# Gaining Free or Low-Cost Transparency with Interpretable Partial Substitute

Tong Wang

University of Iowa
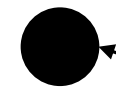
Tippie College of Business

tong-wang@uiowa.edu
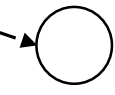
A black-box
model

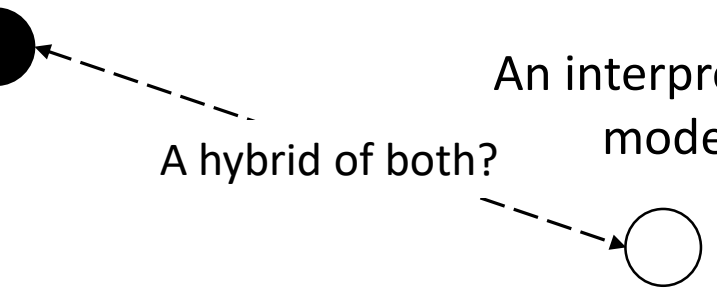+ High predictive performance

- non-interpretable

A hybrid of both?

An interpretable
model

+ interpretable

- lower predictive performance

A black-box model

**+** High predictive performance

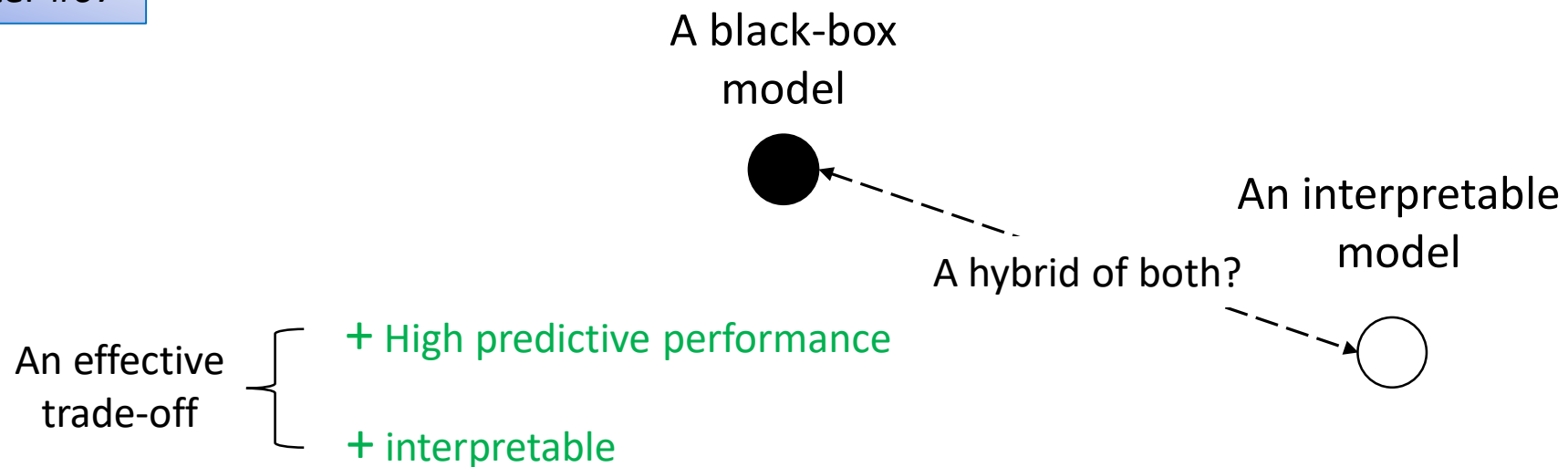**-** non-interpretable

A hybrid of both?
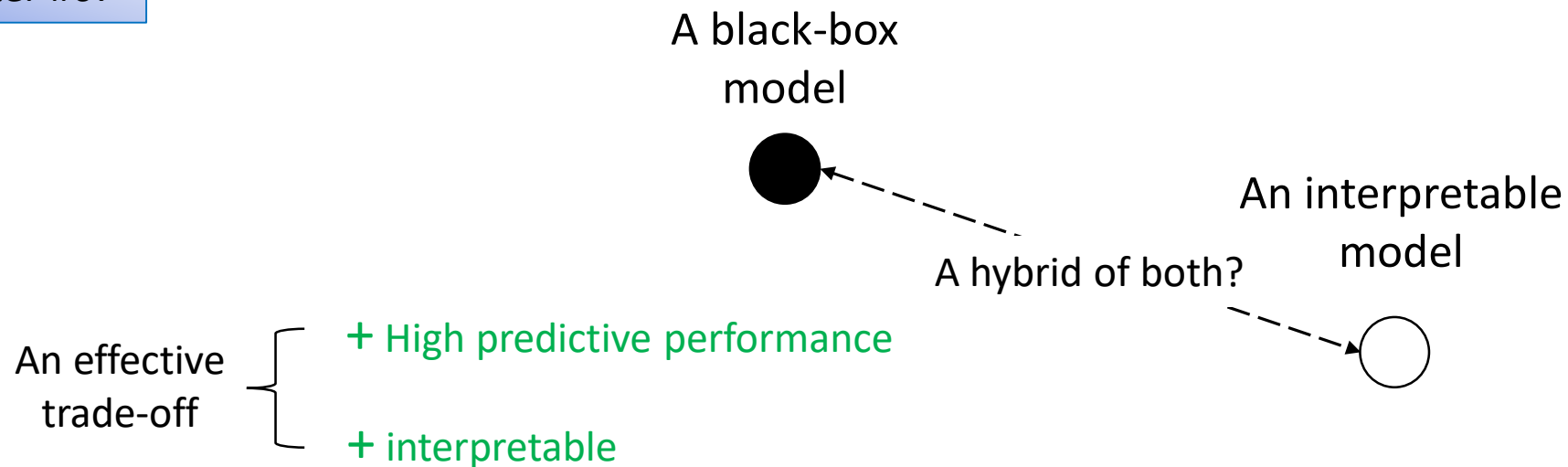
An interpretable model

