



ACM SIGCOMM 2023 Tutorial: Closed-Loop “ML for Networks” Pipelines

Trustee: An Augmented ML Pipeline for Explaining ML Models

Ronaldo A Ferreira
UFMS

September 10, 2023



Traditional AI/ML Pipeline

- The “traditional AI/ML pipeline” consists of:
 - A training task characterized by a model specification
 - A (labelled) training dataset
 - An **independent and identically distributed (IID)** evaluation procedure
 - A (single) metric (e.g., F1-score) that measures the model’s expected predictive performance on **data drawn from the training distribution**

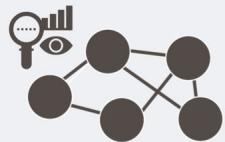
Traditional AI/ML Pipeline

- Main criticisms of the “traditional ML pipeline” include:
 - The pipeline’s output (i.e., best-performing ML model) is prone to suffer from the **problem of underspecification**
 - The evaluation of the pipeline’s output is agnostic to the particular **inductive biases** encoded by the trained model

How does it work?

Traditional AI/ML Development Pipeline

Collect Data

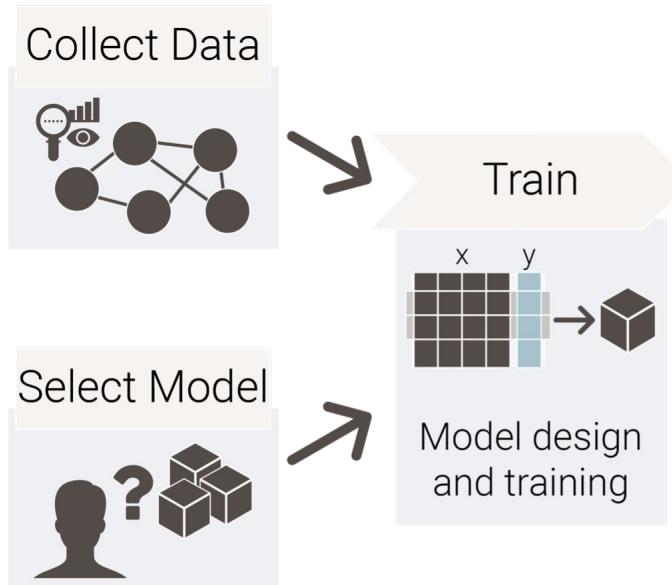


Select Model



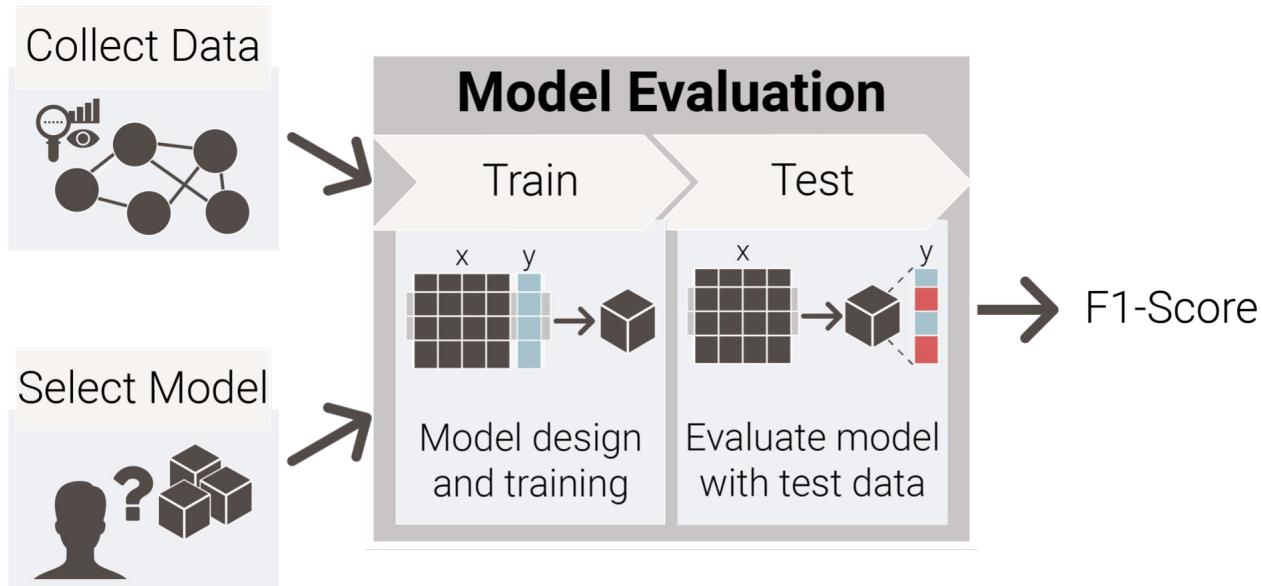
How does it work?

Traditional AI/ML Development Pipeline

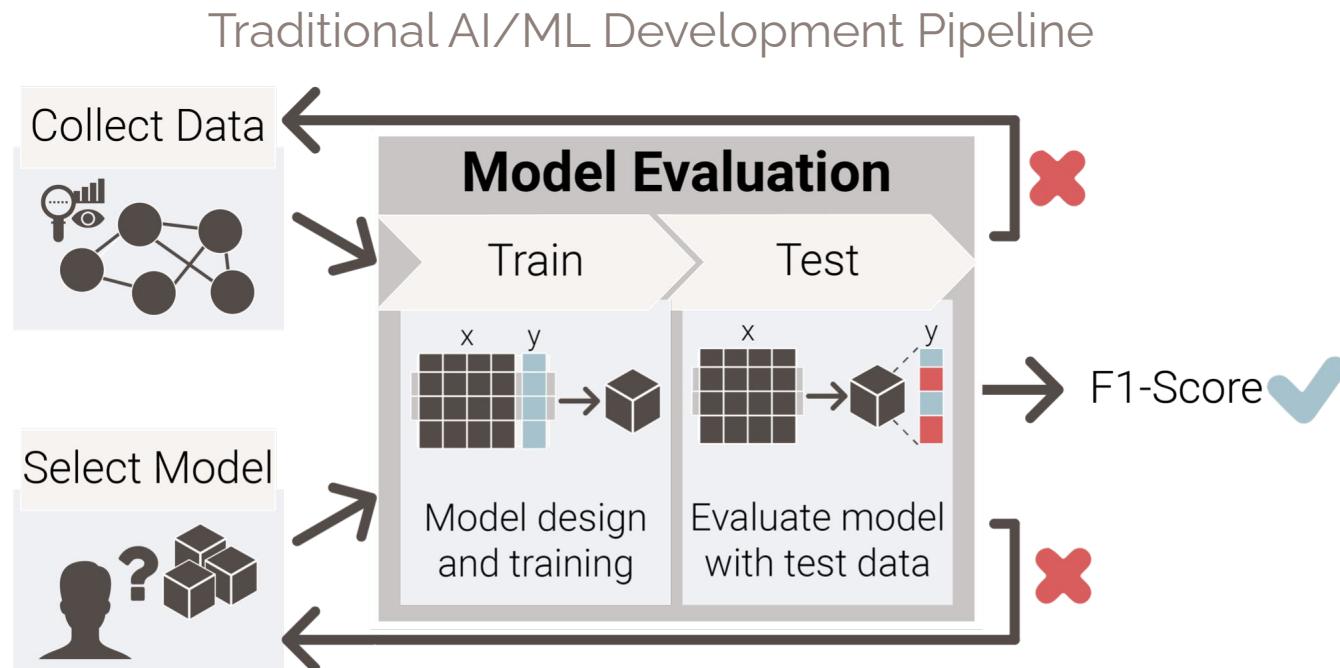


How does it work?

Traditional AI/ML Development Pipeline



How does it work?



What about high-stakes decision making?

Why (and how) does the model work?



Self-driving Cars

When does the model not work?



Network Security

Underspecification issues!

Shortcut Learning

Model takes shortcuts to classify data!

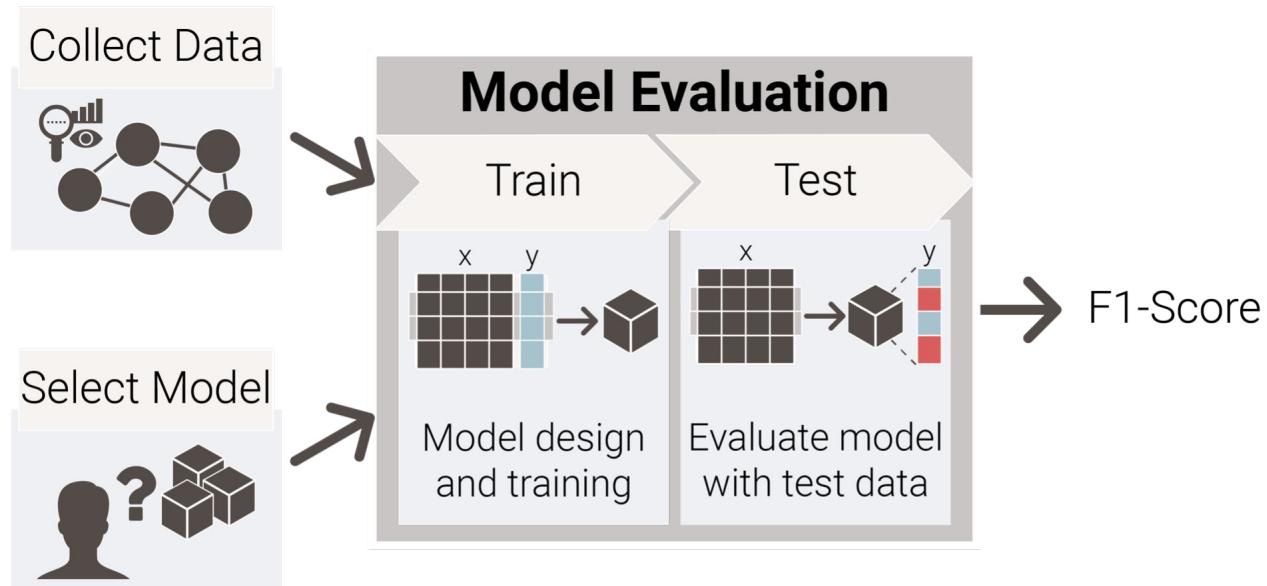
O.O.D. Samples

Model does not generalize!

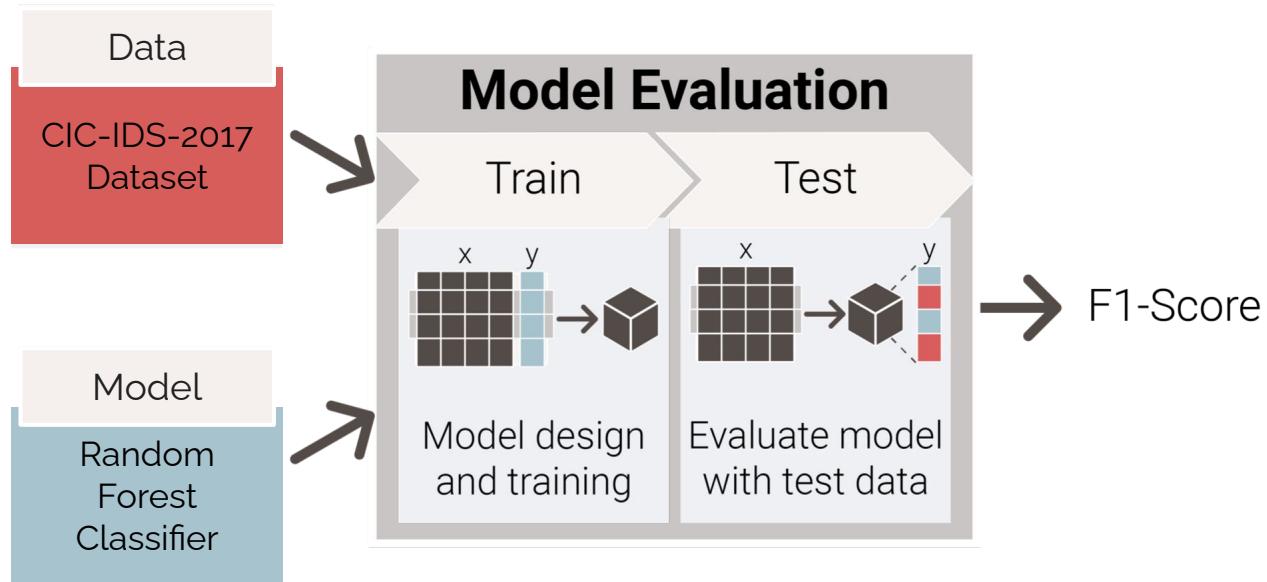
Spurious Correlations

Model picks up wrong correlations in the data!

