

Testing and Reconstruction via Decision Trees

Li-Yang Tan

Joint work with:



Guy Blanc

Stanford



Jane Lange

MIT



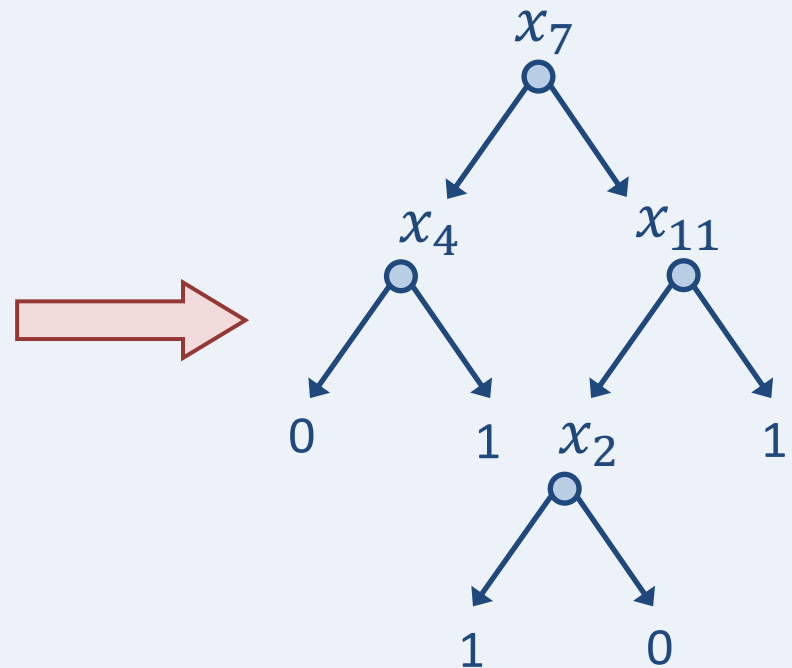
(Slides and preprint available on my webpage)

Decision tree learning

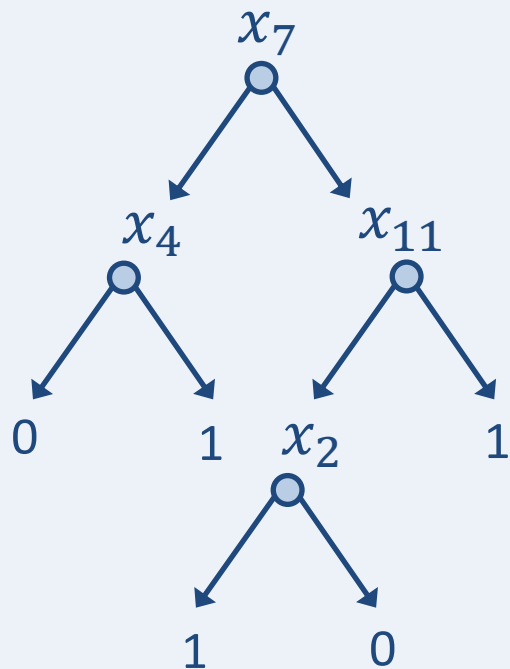
Labeled data

| Example $x \in \{0,1\}^n$ | Label $y \in \{0,1\}$ |
|---------------------------|-----------------------|
| 1 0 0 0 1 0 1 0 0 1 0 1 | 1 |
| 0 1 0 0 1 1 0 0 1 0 1 1 | 0 |
| 1 1 1 0 0 0 1 0 0 1 0 1 | 1 |
| 0 0 1 0 1 0 1 1 1 0 1 0 | 1 |
| 0 0 1 0 0 1 1 0 0 1 0 1 | 0 |
| 1 1 1 1 1 0 0 1 0 0 0 0 | 1 |
| 0 0 1 1 0 0 1 0 0 1 0 0 | 1 |
| 1 0 0 1 0 0 1 1 1 0 1 0 | 0 |
| 1 0 0 1 0 1 0 1 0 1 1 0 | 1 |

Decision tree representation



Decision trees: simple and effective

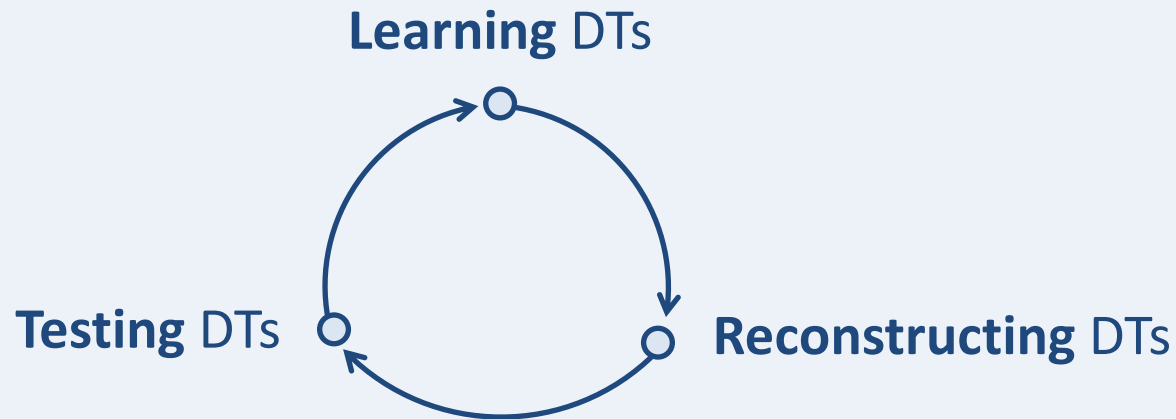


- Fast to evaluate
- Easy to understand, easy to explain predictions
- Algorithms widely employed, empirically successful

This talk:

Testing and Reconstructing decision trees

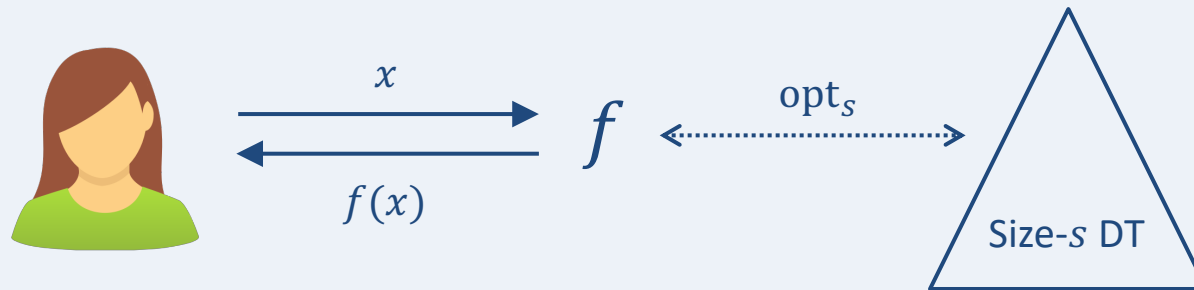
- Both tasks easier than learning
 - We draw on recent techniques from learning DTs
 - Our results have new implications for learning DTs



This talk: Surprisingly rich web of connections for DTs

Reconstruction: On-the-fly learning

Given query access to f , promised to be close to small decision tree:



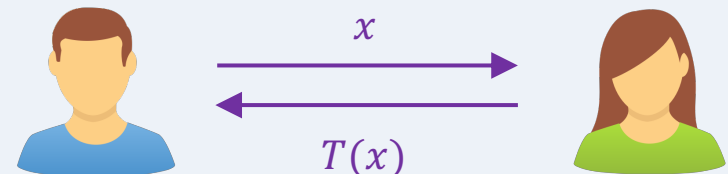
Traditional (Proper) Learning

Construct DT hypothesis:



Reconstruction: On-the-fly learning

Support queries to DT hypothesis:

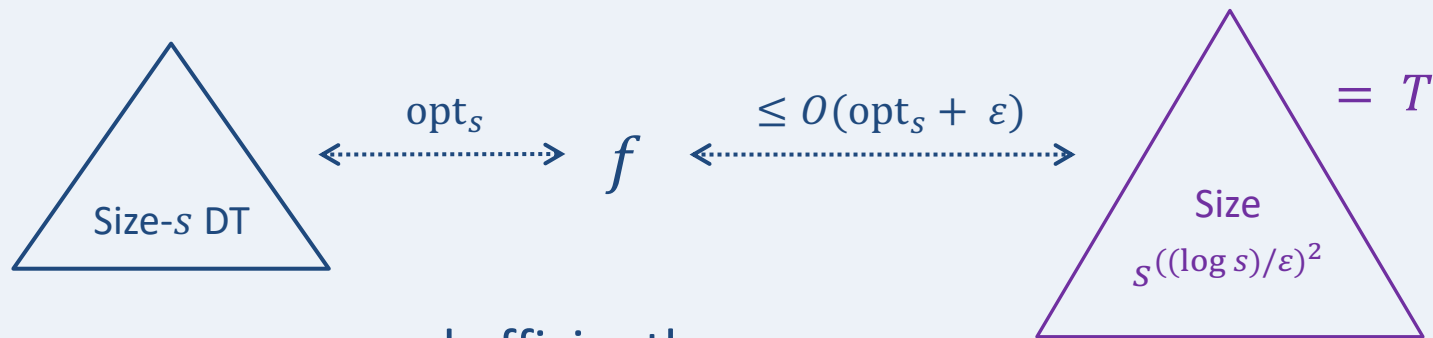


"Property-preserving data reconstruction" Ailon, Chazelle, Seshadhri, Liu, 2004

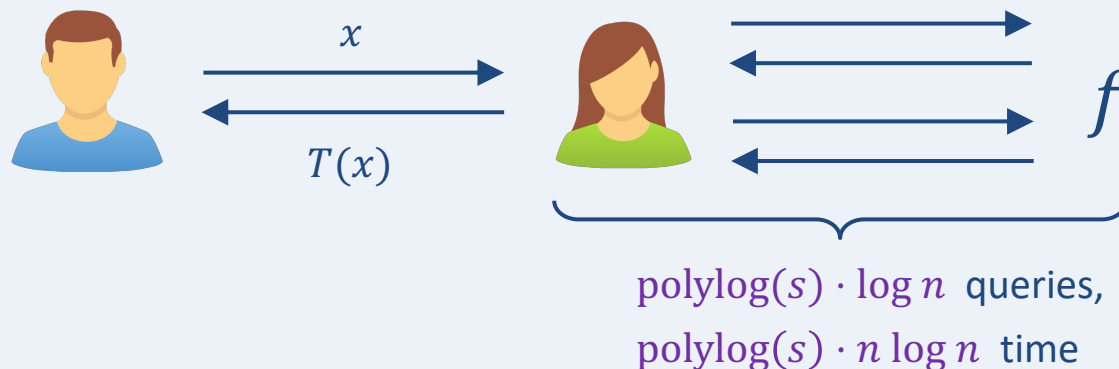
"On the testability and repair of hereditary hypergraph properties" Austin, Tao, 2008

Main result: Reconstruction algorithm for DTs

Given query access to $f: \{0,1\}^n \rightarrow \{0,1\}$, promised to be opt_s -close to size- s DT. We support queries to a DT hypothesis T :



Every query answered efficiently:

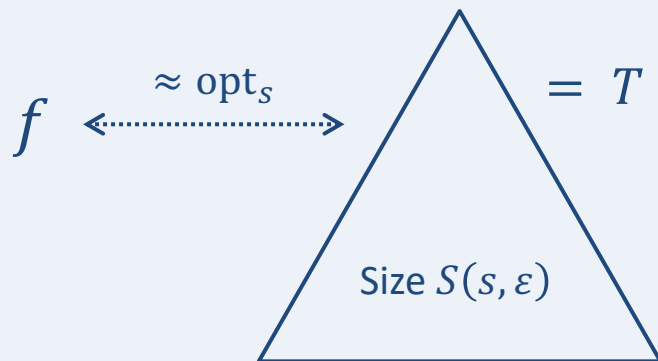


Traditional vs. On-the-fly learning of DTs

Both cases: Given query access to f , promised opt_s -close to size- s DT

Traditional (Proper) Learning

Construct DT hypothesis:



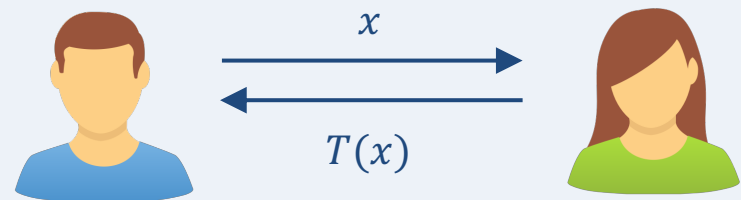
Fact: Need

$\Omega(s)$ queries to f

$\Omega(s) \cdot n$ time

Reconstruction: On-the-fly learning

Support queries to DT hypothesis:



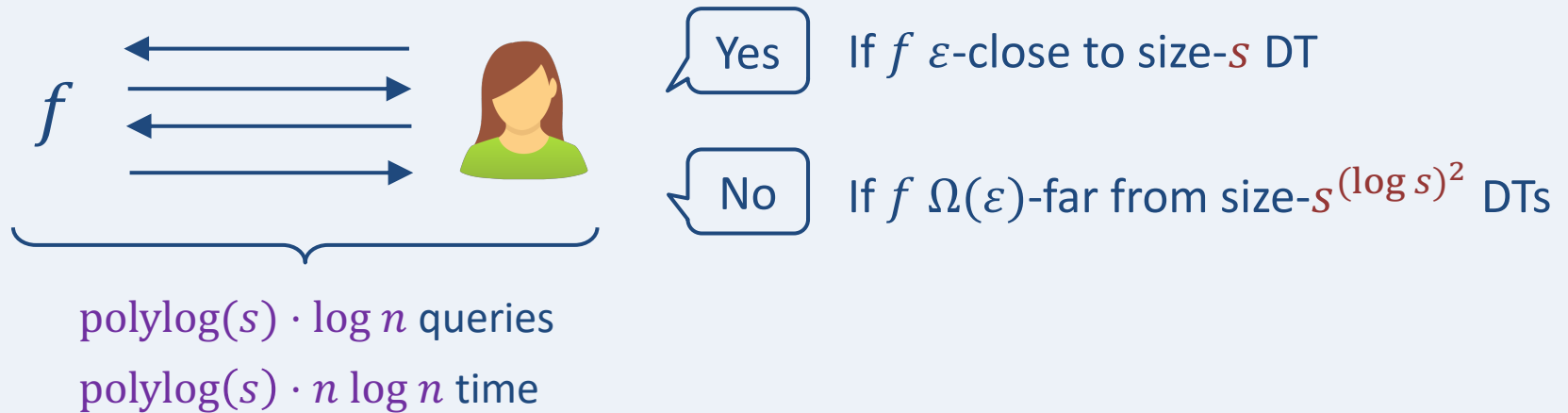
Our result: Each query to T answered with

$\text{polylog}(s) \cdot \log n$ queries to f

$\text{polylog}(s) \cdot n \log n$ time

Corollary: New tester for DTs

Given query access to unknown $f: \{0,1\}^n \rightarrow \{0,1\}$ and $s \in \mathbb{N}$,

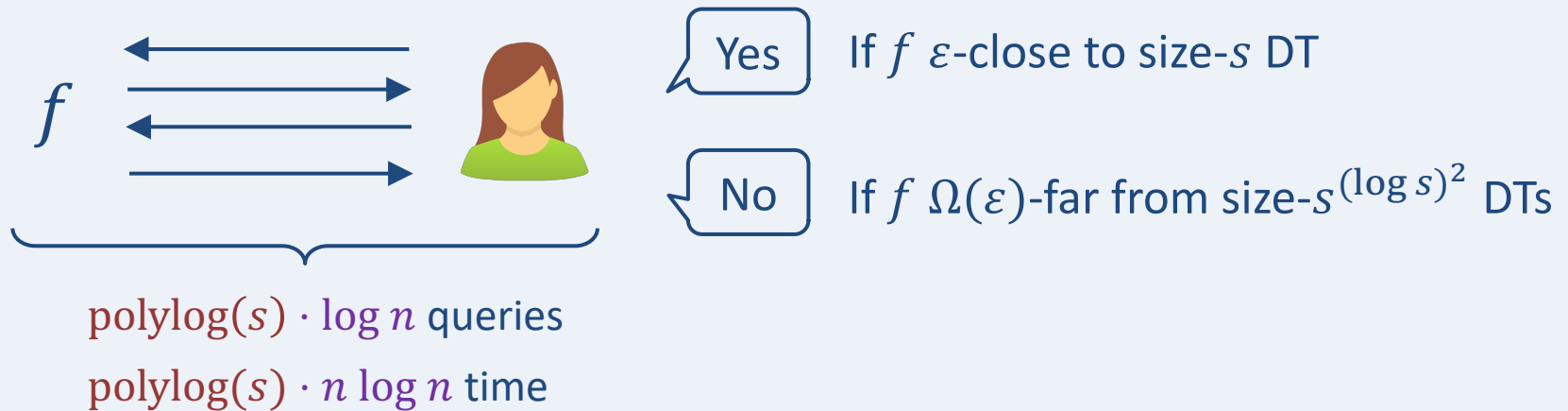


Adds to a long line of work on testing DTs:

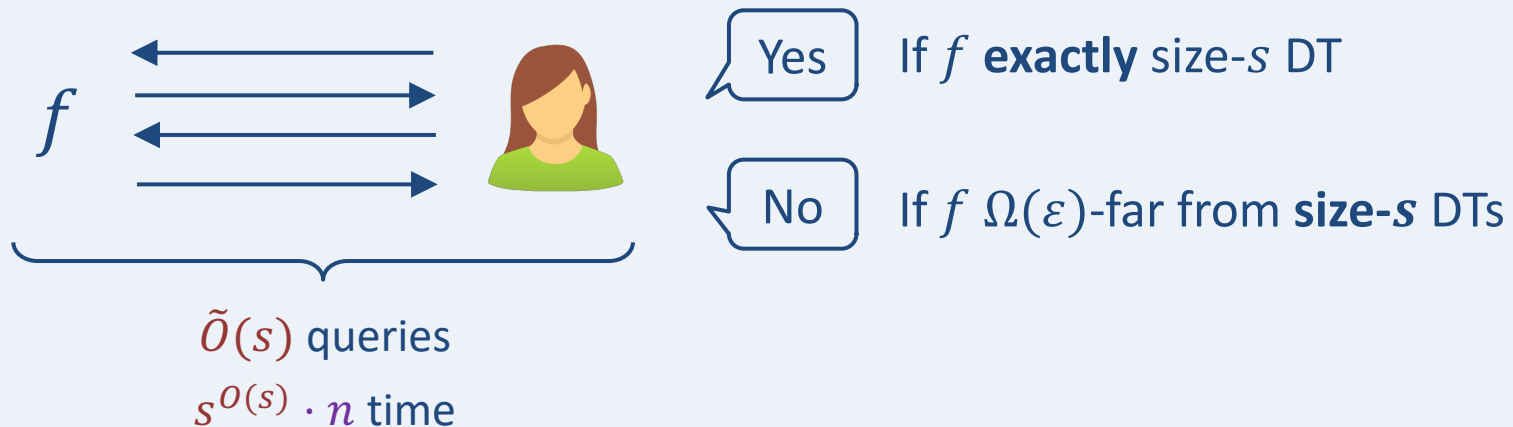
[KR00, DLMORSW07, CGSM11, BBM12, Bsh20]

Comparison with prior work

Our tester:

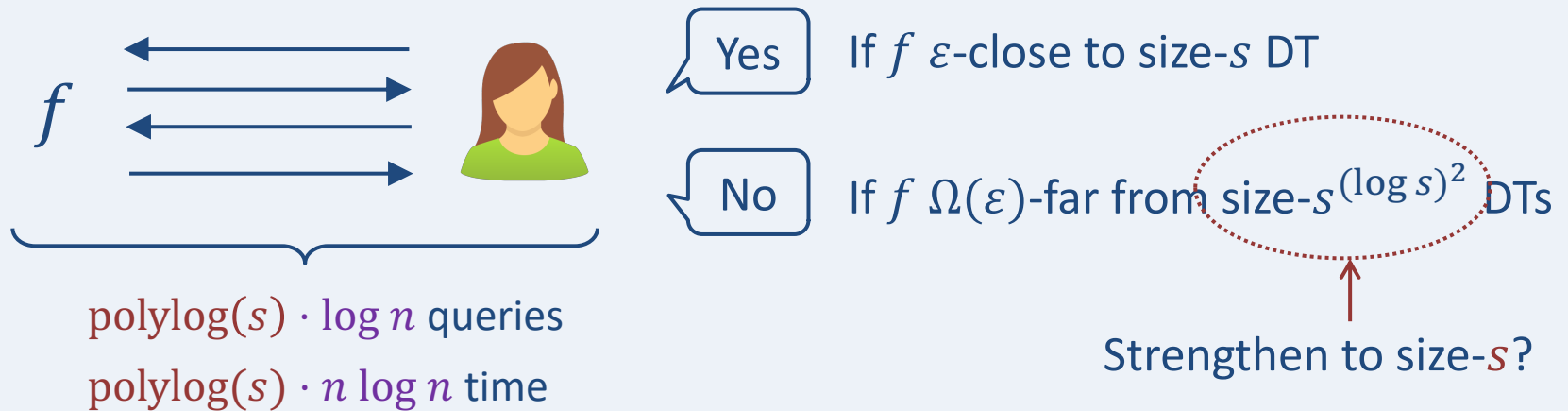


Existing testers: [DLMORSW07, CGSM11, Bsh20]

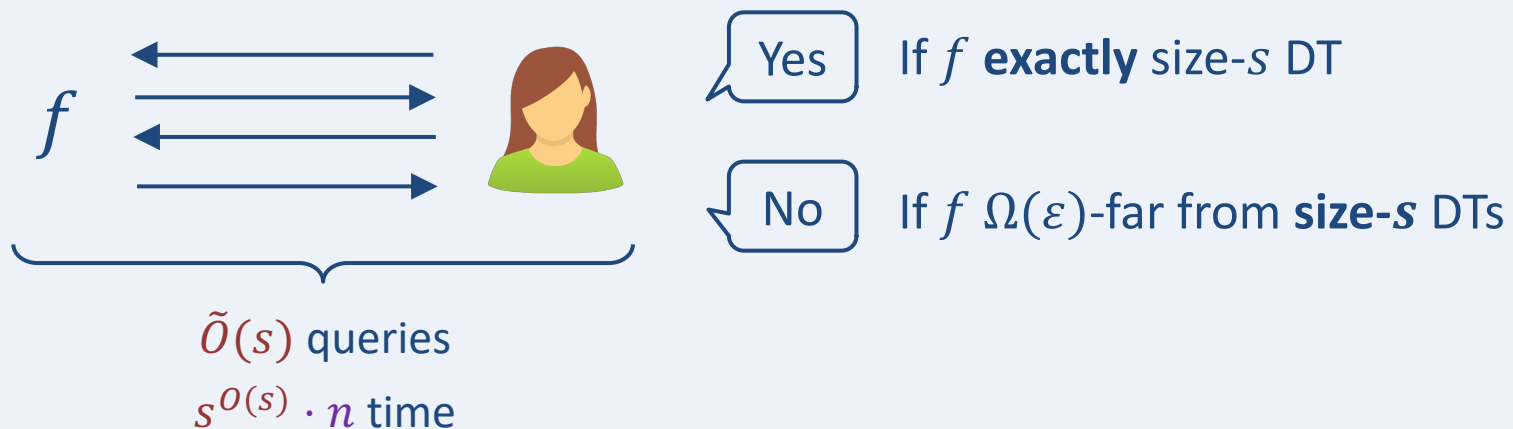


Comparison with prior work

Our tester:

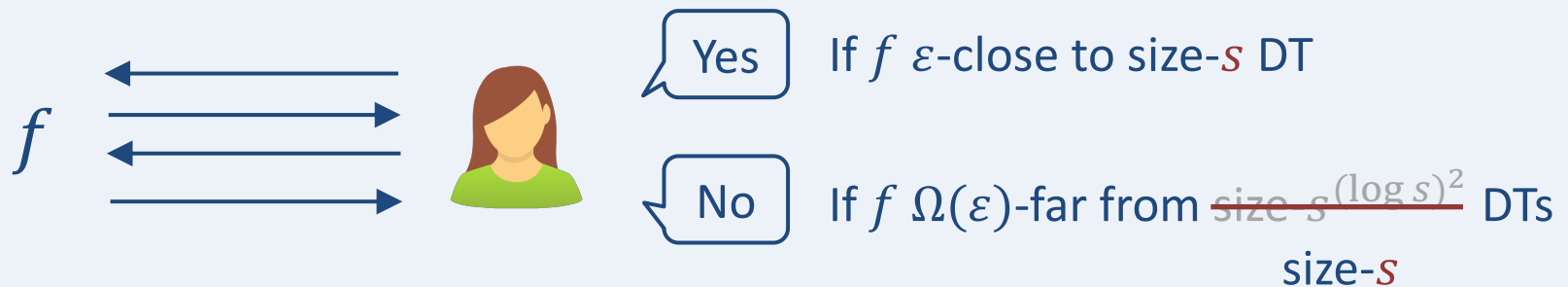


Existing testers: [DLMORSW07, CGSM11, Bsh20]



Improved tester \Rightarrow New learning algorithm

Suppose our **tester** can be improved to:

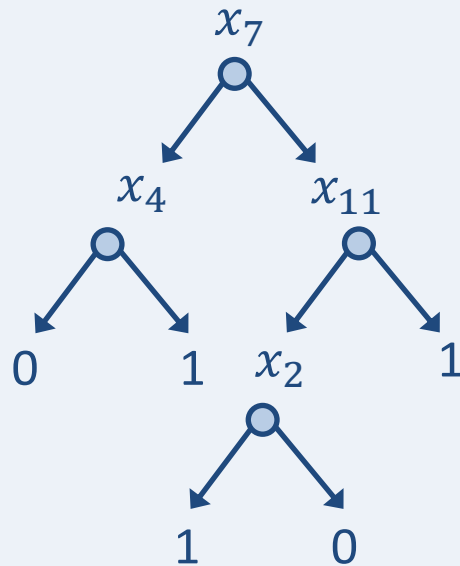


Then exists $\text{poly}(n, s, 1/\varepsilon)$ -time algorithm for **proper learning** of size- s DTs.

- Runtime of current best algorithm: $\text{poly}(n^{\log s}, 1/\varepsilon)$ [EH89]
- “Proper Learning \Rightarrow Testing” standard and long known [GGR98]
- This gives an example of “Testing \Rightarrow Proper Learning”

Reconstructors and testers for other properties

DT complexity closely related to many other measures:



- Fourier degree
- Approximate degree
- Randomized query complexity
- Quantum query complexity
- Sensitivity
- ...