**+** interpretable

**-** lower predictive performance

A key observation: there might exist a subspace where a black-box is *overkill* and a simple interpretable model can perform just as well as the black-box

A black-box
model

An interpretable
model

A hybrid of both?

An effective
trade-off
+ High predictive performance

+ interpretable

A key observation: there might exist a subspace where a black-box is *overkill* and a simple interpretable model can perform just as well as the black-box

The proposed solution: to substitute the black-box model with an interpretable model, where there is no or low-cost of predictive performance

A black-box
model

An interpretable
model

A hybrid of both?

An effective
trade-off
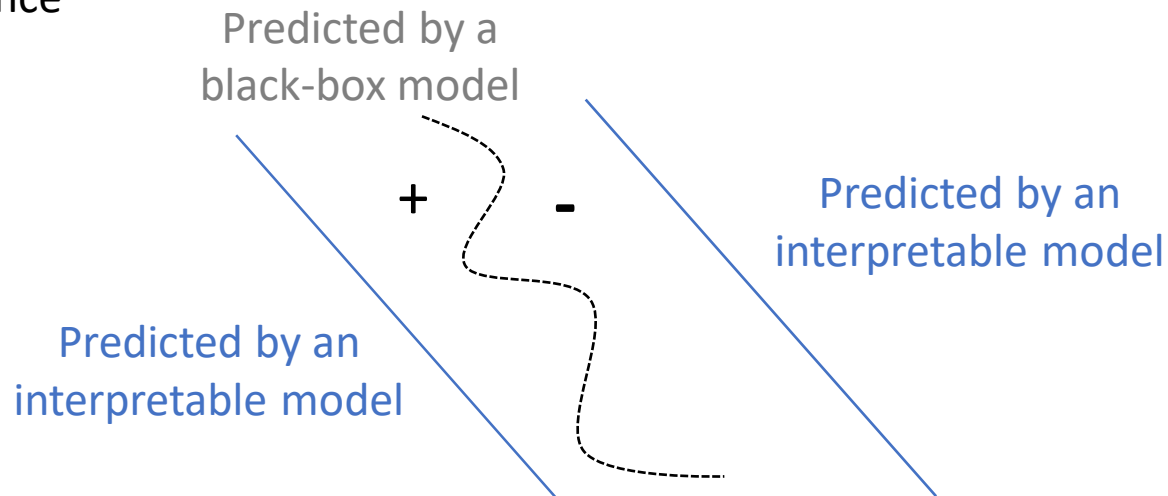
+ High predictive performance