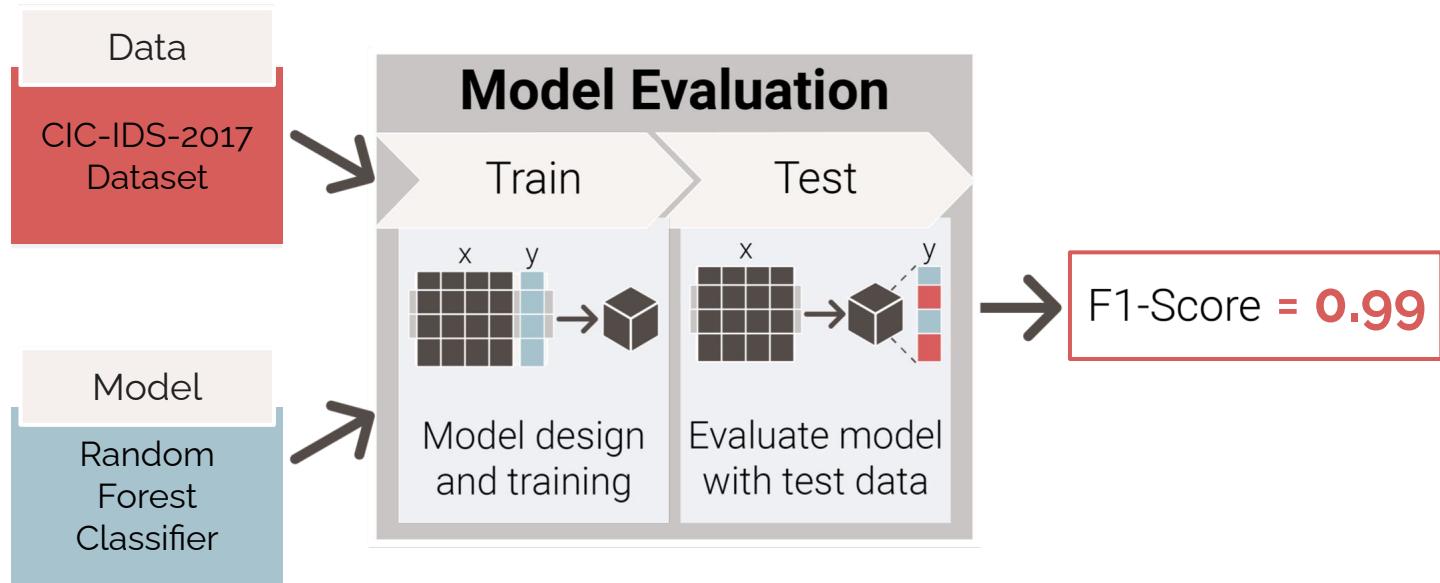
Consider this example...



Consider this example...



Consider this example...



Can you answer these questions?

Why (and how) does the model work?

When does the model not work?

Can you answer these questions?

Why (and how) does the model work?



When does the model not work?

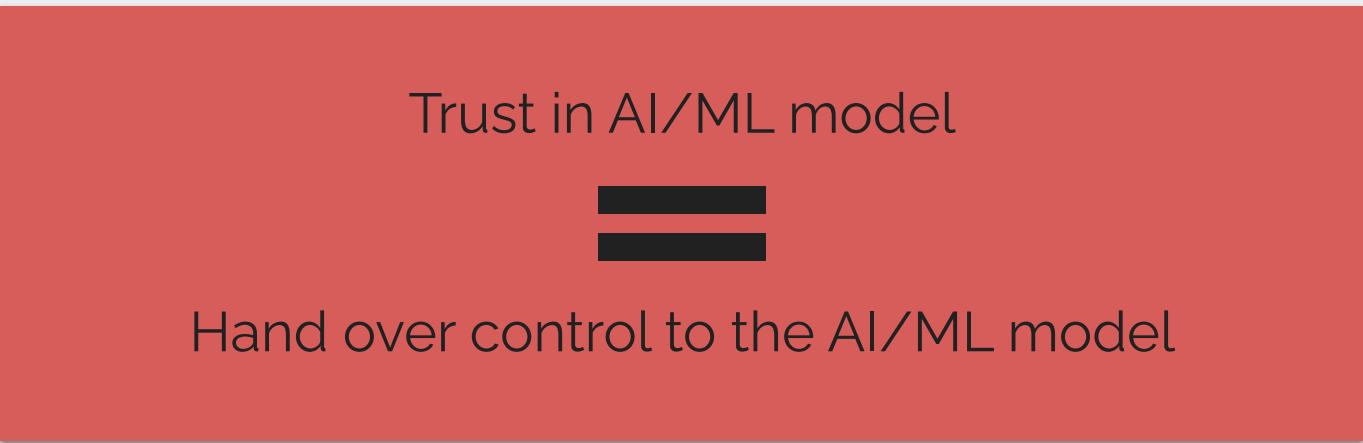




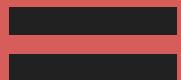
Can you **trust** this model?



Can you **trust** this model?



