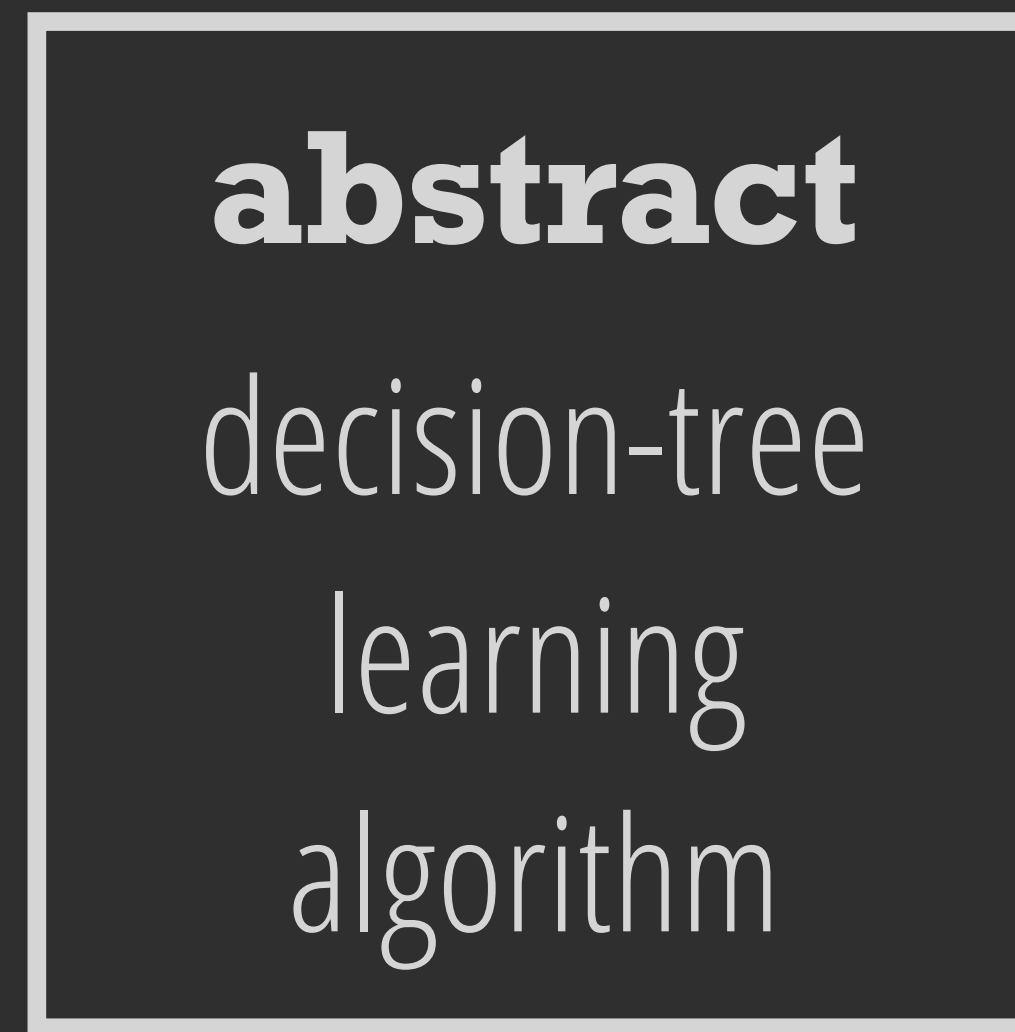
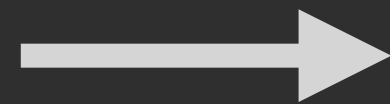
large set of  
training datasets



large set of  
trained models



large set of  
training datasets



large set of  
decision trees



(Very) simplified decision tree algorithm

1. Choose the predicate that minimizes entropy on the data
2. Split the data according to that predicate
3. Repeat on child nodes until entropy is 0, or maximum depth is reached