+ interpretable

A key observation: there might exist a subspace where a black-box is *overkill* and a simple interpretable model can perform just as well as the black-box

The proposed solution: to substitute the black-box model with an interpretable model, where there is no or low-cost of predictive performance
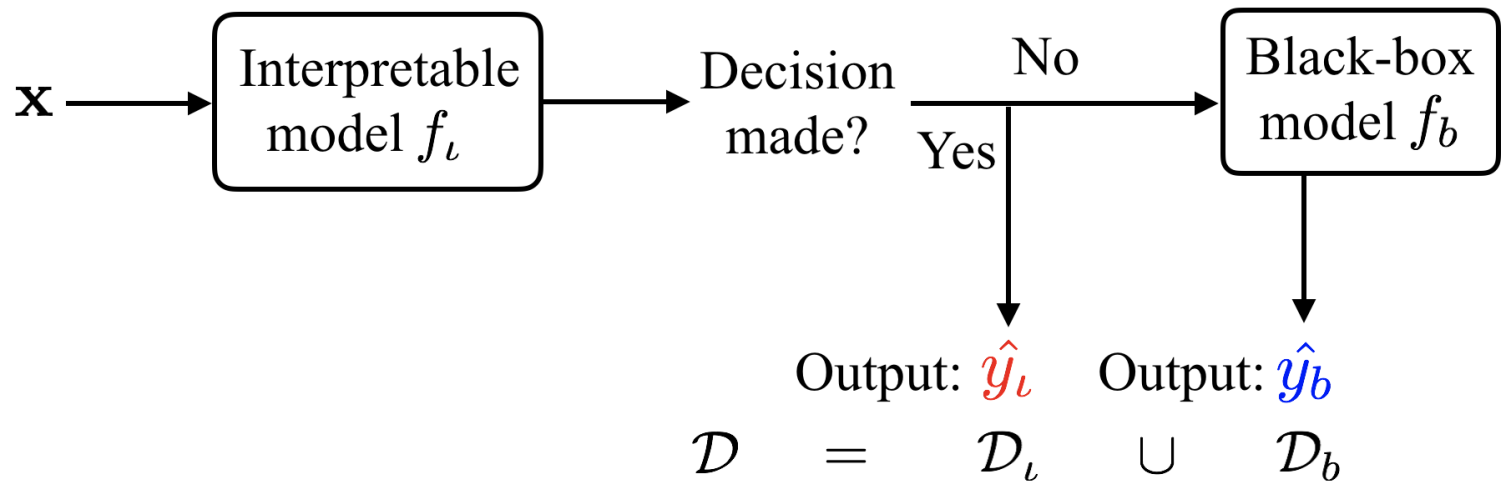
Predicted by a
black-box model

Predicted by an
interpretable model

+

-

Predicted by an
interpretable model

A hybrid predictive model

$\mathbf{x} \longrightarrow$ Interpretable model $f_\iota$ $\longrightarrow$ Decision made?