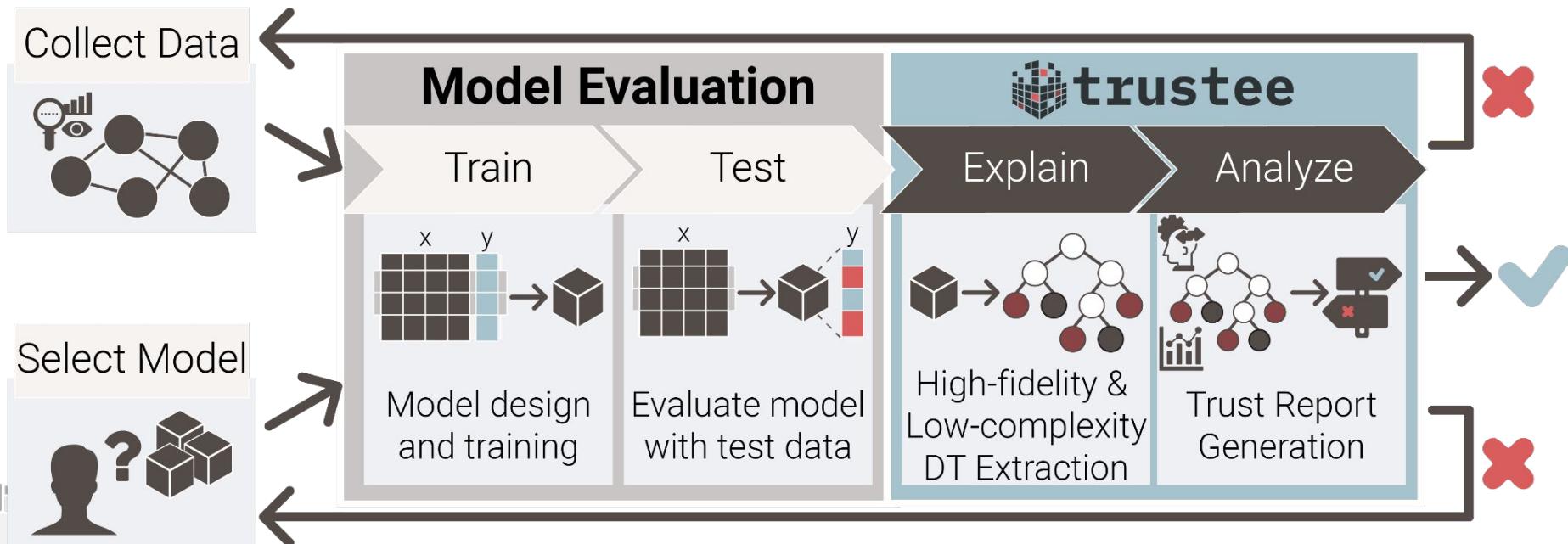
Trust in AI/ML model



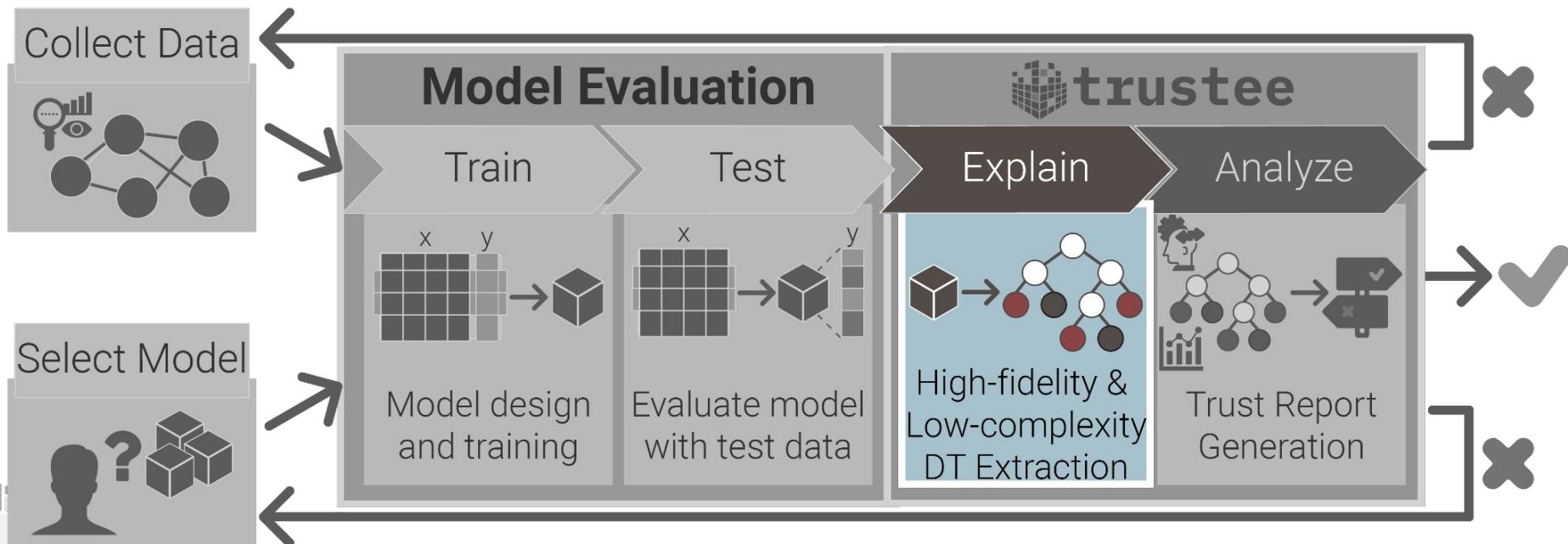
Hand over control to the AI/ML model



Augmented AI/ML Development Pipeline



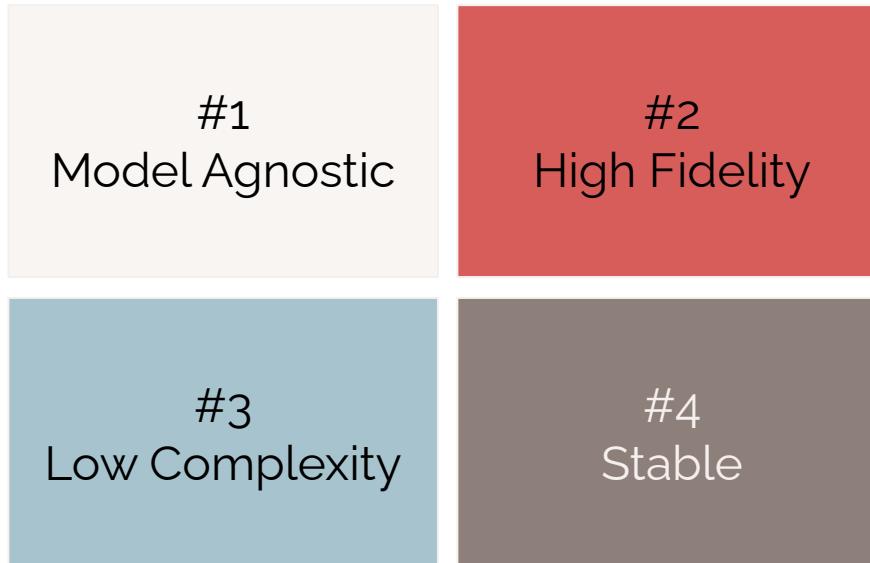
Augmented AI/ML Development Pipeline

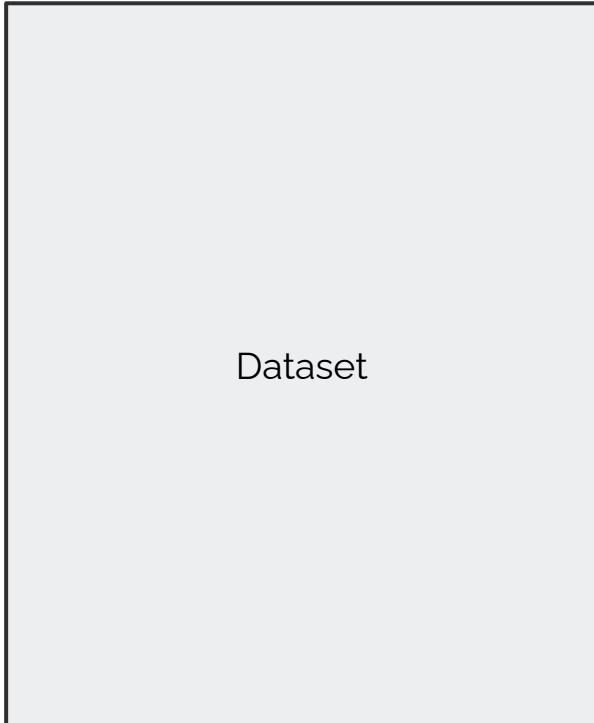




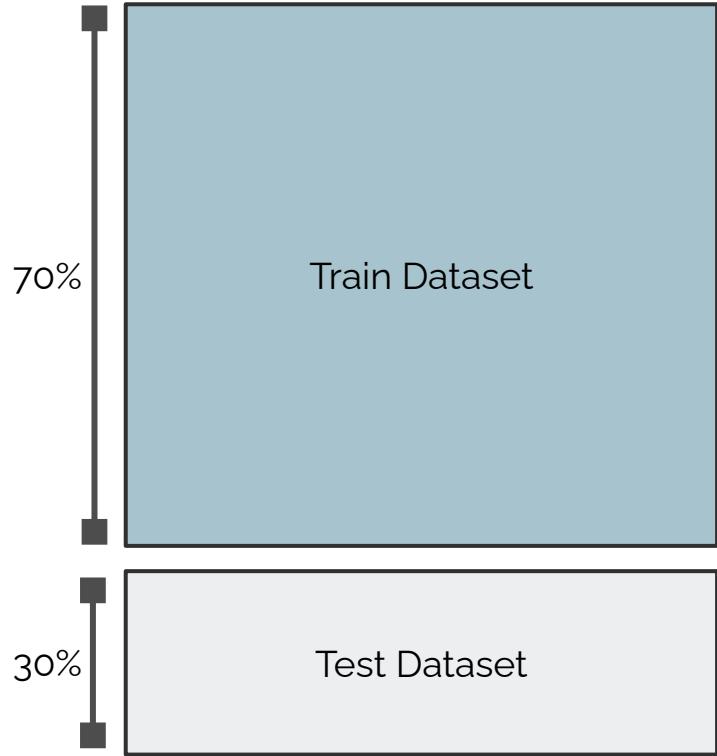
- Why do we need another XAI method?

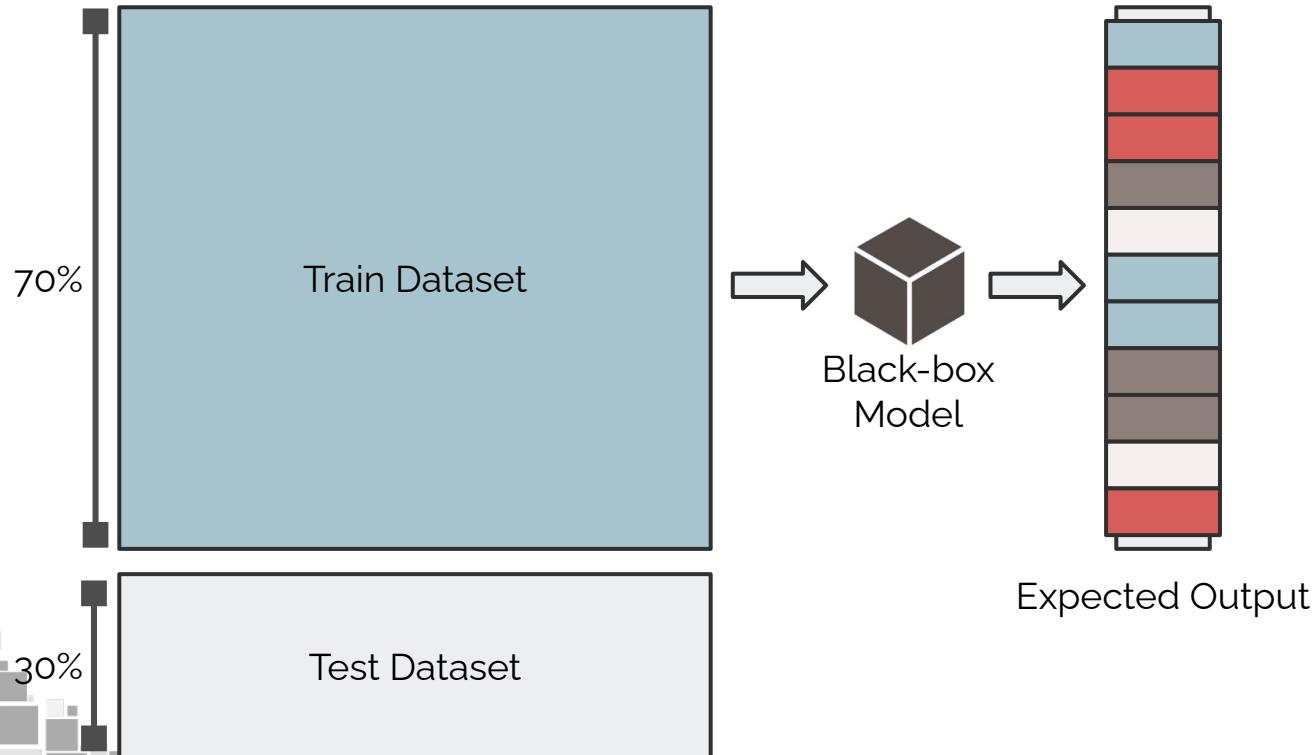
- Why do we need another XAI method?
- No existing global XAI method satisfies four requirements we desire