(Very) simplified decision tree algorithm

1. Choose the predicate that minimizes entropy on the data
2. Split the data according to that predicate
3. Repeat on child nodes until entropy is 0, or maximum depth is reached

# Dataset $D$



Dataset  $D$



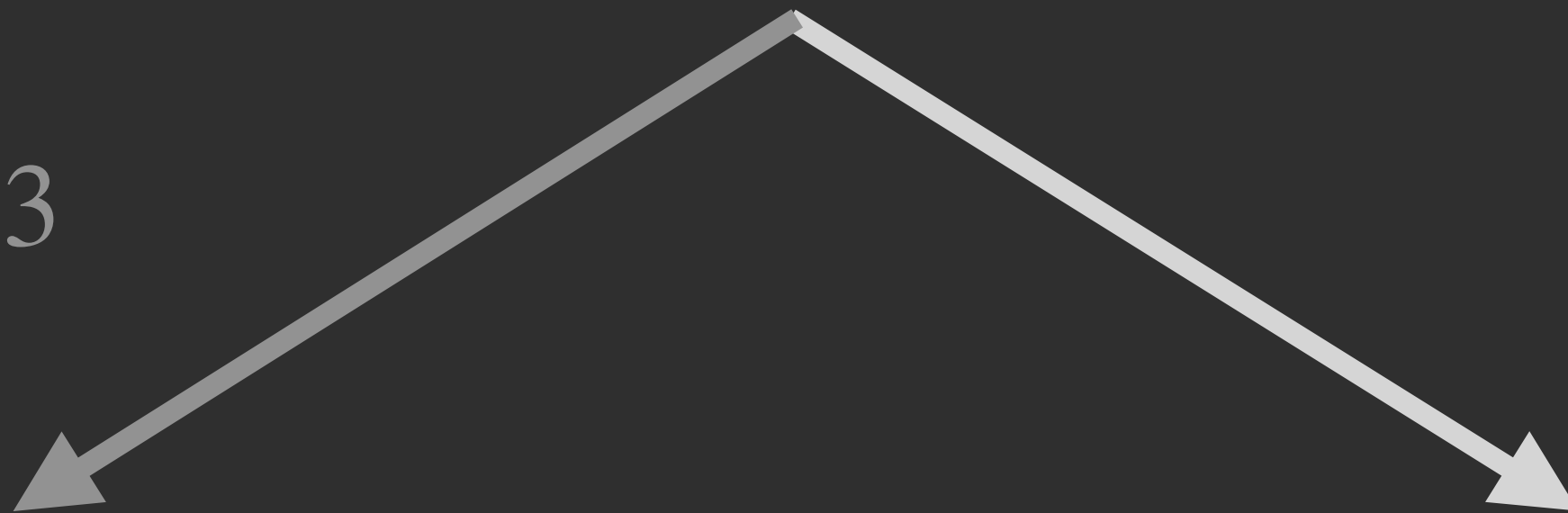
$\phi := \text{value} \leq 3$



Dataset  $D$



$\phi := \text{value} \leq 3$



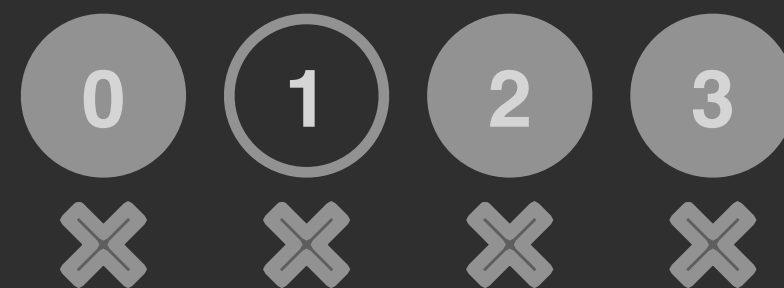
Number ✓ = 4

Number x = 1

Dataset  $D$



$\phi := \text{value} \leq 3$

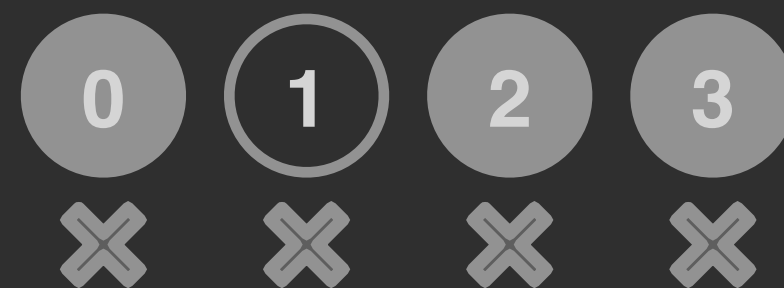


$$\begin{aligned} \text{Gini Impurity} &= \checkmark \cdot (1 - \checkmark) + \text{X} \cdot (1 - \text{X}) \\ &= \frac{4}{5} \left(1 - \frac{4}{5}\right) + \frac{1}{5} \left(1 - \frac{1}{5}\right) = 0.32 \end{aligned}$$

# Abstraction of Dataset $D$



$\phi := \text{value} \leq 3$



Number ✓ = 4  
Number ✗ = 1

Number ✓ = [3, 5]  
Number ✗ = [0, 2]

## Abstraction of Dataset $D$



$\phi := \text{value} \leq 3$



$$\text{Gini Impurity} = \checkmark \cdot (1 - \checkmark) + \text{✗} \cdot (1 - \text{✗})$$

$$= \frac{[3,5]}{5} \left(1 - \frac{[3,5]}{5}\right) + \frac{[0,2]}{5} \left(1 - \frac{[0,2]}{5}\right) = [0, 0.8]$$





$$\begin{aligned}\text{Gini Impurity} &= \checkmark \cdot (1 - \checkmark) + \times \cdot (1 - \times) \\ &= \frac{[3,5]}{5} \left(1 - \frac{[3,5]}{5}\right) + \frac{[0,2]}{5} \left(1 - \frac{[0,2]}{5}\right) = [0, 0.8]\end{aligned}$$

Aside: We can be more precise!