Our results for DTs

⇒ Reconstructors and testers for these properties

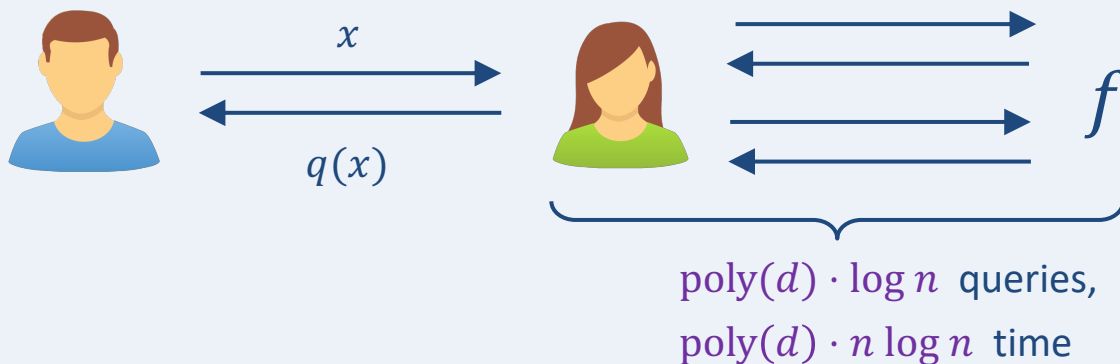
Example: Reconstructor for Fourier degree

Fact: For all f , we have $\deg(f) \leq \text{DT depth}(f) \leq \deg(f)^3$.

DT Reconstructor + Fact \uparrow : Given query access to $f: \{0,1\}^n \rightarrow \{0,1\}$, promised to be opt_d -close to a degree- d polynomial $p: \{0,1\}^n \rightarrow \{0,1\}$. We support queries to a **degree- $O(d^7)$** polynomial $q: \{0,1\}^n \rightarrow \{0,1\}$,

$$p \xleftrightarrow{\text{opt}_d} f \xleftrightarrow{\leq O(\text{opt}_d + \varepsilon)} q$$

Every query answered efficiently:



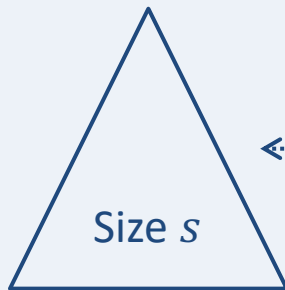
Outline of the rest of this talk

- Overview of our results
- **Key structural result and its proof**
- Our reconstruction algorithm
- Avenues for future work



Key structural result in a nutshell

Suppose f is opt_s -close to a size- s decision tree

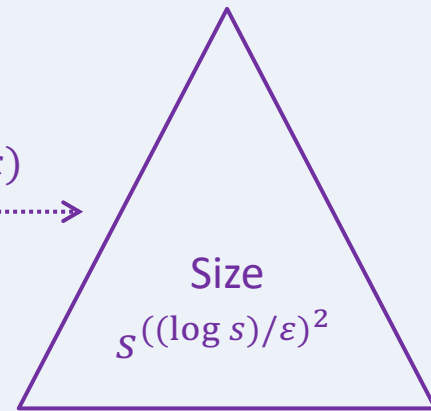


No idea what this
tree looks like

$\longleftrightarrow \text{opt}_s$

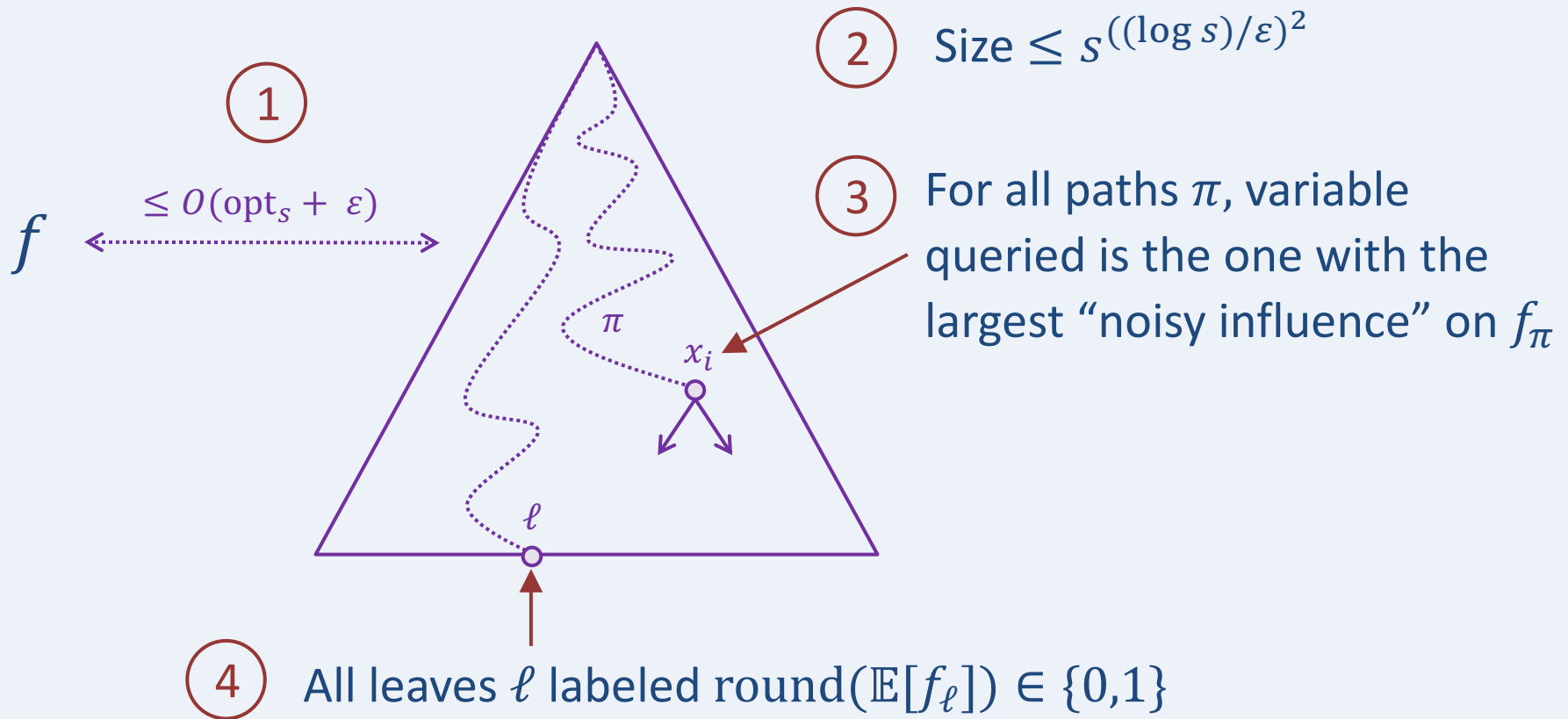
f

$\longleftrightarrow \leq O(\text{opt}_s + \varepsilon)$



Very specific structure,
Many enjoyable properties

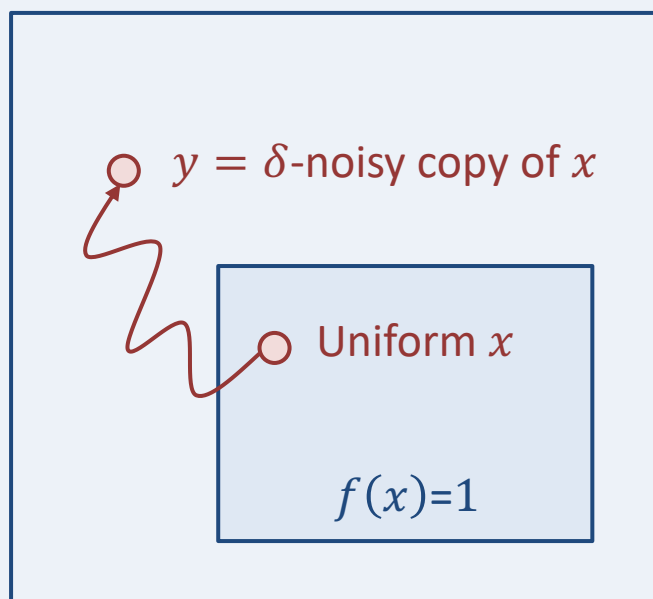
Structure and properties of this tree



Noise sensitivity and noisy influence

Def. Noise sensitivity of $f: \{0,1\}^n \rightarrow \{0,1\}$ at noise rate δ is the quantity:

$$\text{NS}_\delta(f) := \mathbb{P}_{y \sim_\delta x} [f(x) \neq f(y)]$$

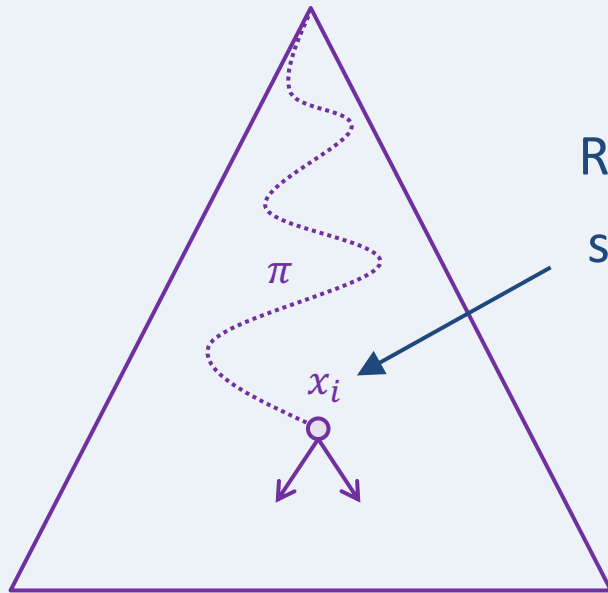


$\{0,1\}^n$

Def. The noisy influence of $i \in [n]$ on f is the quantity:

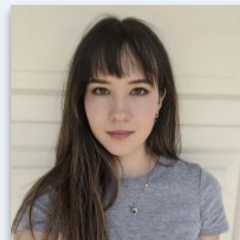
$$\text{NS}_\delta(f) - \mathbb{E}_{b \sim \{0,1\}} [\text{NS}_\delta(f_{x_i=b})]$$

Context: Splitting criteria of DT learning heuristics



Real-world heuristics (e.g. ID3, C4.5, CART)
split on x_i with largest **correlation** with f_π

Noisy influence = higher-order generalization of correlation
(Structure theorem false for correlation.)

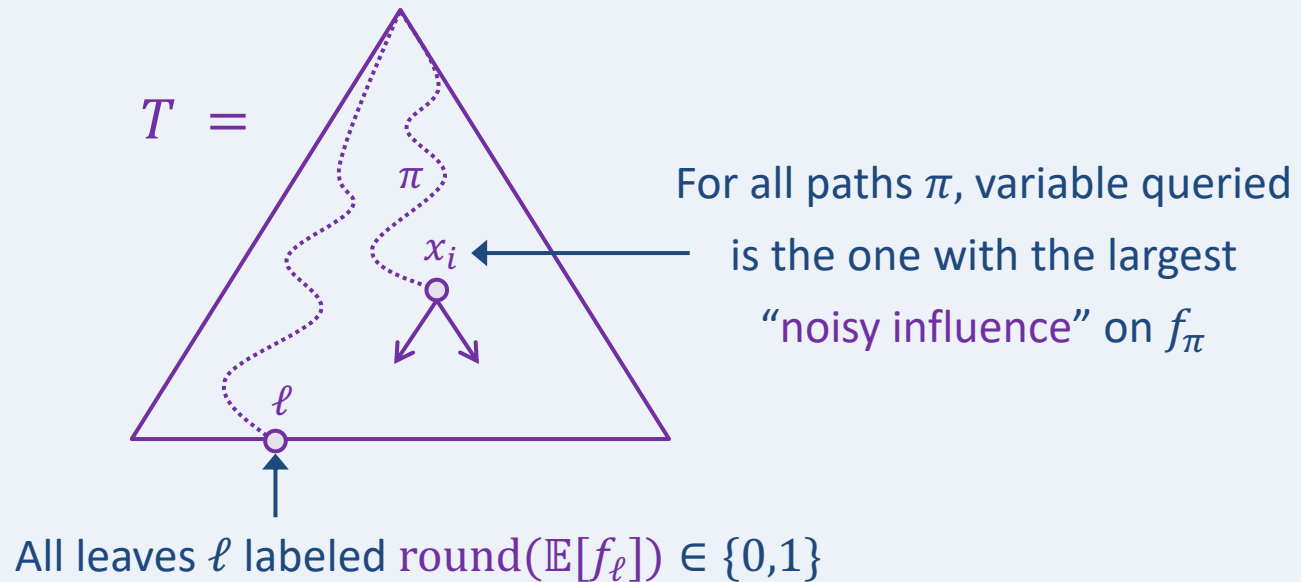


Blanc, Lange, Tan. ICML 2020
Blanc, Gupta, Lange, Tan. NeurIPS 2020

Our structural result, restated

Let f be opt_s -close to a size- s DT.

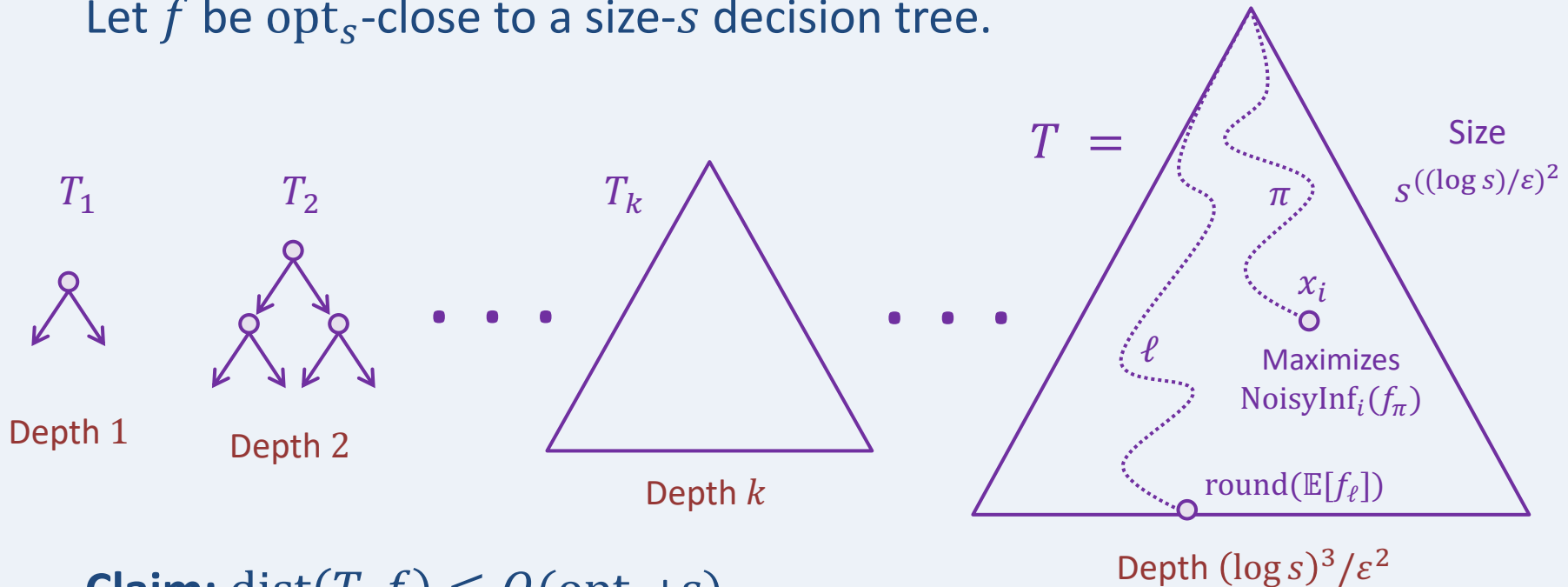
Consider the tree T of size $s^{((\log s)/\varepsilon)^2}$ defined as follows:



This tree is $O(\text{opt}_s + \varepsilon)$ -close to f .

Proof overview

Let f be opt_s -close to a size- s decision tree.



Claim: $\text{dist}(T, f) \leq O(\text{opt}_s + \epsilon)$

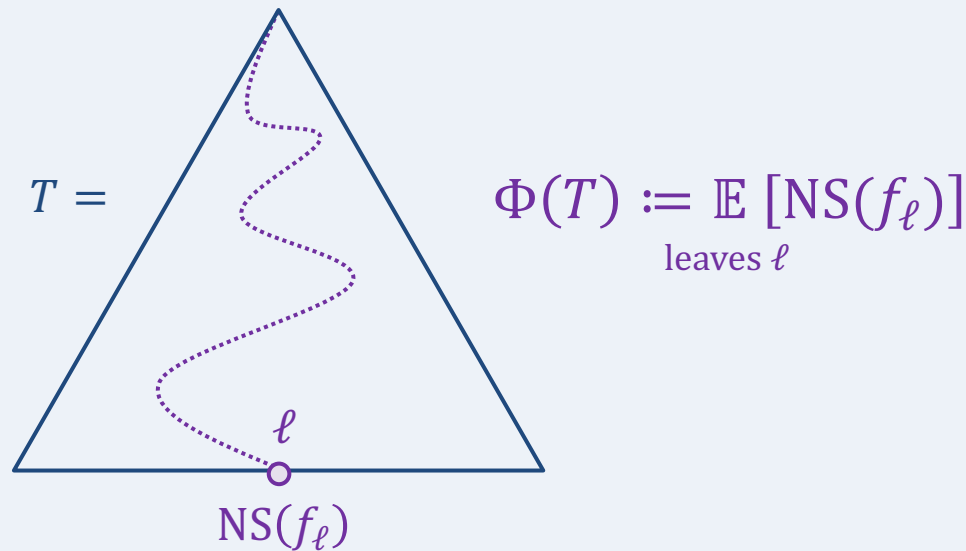
Proof strategy: Define potential function $\Phi : \text{Trees} \rightarrow [0, 1]$