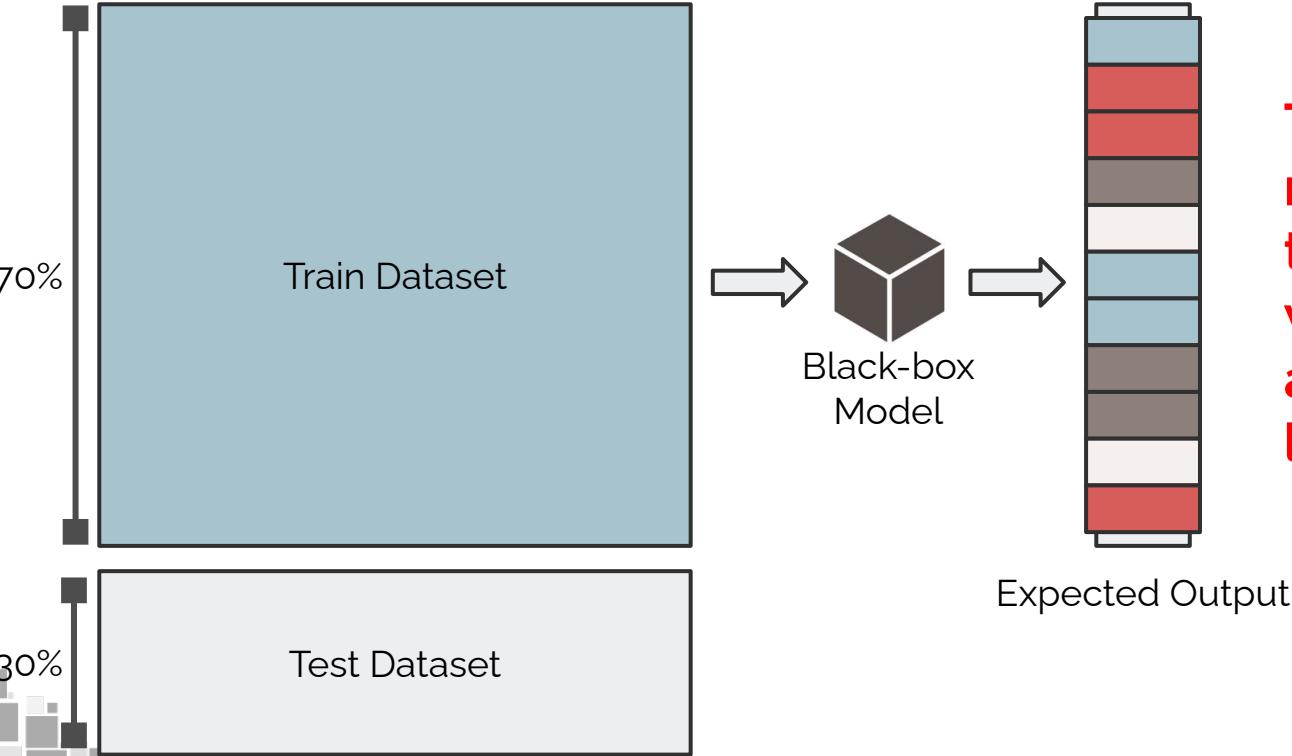




Black-box
Model

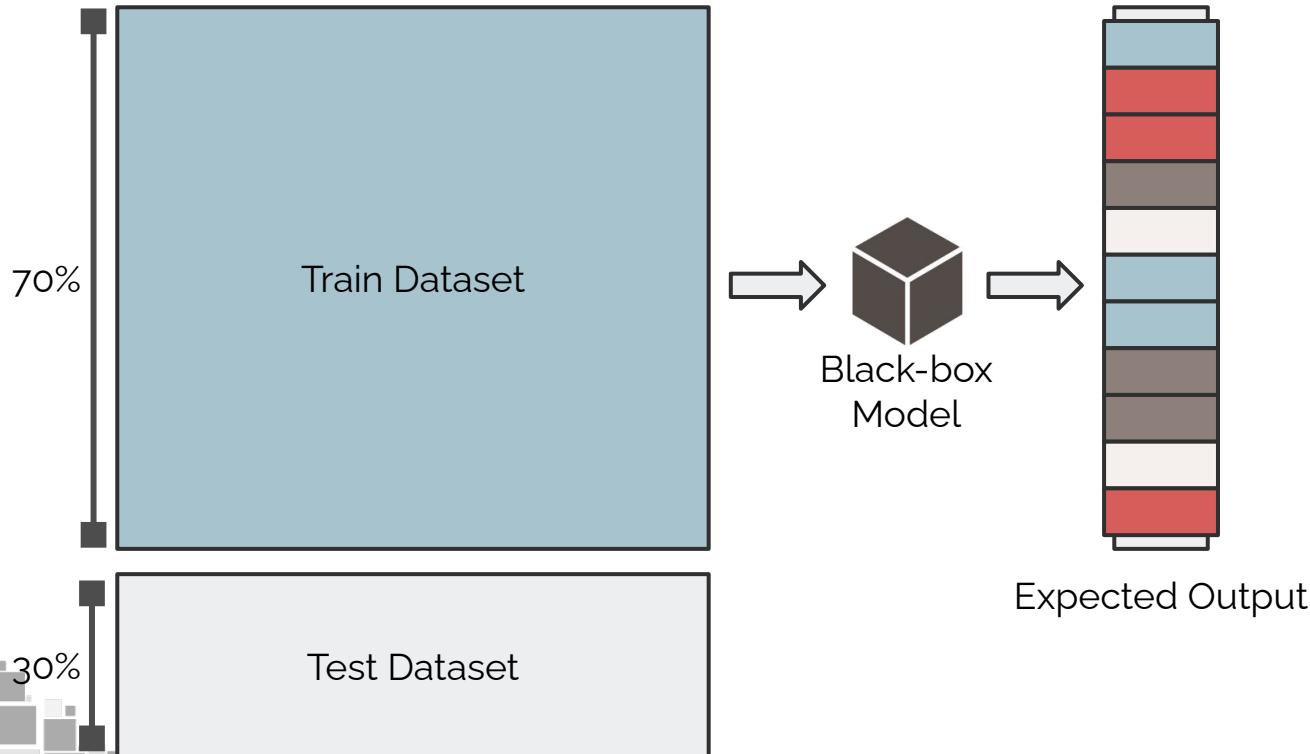


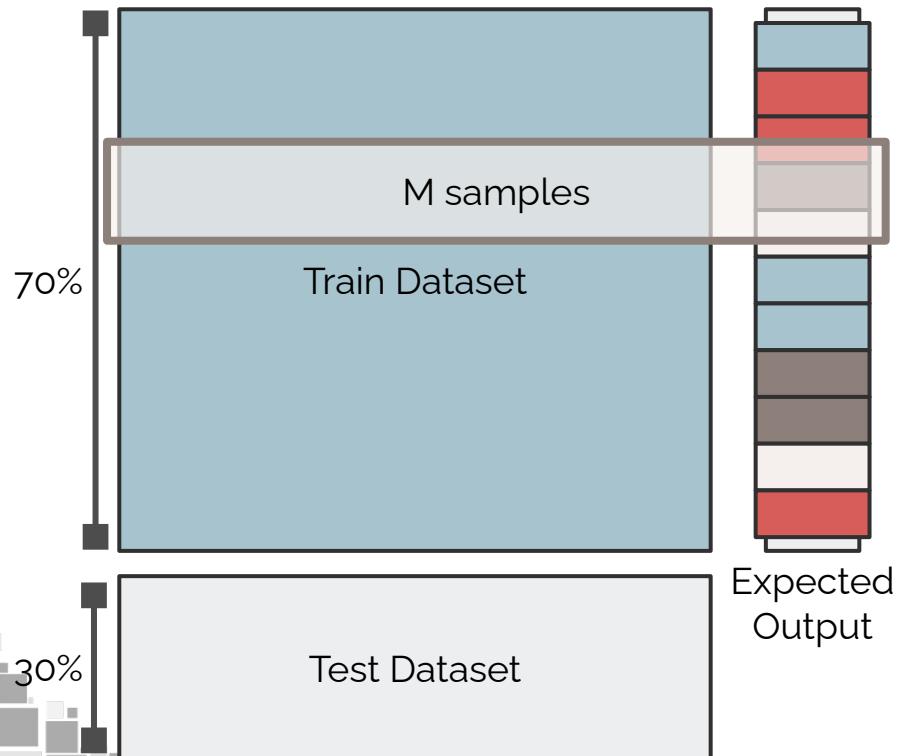


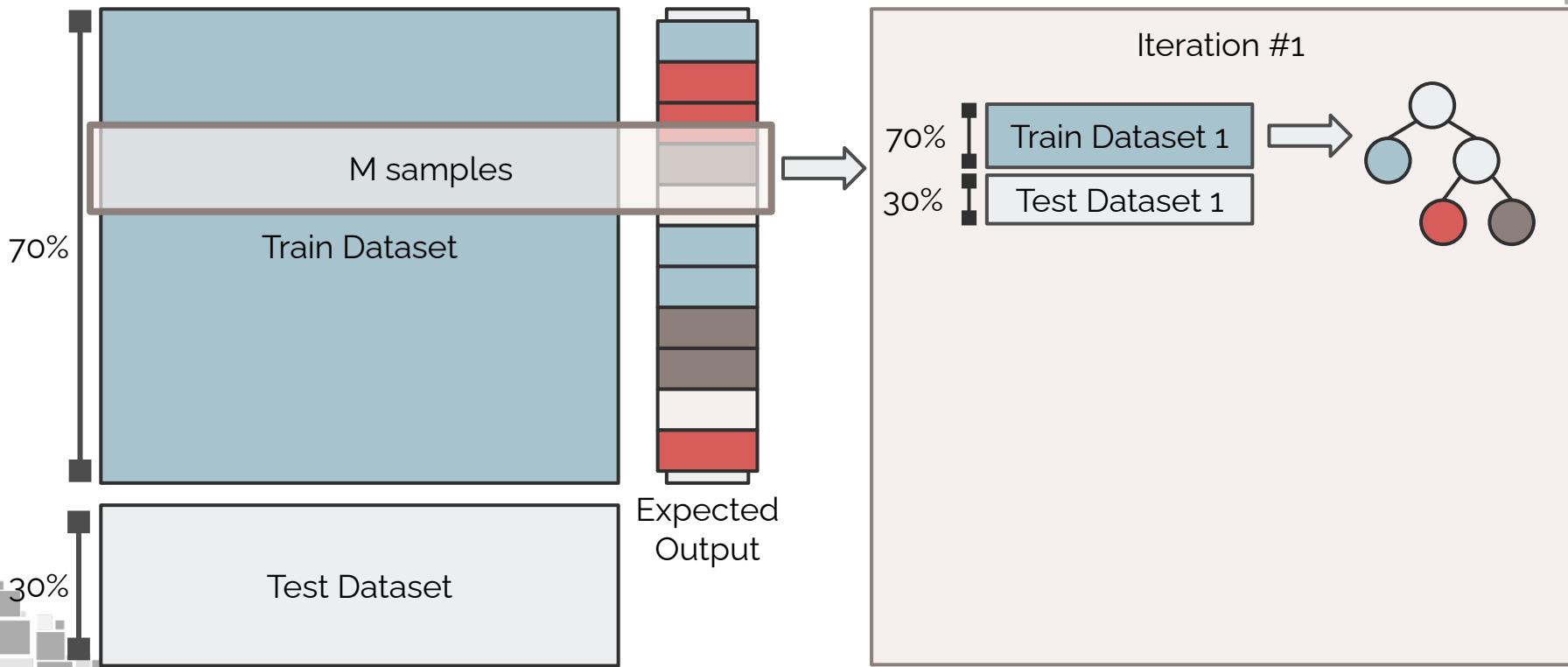


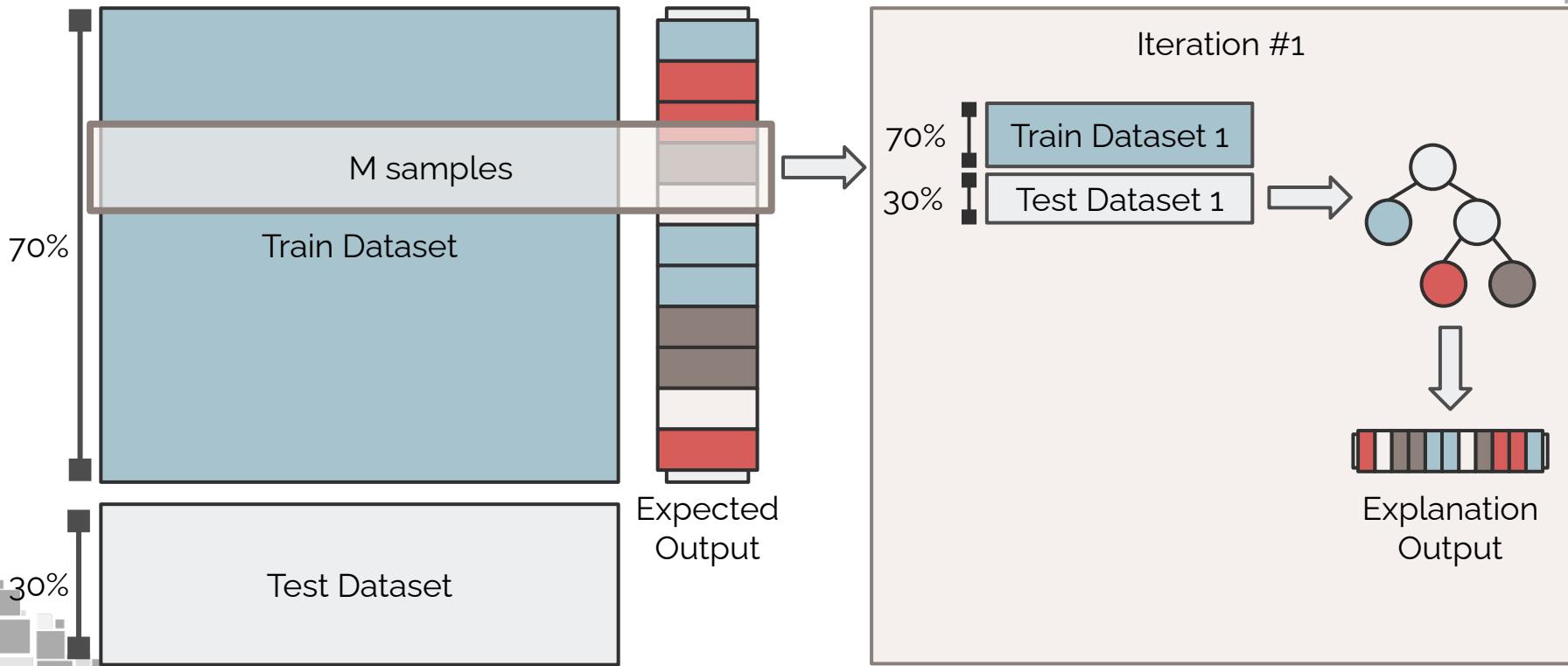
Trustee uses this model as an expert to train a student (a white-box model), akin to imitation learning.