E.g., we can't simultaneously have 5  $\checkmark$  and 2  $\times$ .

—> details omitted from this presentation, but in our experiments we use the precise version

# Experimental results

# Certification rate

Given n% bias, what percentage of test data points are certifiably robust?

Bias type	Dataset	Bias amount as a percentage of training set					
		0.05	0.1	0.2	0.4	0.7	1.0
FLIP (label-flipping)	Drug Consumption	94.5	94.5	94.5	94.5	92.1	85.1
	COMPAS	100.0	89.0	81.5	71.5	47.8	39.7
	Adult Income	98.4	96.6	85.2	64.4	42.2	14.8
	COMPAS targeted	100.0	89.0	89.0	81.9	76.2	53.0
	AI targeted	98.8	98.6	96.6	81.6	69.0	52.8

# Certification rate

Given n% bias, what percentage of test data points are certifiably robust?

Bias type	Dataset	Bias amount as a percentage of training set					
		0.05	0.1	0.2	0.4	0.7	1.0
FLIP (label-flipping)	Drug Consumption	94.5	94.5	94.5	94.5	92.1	85.1
	COMPAS	100.0	89.0	81.5	71.5	47.8	39.7
	Adult Income	98.4	96.6	85.2	64.4	42.2	14.8
	COMPAS targeted	100.0	89.0	89.0	81.9	76.2	53.0
	AI targeted	98.8	98.6	96.6	81.6	69.0	52.8

Bias-set size color scheme	< 10 <sup>10</sup>	< 10 <sup>50</sup>	< 10 <sup>100</sup>	< 10 <sup>500</sup>	> 10 <sup>500</sup>	infinite
----------------------------	--------------------	--------------------	---------------------	---------------------	---------------------	----------

# Future work

- Extensions to other ML algorithms
- Counter-examples to robustness