No $\longrightarrow$ Black-box model $f_b$

Yes

Output: $\hat{y}_\iota$     Output: $\hat{y}_b$
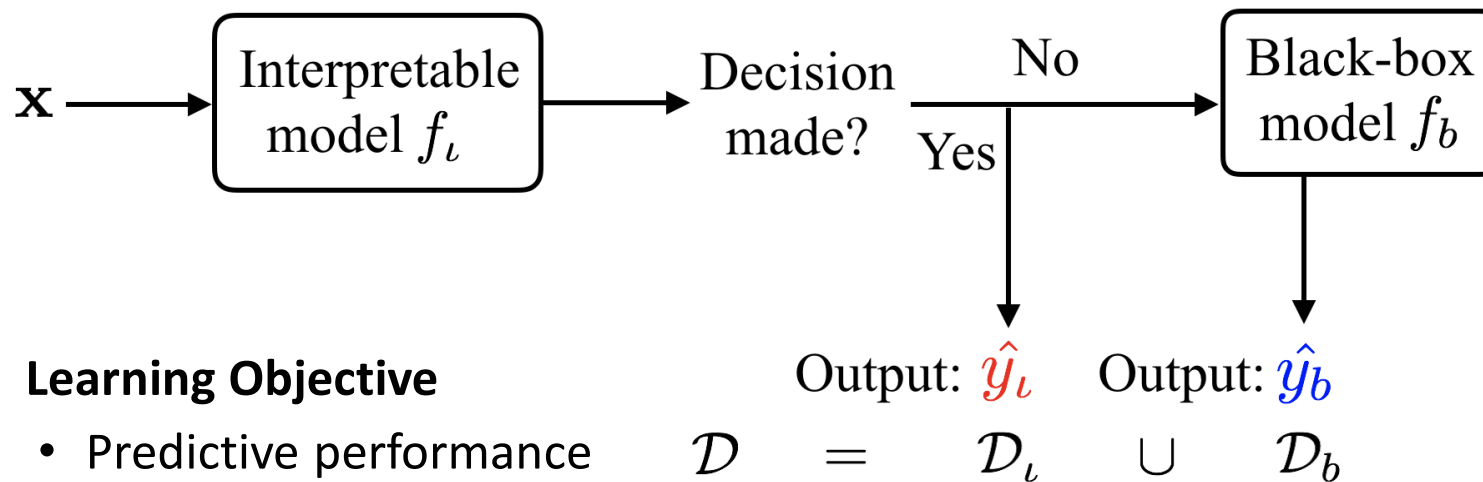
$\mathcal{D}$     $=$     $\mathcal{D}_\iota$     $\cup$     $\mathcal{D}_b$

Define transparency of model: $\frac{D_i}{D}$

**A hybrid predictive model**

**Learning Objective**

- Predictive performance
- Interpretability of $f_i$
- Transparency

$$\mathbf{x} \rightarrow \boxed{\begin{array}{c}\text{Interpretable} \\ \text{model } f_\iota\end{array}} \rightarrow \begin{array}{c}\text{Decision} \\ \text{made?}\end{array}$$

No $\rightarrow$ $\boxed{\begin{array}{c}\text{Black-box} \\ \text{model } f_b\end{array}}$

Yes

Output: $\hat{y}_\iota$    Output: $\hat{y}_b$

$$\mathcal{D} = \mathcal{D}_\iota \cup \mathcal{D}_b$$

Define transparency of model: $\frac{D_i}{D}$

A hybrid predictive model

$\mathbf{x}$ → Interpretable model $f_\iota$ → Decision made? → No → Black-box model $f_b$