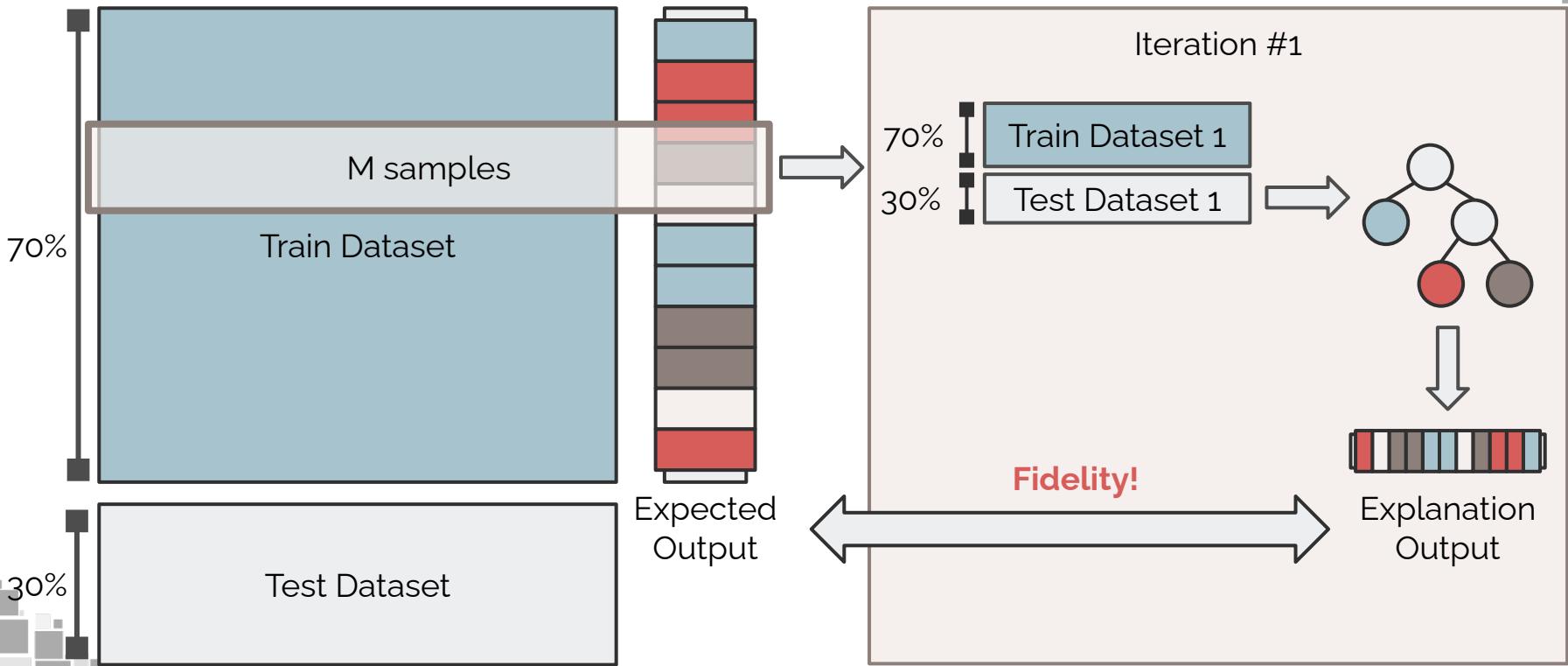
#1
Model
Agnostic

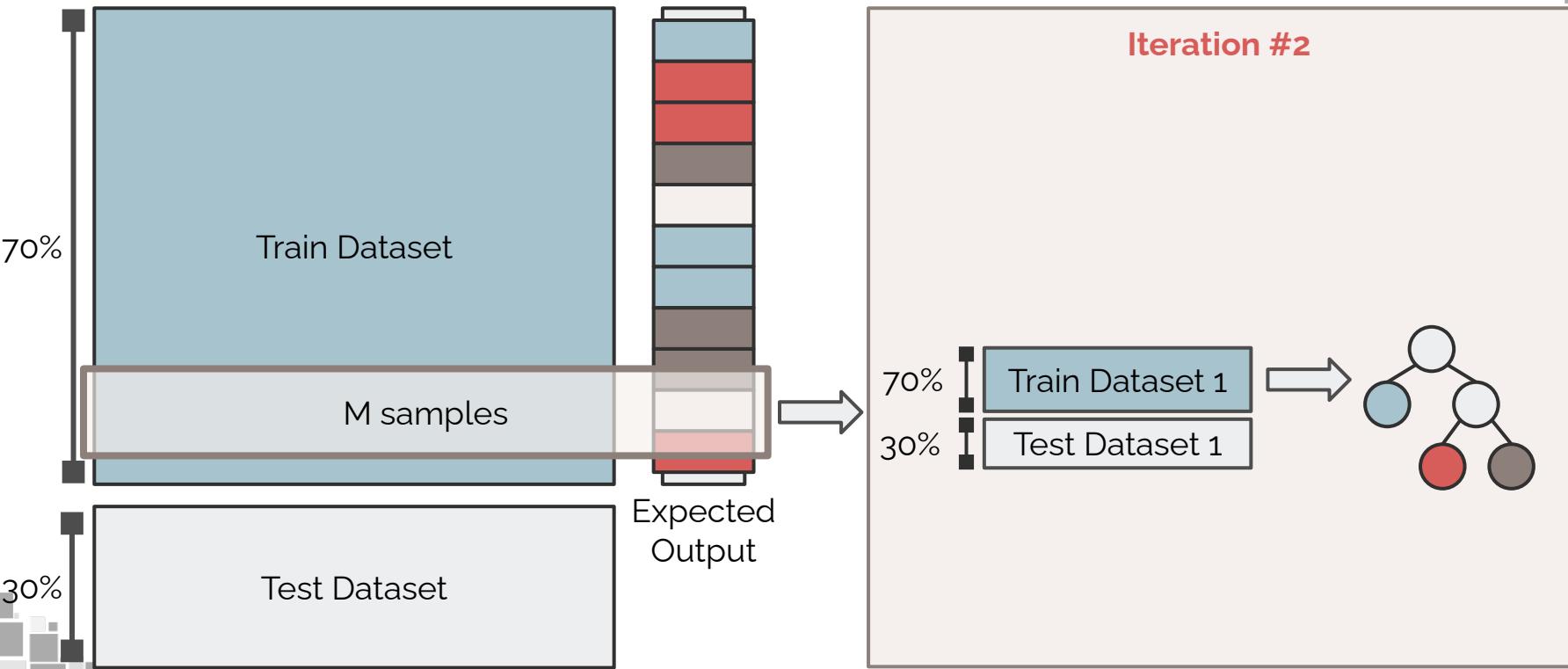


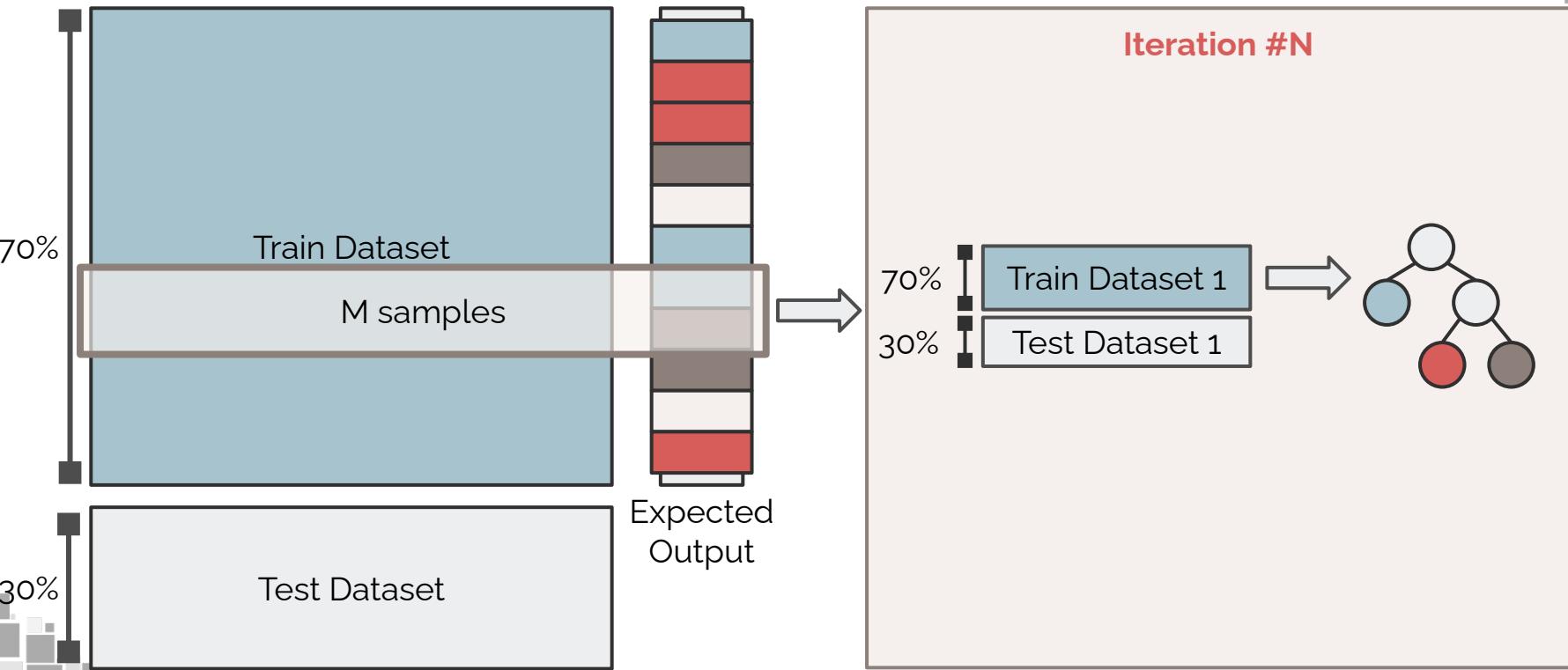


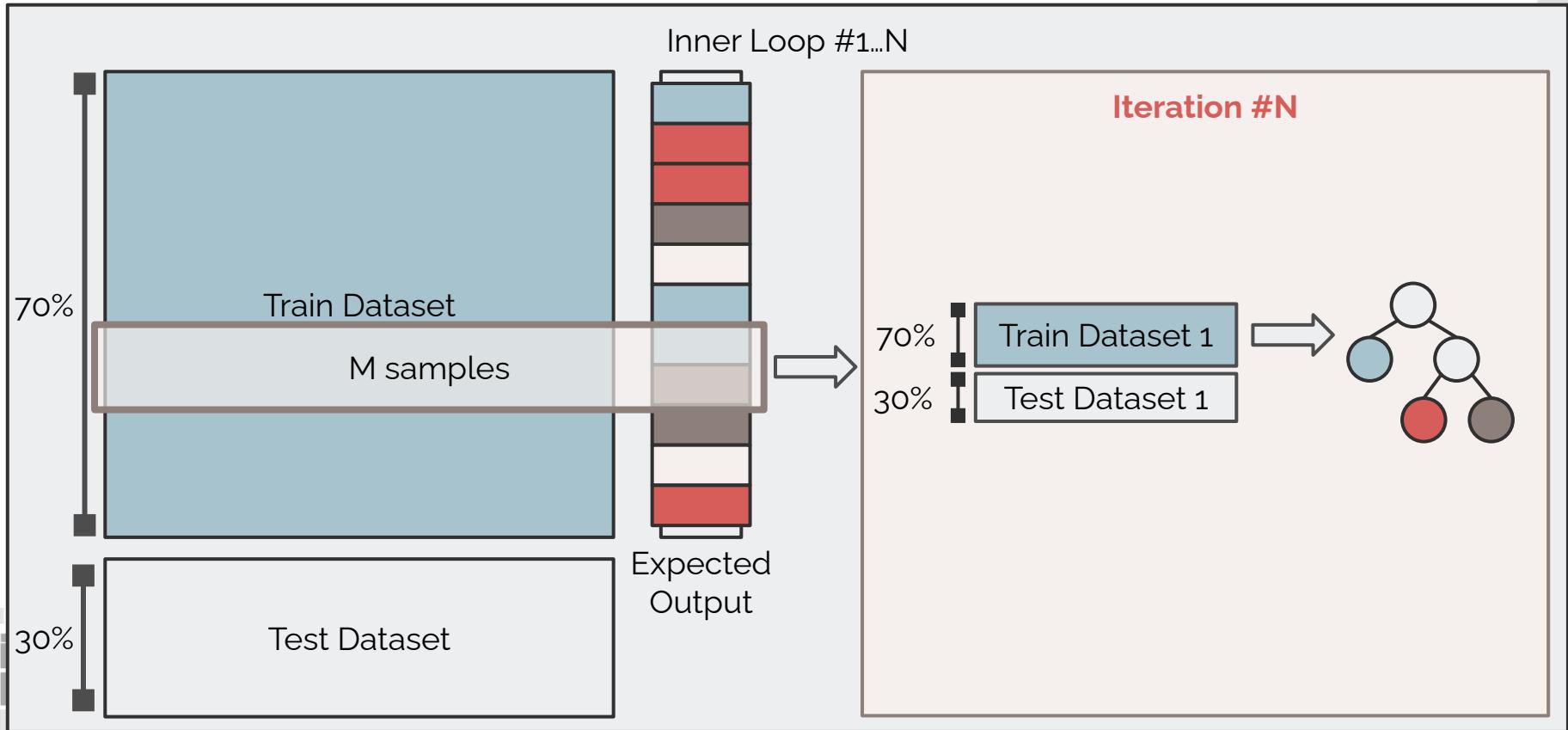


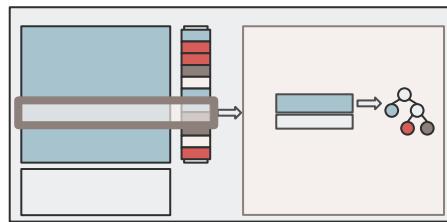


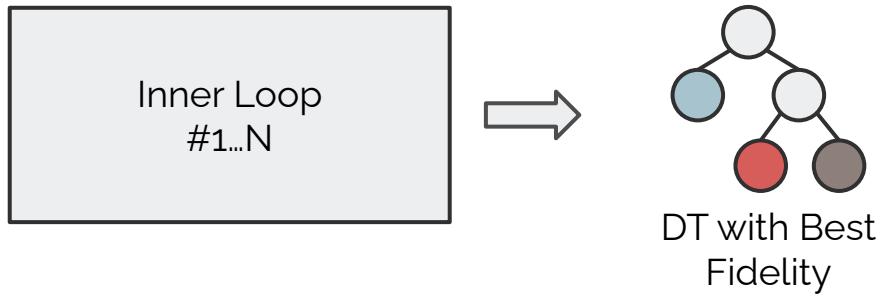




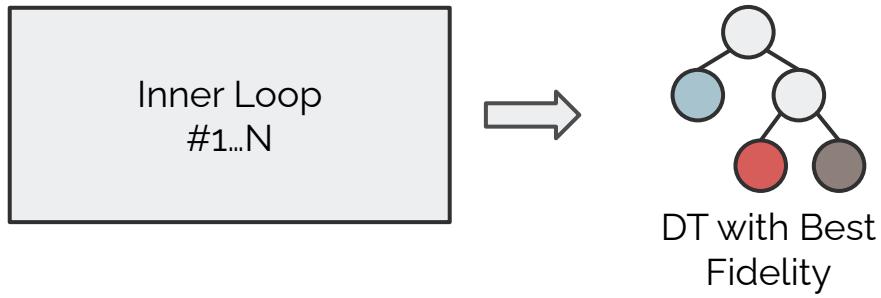








#2
High
Fidelity





How does Trustee achieve a high-fidelity explanation?