Argue that for all k , either:

- Already done: $\text{dist}(T_k, f) \leq O(\text{opt}_s + \epsilon)$
- \Downarrow in potential: $\Phi(T_{k+1}) \leq \Phi(T_k) - \epsilon^2/(\log s)^3$

The potential function

$\Phi: \text{Trees} \rightarrow [0,1]$, $\Phi(T) = \text{Noise sensitivity of } f \text{ with respect to } T$

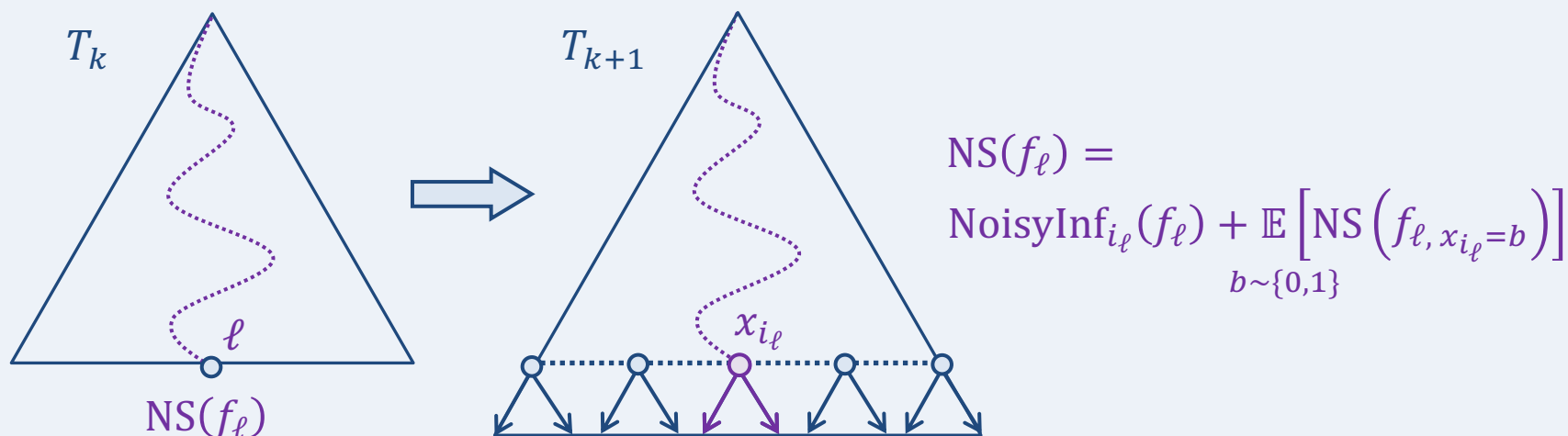


- Observations:
- $\Phi(\text{empty tree}) = \text{NS}(f) \leq 1$
 - $\Phi(T) \geq 0$ for all trees T

"A regularity lemma for low noisy-influences" O'Donnell, Servedio, Tan, Wan, 2010

"A noisy-influence regularity lemma for Boolean functions" Jones, 2016

Our potential function and our splitting criterion



$$\begin{aligned}
 & \Phi(T_k) - \Phi(T_{k+1}) \\
 &= \mathbb{E}_{\ell \sim T_k} [NS(f_\ell)] - \mathbb{E}_{\ell^* \sim T_{k+1}} [NS(f_{\ell^*})] \\
 &= \mathbb{E}_{\ell \sim T_k} [NoisyInf_{i_\ell}(f_\ell)]
 \end{aligned}$$

Our **splitting criterion** greedily drives down our **potential function**

Key lemma: Lower bound on noisy influence

Let f be opt_s -close to a size- s DT. Then:

\tilde{f} = smoothed version of f

$$\max_{i \in [n]} \{\text{NoisyInf}_i(f)\} \gtrsim \frac{\text{Var}(\tilde{f}) - \text{opt}_s}{(\log s)^2}$$

Hides dependence on noise rate

- Variant of the “OSSS inequality” from analysis of Boolean functions
- Applying this lemma:
 - $\text{Var}(\tilde{f}) > \text{opt}_s + \varepsilon \implies \text{RHS} > \varepsilon/(\log s)^2$
 - $\text{Var}(\tilde{f}) \leq \text{opt}_s + \varepsilon \implies f$ is $O(\text{opt}_s + \varepsilon)$ -close to constant

“Every decision tree has an influential variable” O’Donnell, Saks, Schramm, Servedio. FOCS 2005

“The influence of variables on Boolean functions”, Kahn, Kalai, Linial. FOCS 1988

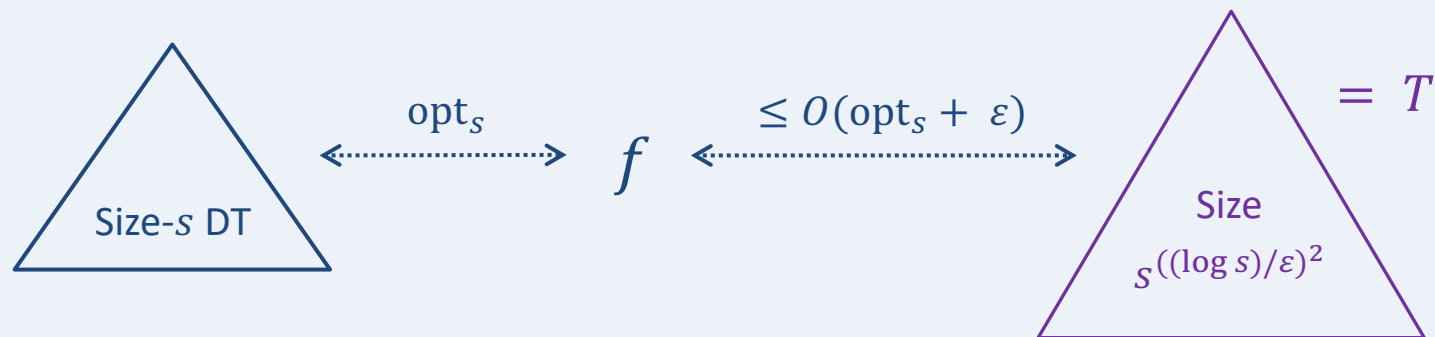
Outline of this talk

- Overview of our results
- Key structural result and its proof
- **Our reconstruction algorithm**
- Avenues for future work