Yes

Output: $\hat{y_\iota}$   Output: $\hat{y_b}$

$\mathcal{D} \quad = \quad \mathcal{D}_\iota \quad \cup \quad \mathcal{D}_b$

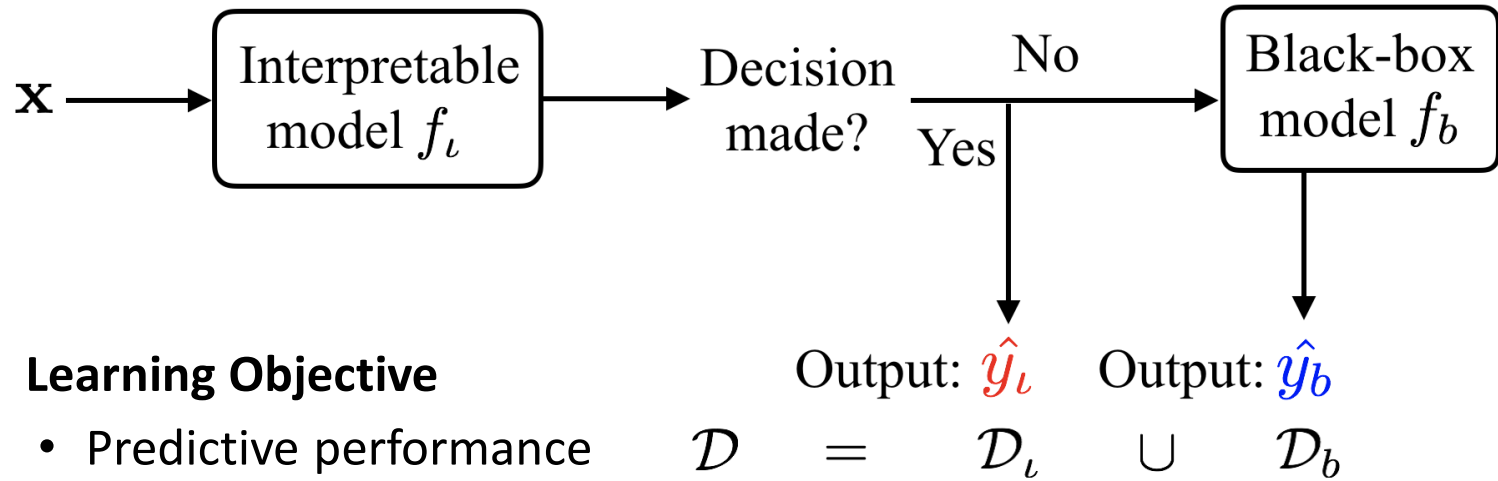**Learning Objective**

- Predictive performance
- Interpretability of $f_i$
- Transparency

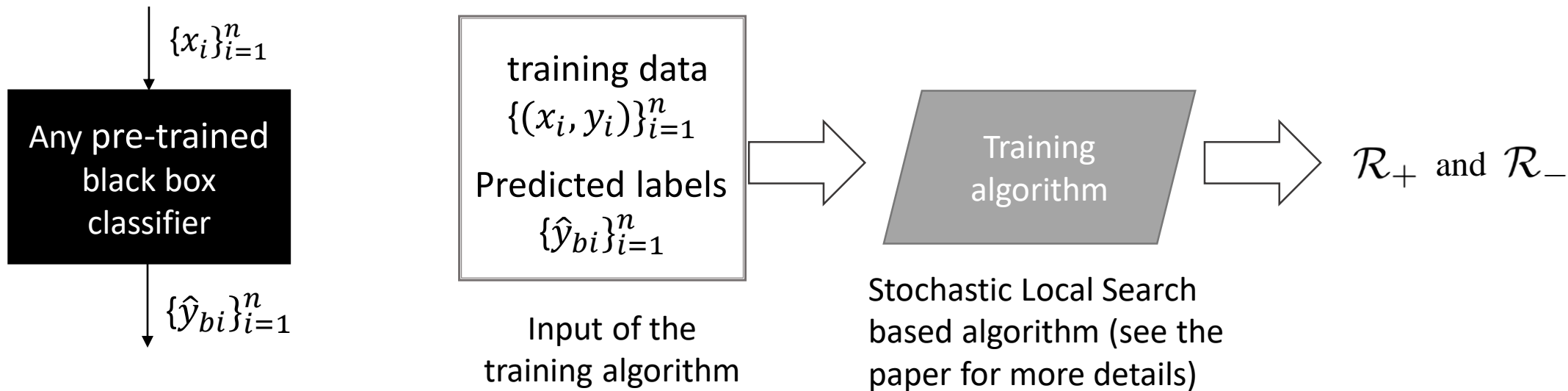Define transparency of model: $\frac{D_i}{D}$