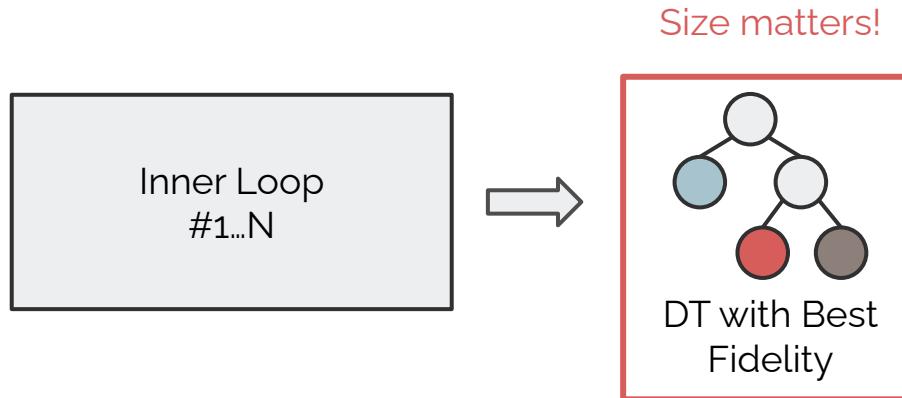


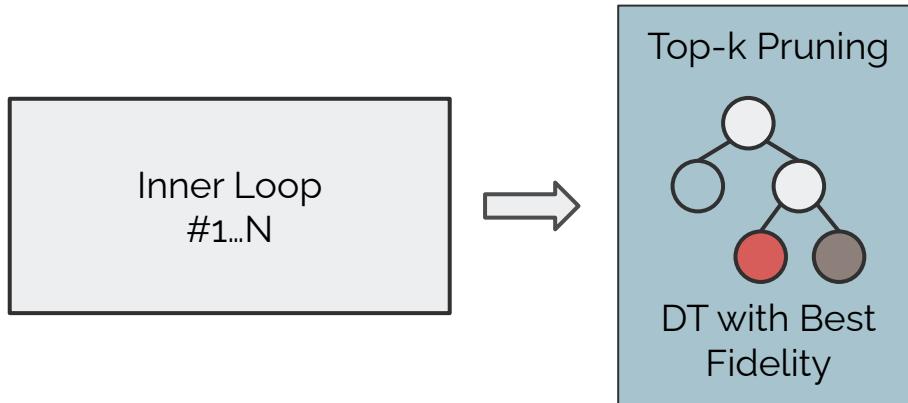
How does Trustee achieve a high-fidelity explanation?

- It optimizes for fidelity
 - The computed explanation is not necessarily the most accurate
- It does data augmentation
 - It overcorrects for the data samples for which the DT model makes wrong decisions



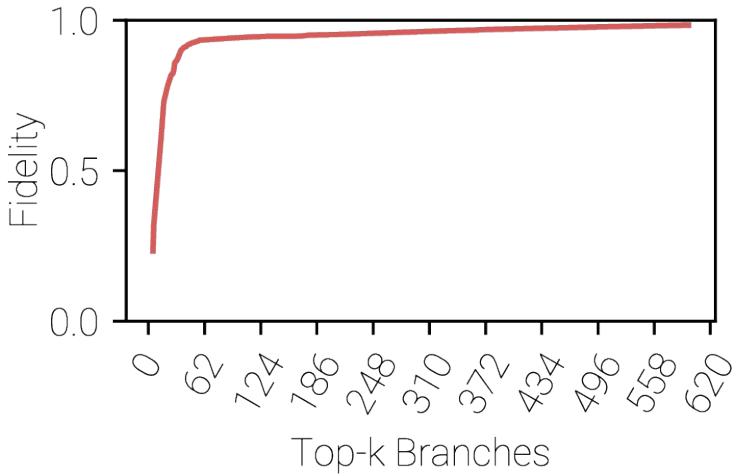
The generated DT may be too large!



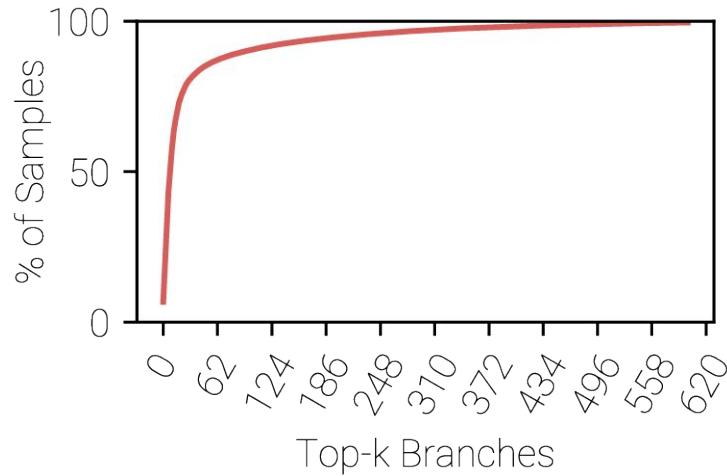


Top-k Pruning

Fidelity

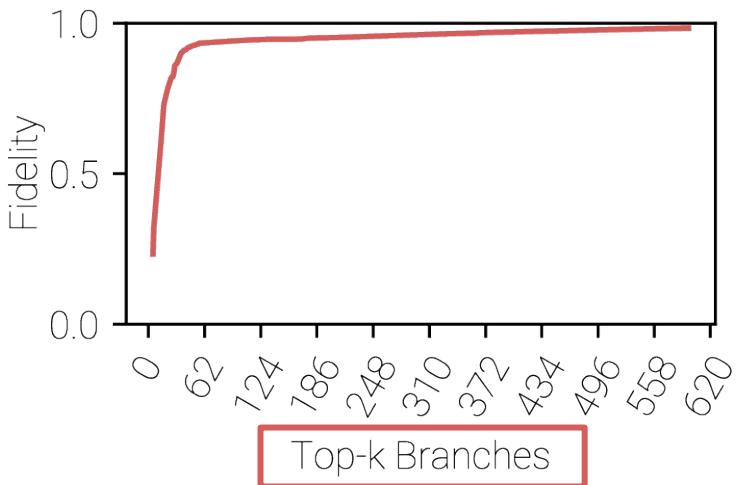


Samples

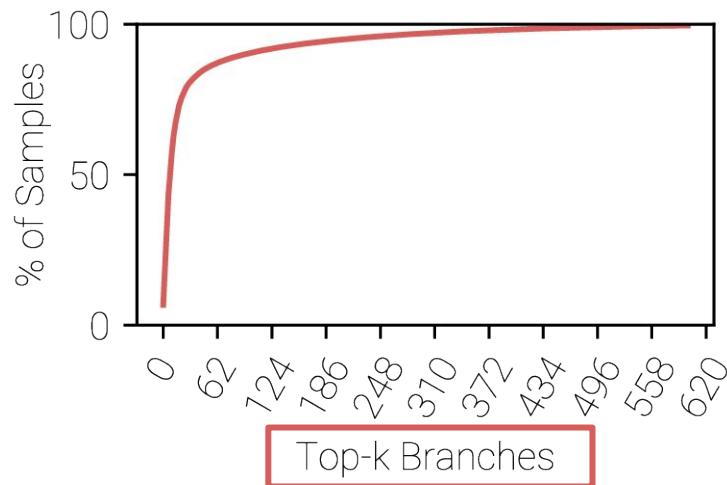


Top-k Pruning

Fidelity



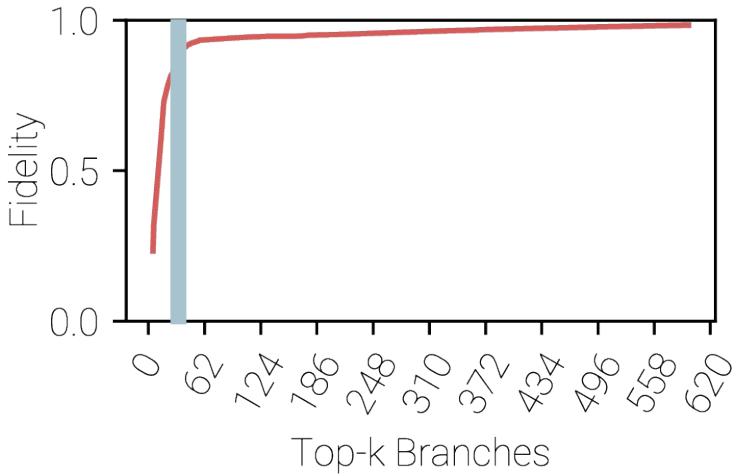
Samples



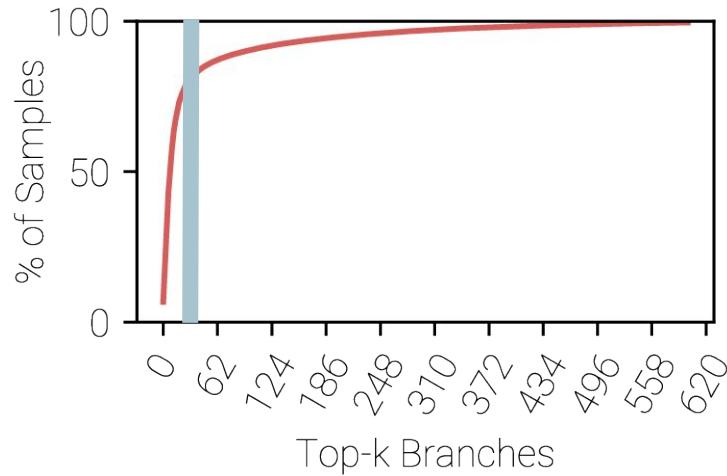
Diminishing returns!