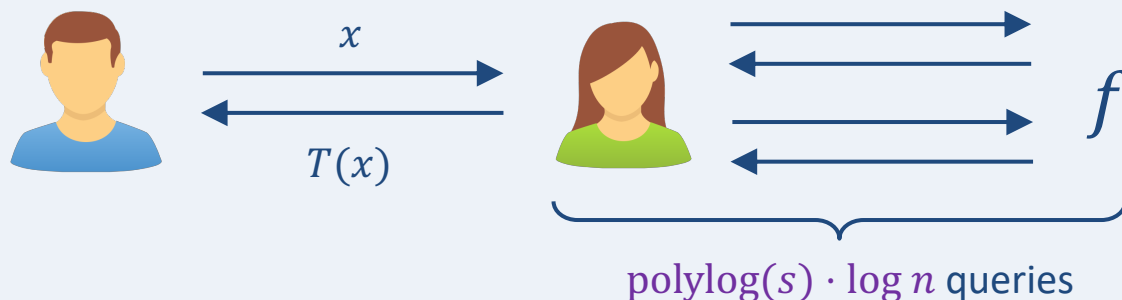


Recall our reconstruction algorithm for DTs

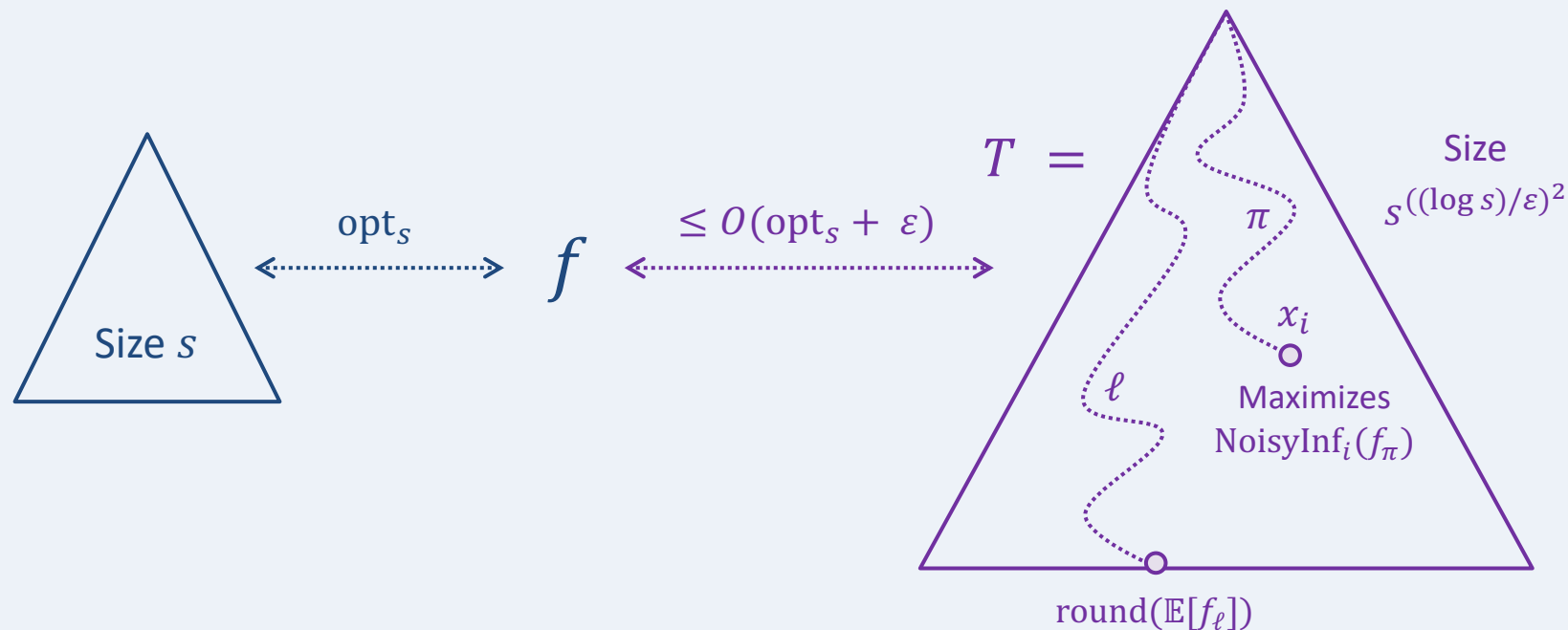
Given query access to $f: \{0,1\}^n \rightarrow \{0,1\}$, promised to be opt_s -close to size- s DT. We support queries to a DT hypothesis T :



Every query answered efficiently:



Algorithmic features of our structural result



Observation: Given query access to f , can construct T efficiently.

$$\text{NoisyInf}_i(f_\pi) := \underbrace{\text{NS}(f_\pi)}_{\mathbb{P}_{y \sim x_\delta}[f_\pi(x) \neq f_\pi(y)]} - \mathbb{E}_b[\text{NS}(f_{\pi, x_i=b})]$$

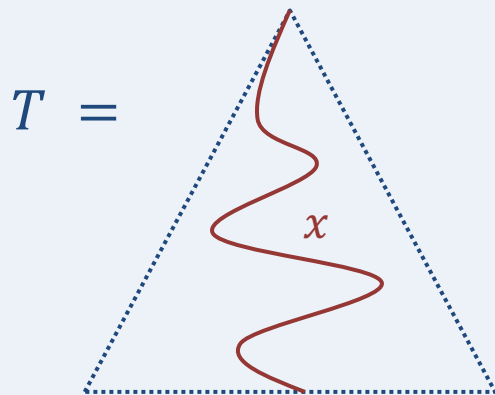
$$\mathbb{P}_{y \sim x_\delta}[f_\pi(x) \neq f_\pi(y)]$$

Evaluating T on a specific input x

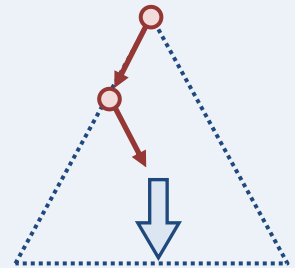
Previous slide: Given query access to f , can construct T — in full.

In fact, given query access to f and an input x , can compute $T(x)$ without constructing T in full.

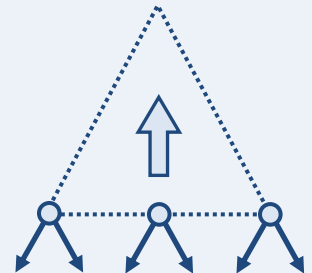
Build only the path in
 T that x follows:



Key enabling feature of T :
top-down, inductive definition



Cf. *bottom-up*, backtracking
DT learning algorithms
(e.g. Ehrenfeucht–Haussler 89)

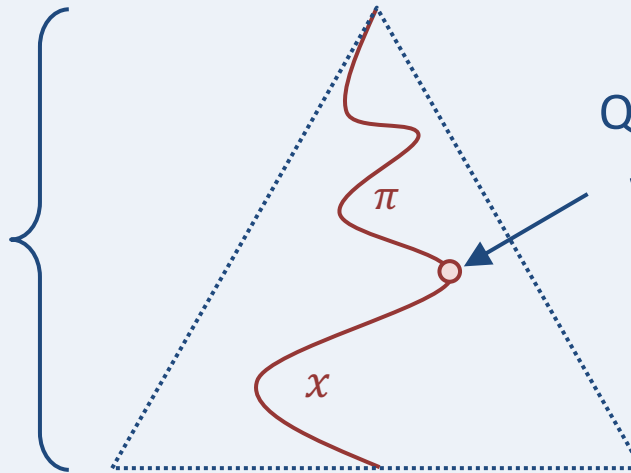


The spirit of Local Computation Algorithms [Rubinfeld et al. 11]

Query complexity of our reconstructor

Claim: For any input x , can compute $T(x)$ using $\text{polylog}(s) \cdot \log n$ queries to f

Structural Lemma
 $\Rightarrow \text{polylog}(s)$ depth



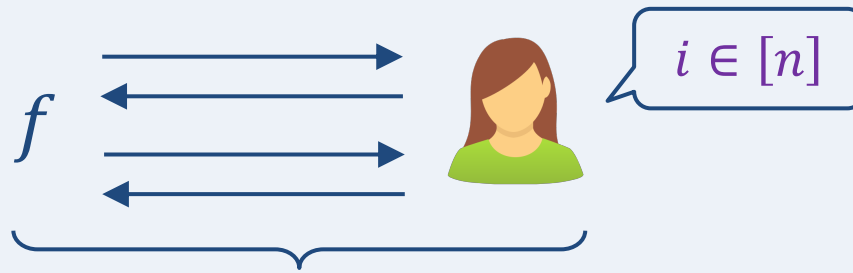
Query complexity of figuring out which variable has the largest noisy influence on f_π ?

Challenge: There are n variables.

Estimating n noisy influences $\Rightarrow \Omega(n)$ query complexity?

Finding the variable with largest noisy influence

Task: Given query access to $f: \{0,1\}^n \rightarrow \{0,1\}$,



As few queries as possible

With high probability,

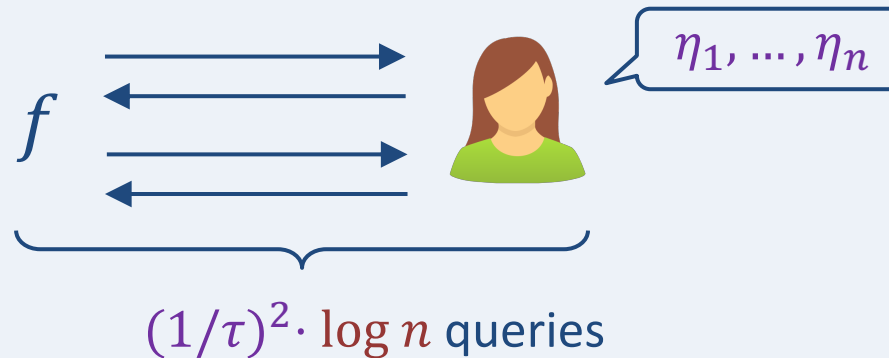
$$\text{NoisyInf}_i(f) \geq \text{NoisyInf}_j(f) \quad \text{for all } j \in [n].$$

Challenge: There are n variables.

Estimating n noisy influences $\Rightarrow \Omega(n)$ query complexity?