A Hybrid Rule Set

**if** $\mathbf{x}_i$ obeys $\mathcal{R}_+, Y = 1$

**else if** $\mathbf{x}_i$ obeys $\mathcal{R}_-, Y = 0$

**else** $Y = f_b(\mathbf{x}_i)$

*Table 1.* An example of a HyRS model

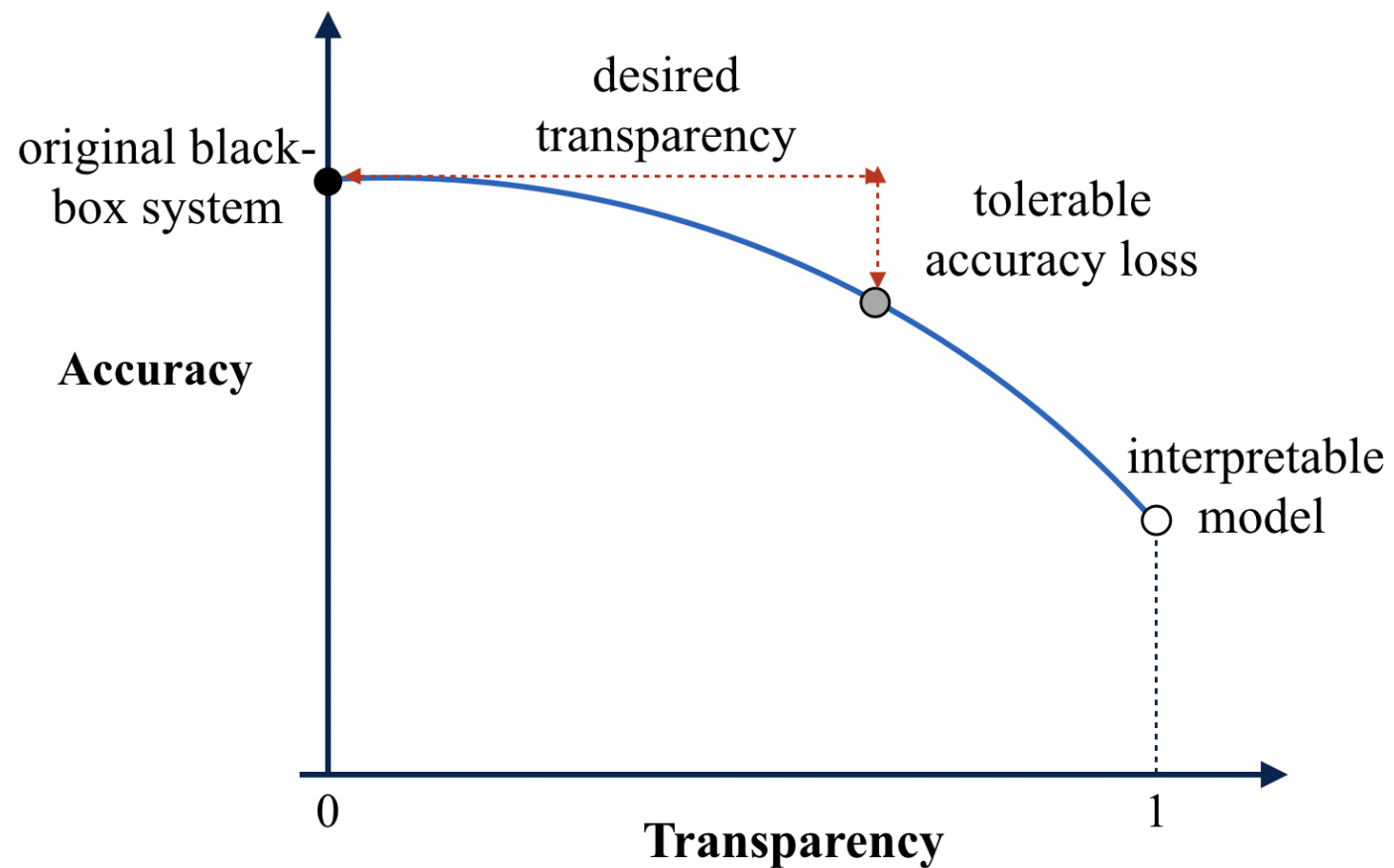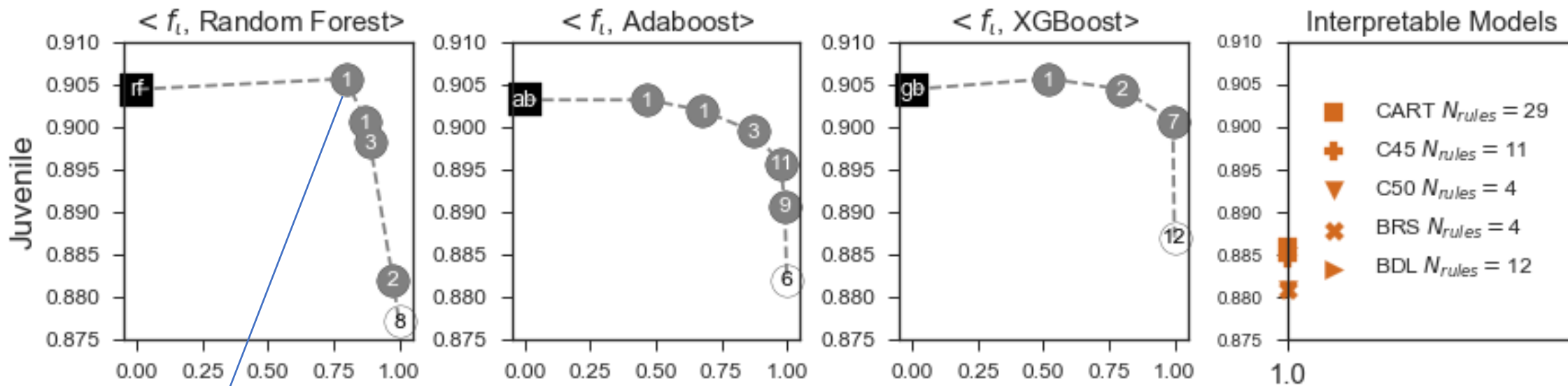| | Rules | Model |
|---|---|---|
| **if** | age$< 35$ *and* maximum heart rate $\geq 178$ OR serum cholestorol $\geq 234$ *and* thal $\neq 3$ *and* the number of vessels $\geq 1$ $\rightarrow Y = 1$ (heart disease) | $\mathcal{R}_+$ |
| **else if** | chest pain type $\neq 4$ *and* age $> 40$ $\rightarrow Y = 0$ ( no heart disease) | $\mathcal{R}_-$ |
| **else** | $\rightarrow Y = f_b(\mathbf{x})$ | $f_b$ |

# Model Training

$\{x_i\}_{i=1}^n$

Any pre-trained
black box
classifier

$\{\hat{y}_{bi}\}_{i=1}^n$

training data
$\{(x_i, y_i)\}_{i=1}^n$

Predicted labels
$\{\hat{y}_{bi}\}_{i=1}^n$

Input of the
training algorithm

Training
algorithm

Stochastic Local Search
based algorithm (see the
paper for more details)

$\mathcal{R}_+$ and $\mathcal{R}_-$
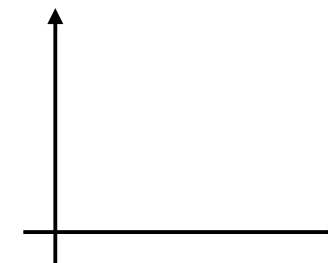
# Performance on Juvenile dataset



**if** Has any of your family members or friends ever attacked you with a weapon $\neq$ Yes *and* Have your friends ever hit or threatened to hit someone without any reason? $\neq$ Yes *and* Have your friends ever broken into a vehicle or building to steal something $\neq$ Yes
**then** $Y = 0$
**else** $Y = f_b(\mathbf{x})$

Thank you!

Poster #67 in Pacific Ballroom