Top-k Pruning

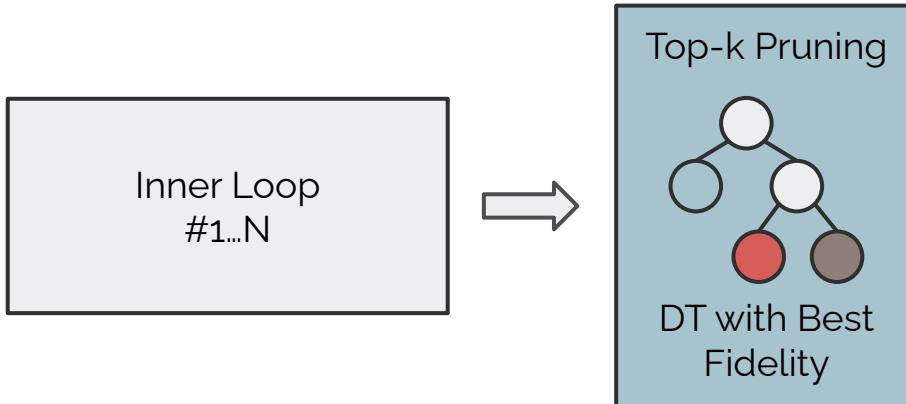
Fidelity



Samples



#3
Low
Complexity



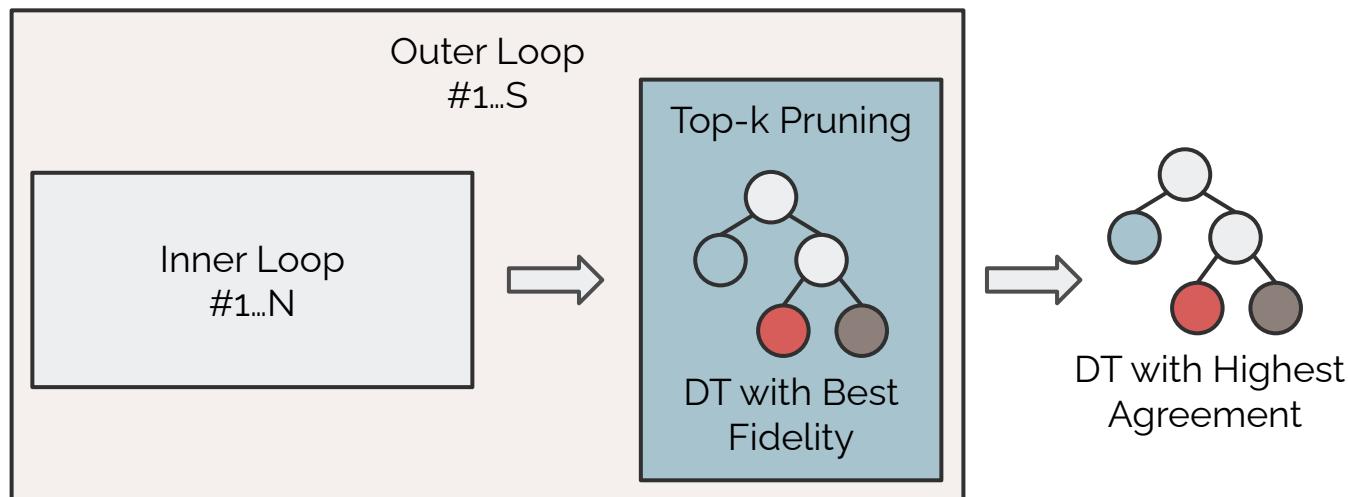


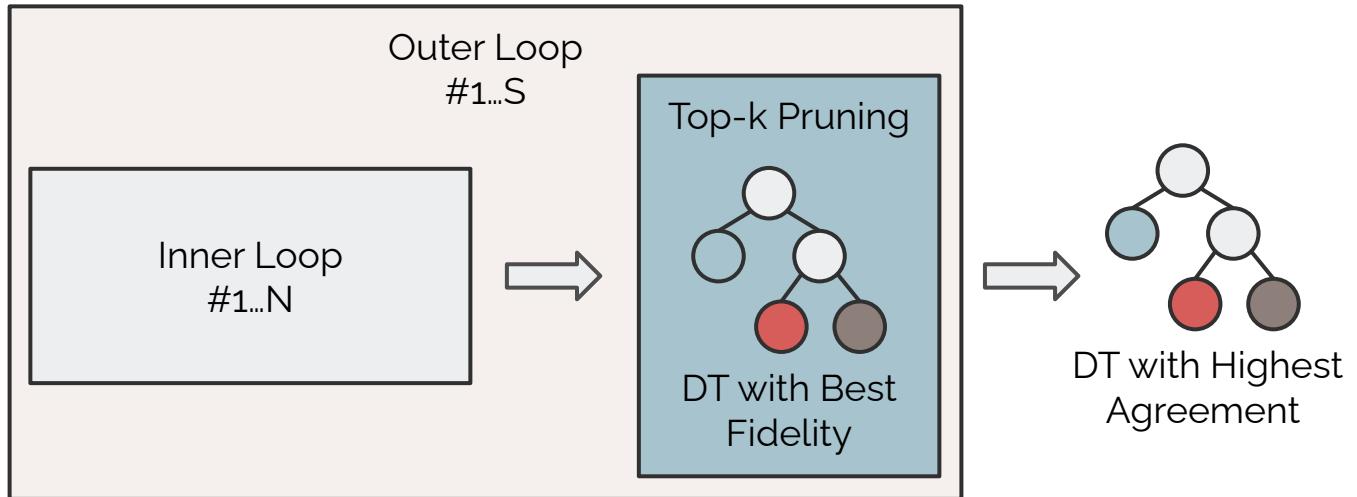
How does Trustee prevent obviously misleading explanations?



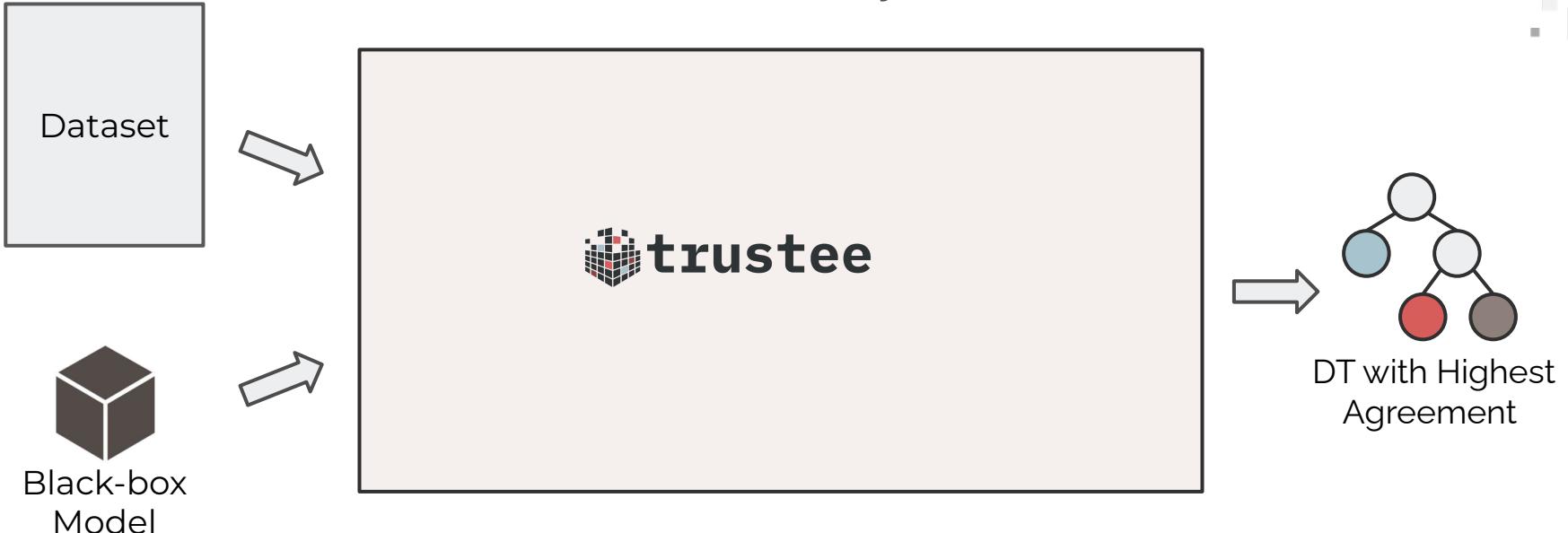
How does Trustee prevent obviously misleading explanations?

- It has an outer loop to select the DT with the highest mean agreement with the other explanations



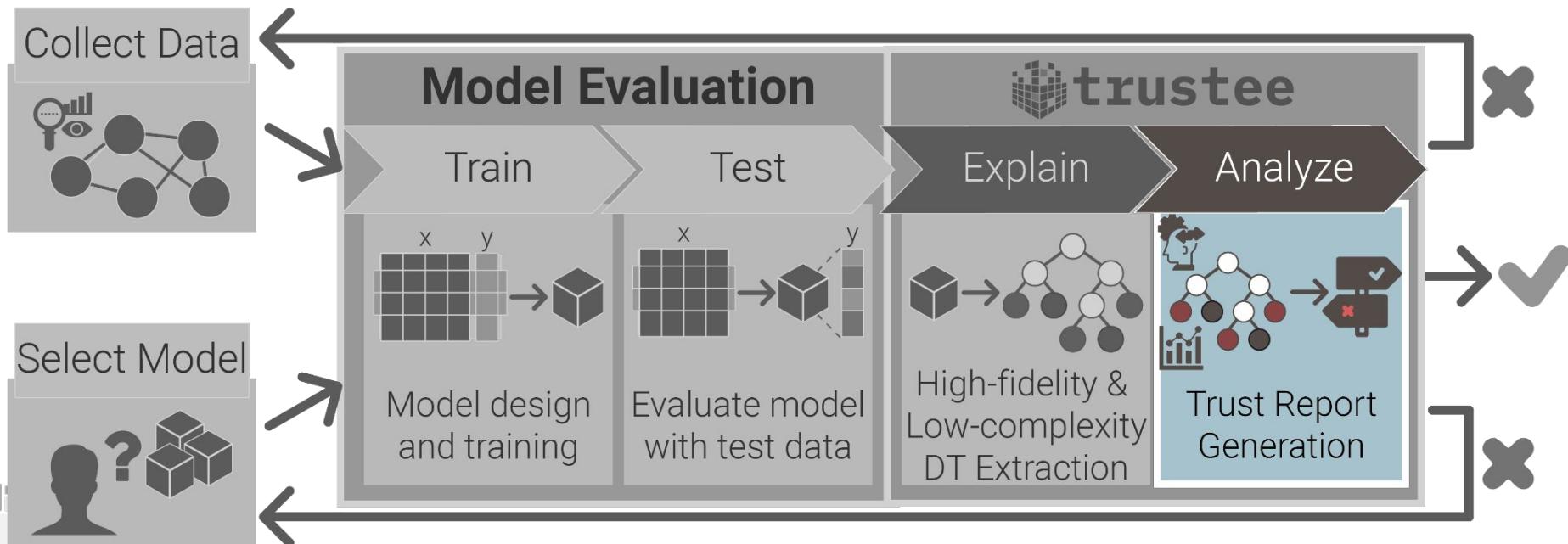


Summary



- Inner loop computes high-fidelity explanations
- Outer loop computes "the most stable explanation," i.e., the one with the highest mean agreement with the other explanations

Augmented AI/ML Development Pipeline



Generating Trust Reports

Underspecification issues! (revisited)

Shortcut Learning

Model takes shortcuts to classify data!

O.O.D. Samples

Model does not generalize!

Spurious Correlations

Model makes the picks up wrong correlations in the data!

Generating Trust Reports

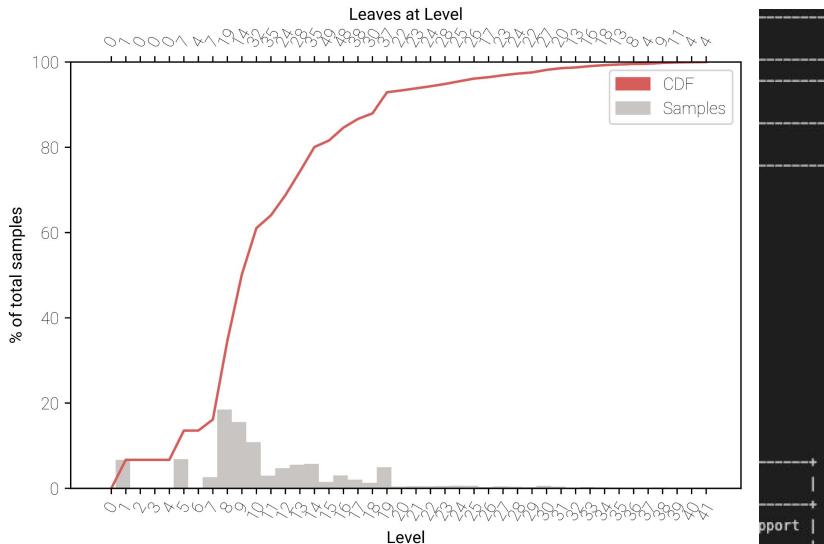
Classification Trust Report															
Summary															
Blackbox				Whitebox				Top-k Whitebox							
Model: RandomForestClassifier				Explanation method: Trustee				Explanation method: Trustee							
Dataset size:	947072			Model:	DecisionTreeClassifier			Model:	DecisionTreeClassifier			Iterations:	1		
Train/Test Split:	70.00% / 30.00%			Iterations:	1			Sample size:	50.00%			Iterations:	1		
# Input features:	61			Decision Tree Info				Decision Tree Info							
# Output classes:	5			Size:	2437			Size:	9			Depth:	31		
				Depth:	1219			Leaves:	5			Leaves:	1		
				# Input features:	18 (29.51%)			# Input features:	-			# Output classes:	5 (100.00%)		
				# Output classes:	5 (100.00%)										
Performance				Fidelity				Fidelity							
	precision	recall	f1-score	support		precision	recall	f1-score	support		precision	recall	f1-score	support	
0	1.000	0.912	0.954	24408		0	1.000	1.000	1.000	22254		0	0.000	0.000	0.000
1	0.752	0.910	0.824	1872		1	1.000	1.000	1.000	2265		1	0.000	0.000	0.000
2	0.929	0.827	0.875	10994		2	0.969	0.965	0.967	9781		2	0.000	0.000	0.000
3	0.997	0.929	0.962	65188		3	0.998	0.998	0.998	60768		3	0.544	0.957	0.694
4	0.958	0.997	0.978	181660		4	0.998	0.998	0.998	189054		4	0.875	0.821	0.847
accuracy			0.967	284122		accuracy			0.997	284122		accuracy			0.751
macro avg	0.927	0.915	0.918	284122		macro avg	0.993	0.992	0.993	284122		macro avg	0.284	0.356	0.308
weighted avg	0.968	0.967	0.967	284122		weighted avg	0.997	0.997	0.997	284122		weighted avg	0.699	0.751	0.712