Query-efficient **simultaneous** estimation of noisy influences

Lemma: Given query access to $f: \{0,1\}^n \rightarrow \{0,1\}$,



With high probability,

$$\eta_i = \text{NoisyInf}_i(f) \pm \tau \quad \text{for all } i \in [n].$$

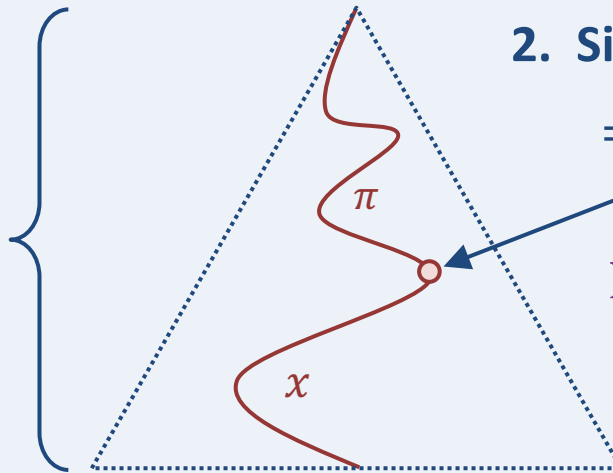
Crux of proof: 2-query unbiased estimator

Zooming out: Two main components of our proof

Claim: For any input x , can compute $T(x)$ using $\text{polylog}(s) \cdot \log n$ queries to f

Proof: Build only the path in T that x follows:

1. Structural lemma
 $\Rightarrow \text{polylog}(s)$ depth

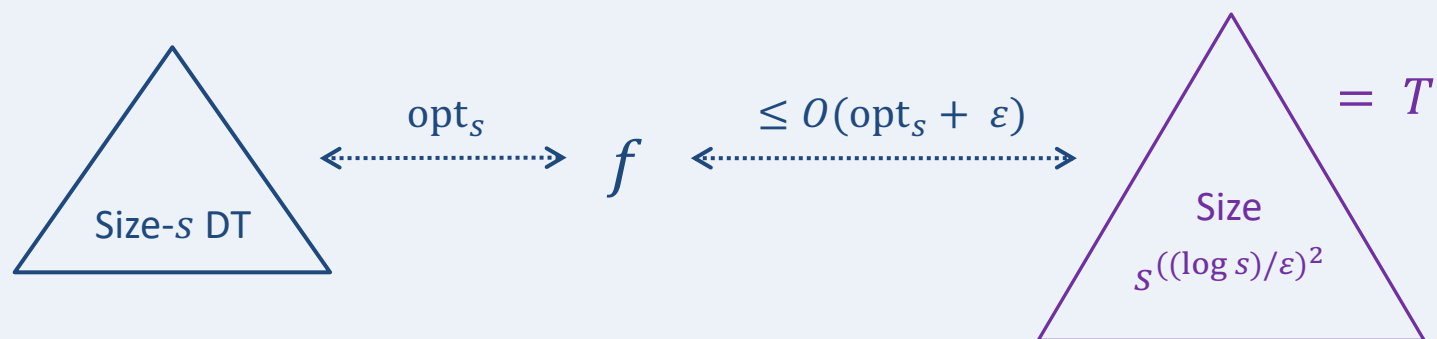


2. Simultaneous estimation algorithm

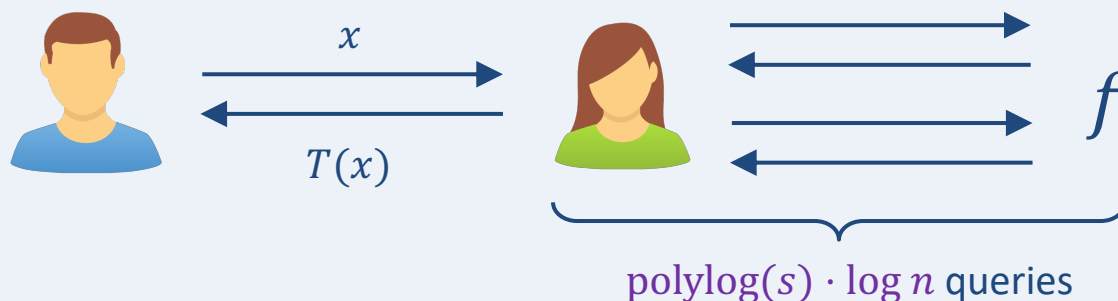
\Rightarrow Identify variable with largest
noisy influence on f_π with
 $\text{polylog}(s) \cdot \log n$ queries to f

Wrapping up: Reconstruction algorithm for DTs

Given query access to $f: \{0,1\}^n \rightarrow \{0,1\}$, promised to be opt_s -close to size- s DT. We support queries to a DT hypothesis T :



Every query answered efficiently:

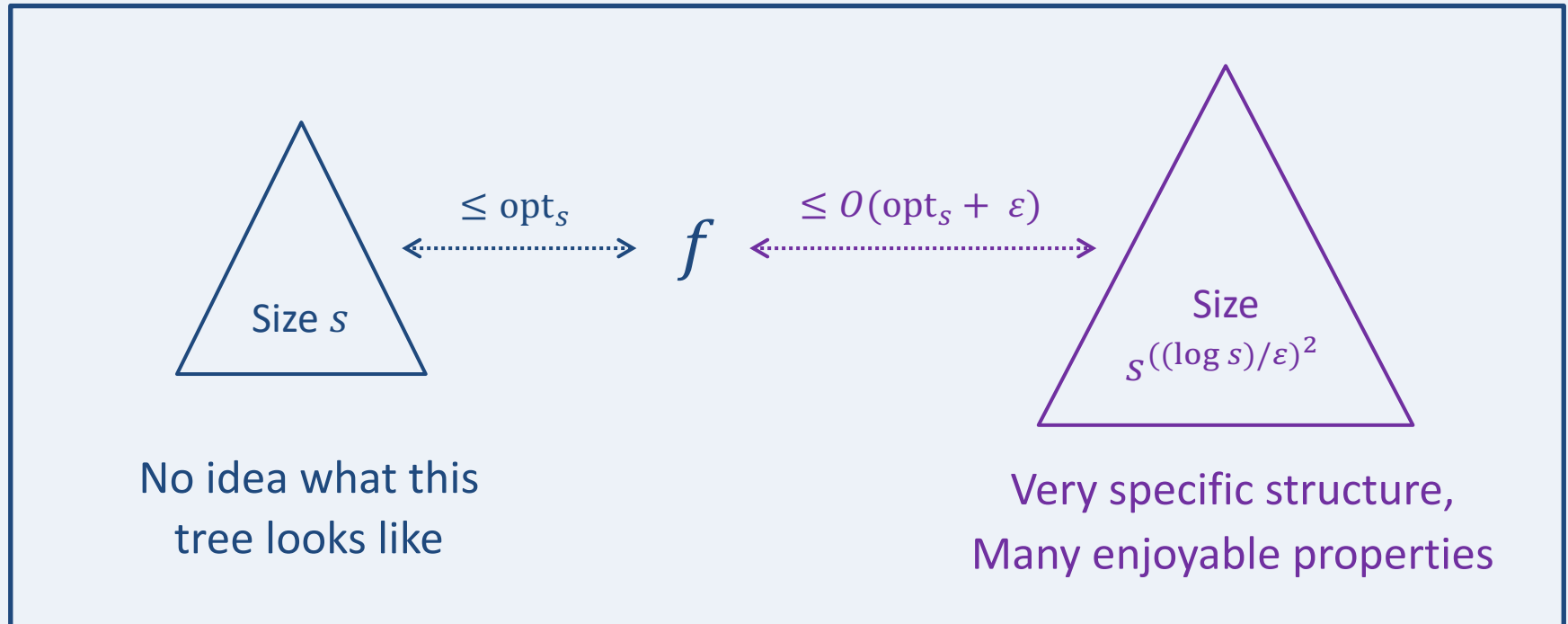


Outline of this talk

- Overview of our results
- Key structural result and its proof
- Our reconstruction algorithm
- **Avenues for future work**



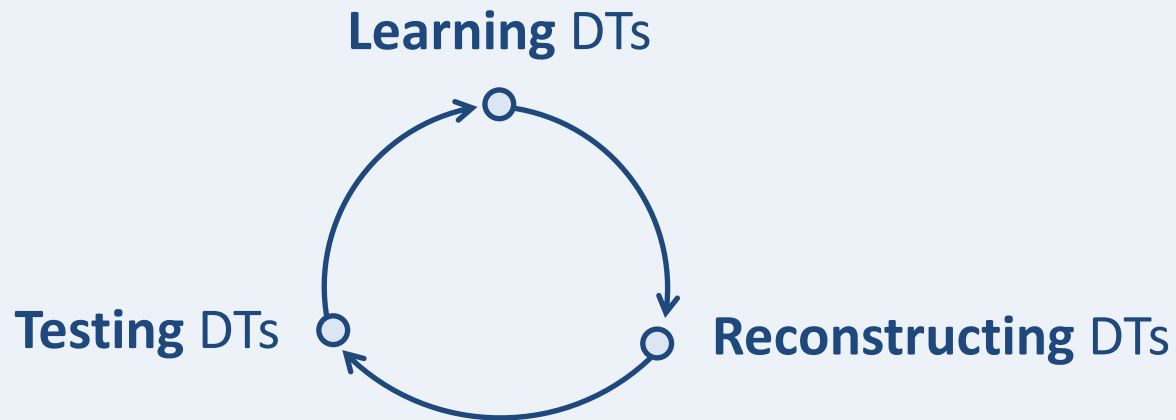
Further applications of our structural result?



- Inspired by real-world DT **learning** heuristics
- This talk: Applications to **reconstruction** and **testing**

Learning, Testing, and Reconstruction

- Generic algorithmic tasks, well-studied for many classes
- Surprisingly rich web of connections for the class of DTs



Much more to be understood, quantitatively and qualitatively

Understanding practical DT learning heuristics

- Rigorous guarantees and inherent limitations?
- Theory of splitting criteria?
- Random forests and boosted DTs?



“In summary, it seems fair to say that despite their other successes, the models of computational learning theory have not yet provided significant insight into the apparent empirical successes of programs like C4.5 and CART.”

– Kearns and Mansour

On the boosting ability of top-down decision tree learning algorithms, STOC 1996

Thank you for listening.