Generating Trust Reports

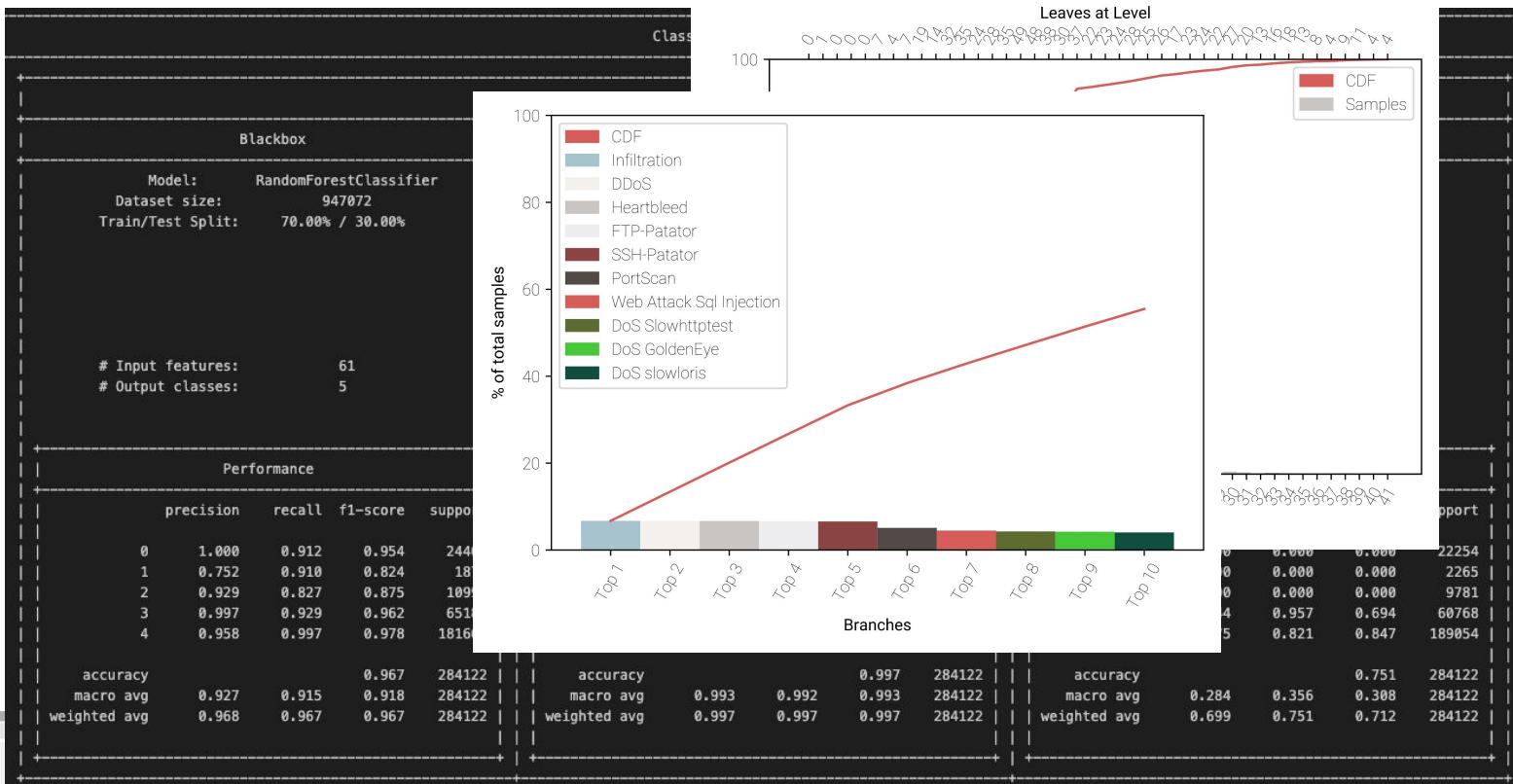
Classification Trust Report											
Summary											
Blackbox				Whitebox				Top-k Whitebox			
Model:	RandomForestClassifier	Explanation method:	Trustee	Model:	DecisionTreeClassifier	Explanation method:	Trustee	Model:	DecisionTreeClassifier	Explanation method:	Trustee
Dataset size:	947072			Model:	DecisionTreeClassifier			Model:	DecisionTreeClassifier		
Train/Test Split:	70.00% / 30.00%	Iterations:	1	Iterations:	1	Iterations:	1	Iterations:	1	Iterations:	1
		Sample size:	50.00%	Sample size:	50.00%	Sample size:	50.00%	Sample size:	50.00%	Sample size:	50.00%
Decision Tree Info											
# Input features: 61				Size: 2437 Depth: 31 Leaves: 1219 # Input features: 18 (29.51%) # Output classes: 5 (100.00%)				Size: 9 Depth: 4 Leaves: 5 Top-k: 1 # Input features: - # Output classes: 5 (100.00%)			
Performance				Fidelity				Fidelity			
	precision	recall	f1-score	support		precision	recall	f1-score	support		precision
0	1.000	0.912	0.954	24408	0	1.000	1.000	1.000	22254	0	0.000
1	0.752	0.910	0.824	1872	1	1.000	1.000	1.000	2265	1	0.000
2	0.929	0.827	0.875	10994	2	0.969	0.965	0.967	9781	2	0.000
3	0.997	0.929	0.962	65188	3	0.998	0.998	0.998	60768	3	0.544
4	0.958	0.997	0.978	181660		0.998	0.998	0.998	105054	4	0.875
accuracy					accuracy					accuracy	
macro avg	0.927	0.915	0.918	284122	macro avg	0.993	0.992	0.993	284122	macro avg	0.284
weighted avg	0.968	0.967	0.967	284122	weighted avg	0.997	0.997	0.997	284122	weighted avg	0.699

Generating Trust Reports

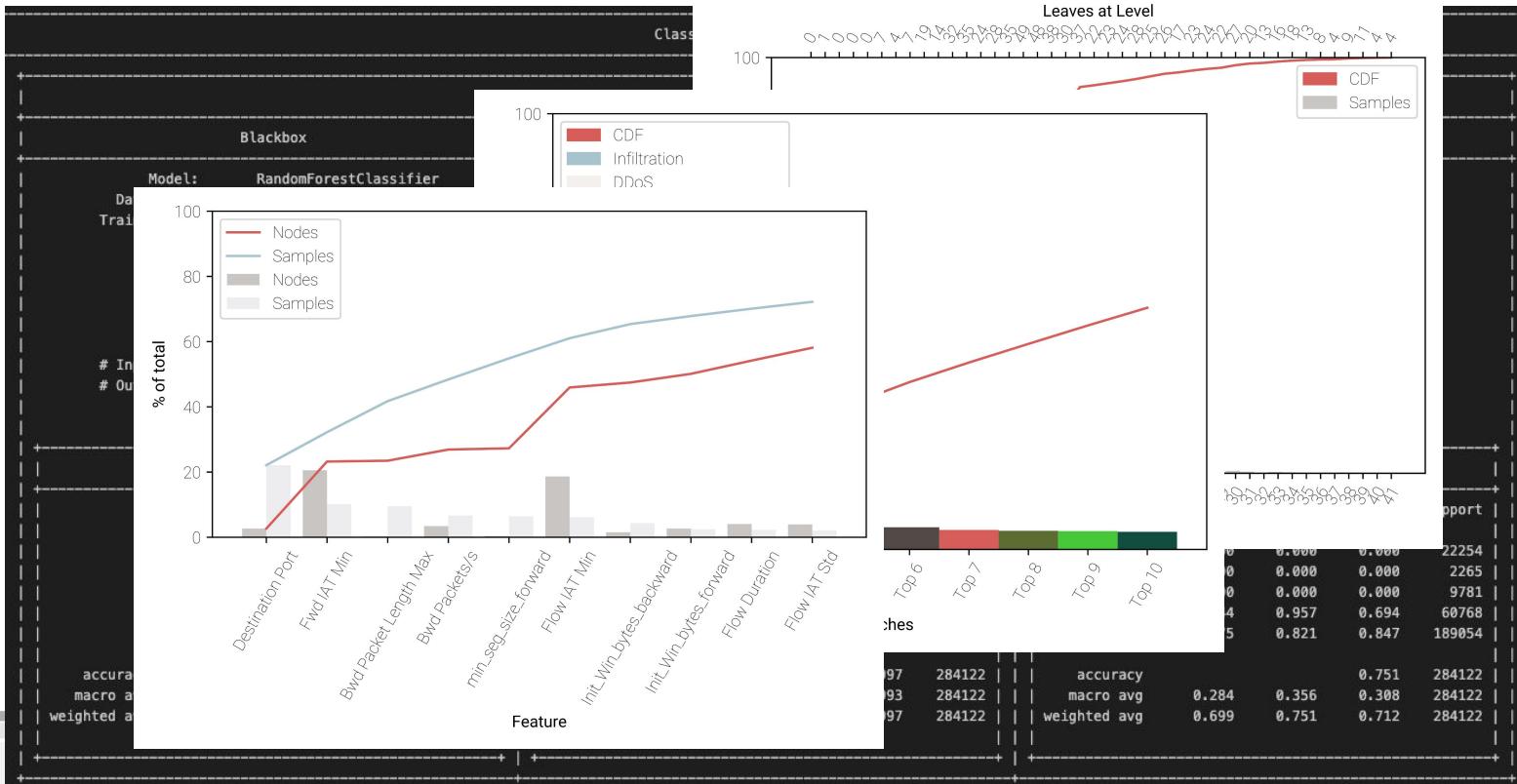
```
Class
+-----+
| Blackbox |
+-----+
Model: RandomForestClassifier
Dataset size: 947072
Train/Test Split: 70.00% / 30.00%
# Input features: 61
# Output classes: 5
+-----+
Explanation m
Model:
Iterations:
Sample si
Decision Tre
Size:
Depth:
Leaves:
# Input fea
# Output cla
+-----+
Performance
precision recall f1-score support precision recall f1-score support
0 1.000 0.912 0.954 24408 0 1.000 1.000 1.000 22254
1 0.752 0.910 0.824 1872 1 1.000 1.000 1.000 2265
2 0.929 0.827 0.875 10994 2 0.969 0.965 0.967 9781
3 0.997 0.929 0.962 65188 3 0.998 0.998 0.998 60768
4 0.958 0.997 0.978 181660 4 0.998 0.998 0.998 189054
accuracy 0.967 284122 accuracy 0.997 284122 accuracy 0.751 284122
macro avg 0.927 0.915 0.918 284122 macro avg 0.993 0.992 0.993 284122 macro avg 0.284 0.356 0.308 284122
weighted avg 0.968 0.967 0.967 284122 weighted avg 0.997 0.997 0.997 284122 weighted avg 0.699 0.751 0.712 284122
+-----+
```



Generating Trust Reports



Generating Trust Reports



Use Cases

Use Case #1: Detecting VPN vs Non-VPN Traffic

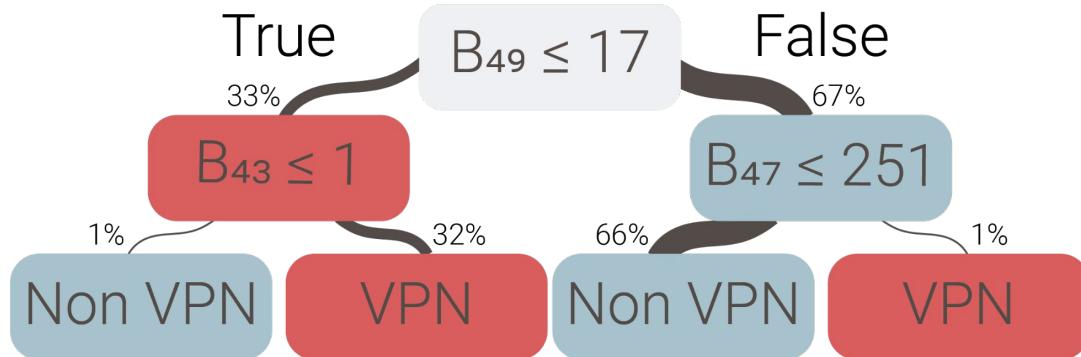
Problem Setup

- **Selected publication:**
 - “*End-to-end encrypted traffic classification with one-dimensional convolution neural networks*” — Wang et al., 2017
- **Proposal:**
 - **Model:** 1D-CNN to classify traffic between encrypted VPN traffic and non-encrypted traffic (i.e. VPN vs Non-VPN)
 - **Features:** first 784 raw bytes of each PCAP file
 - **Dataset:** ISCX VPN-nonVPN 2016 [<https://www.unb.ca/cic/datasets/vpn.html>]
- **Results:**
 - Reported F1-score: 0.99
 - Reproduced F1-score: 0.959

Use Case #1: Detecting VPN vs Non-VPN Traffic

Explanation

Fidelity: 1.000
No pruning
7 nodes



Use Case #1: Detecting VPN vs Non-VPN Traffic

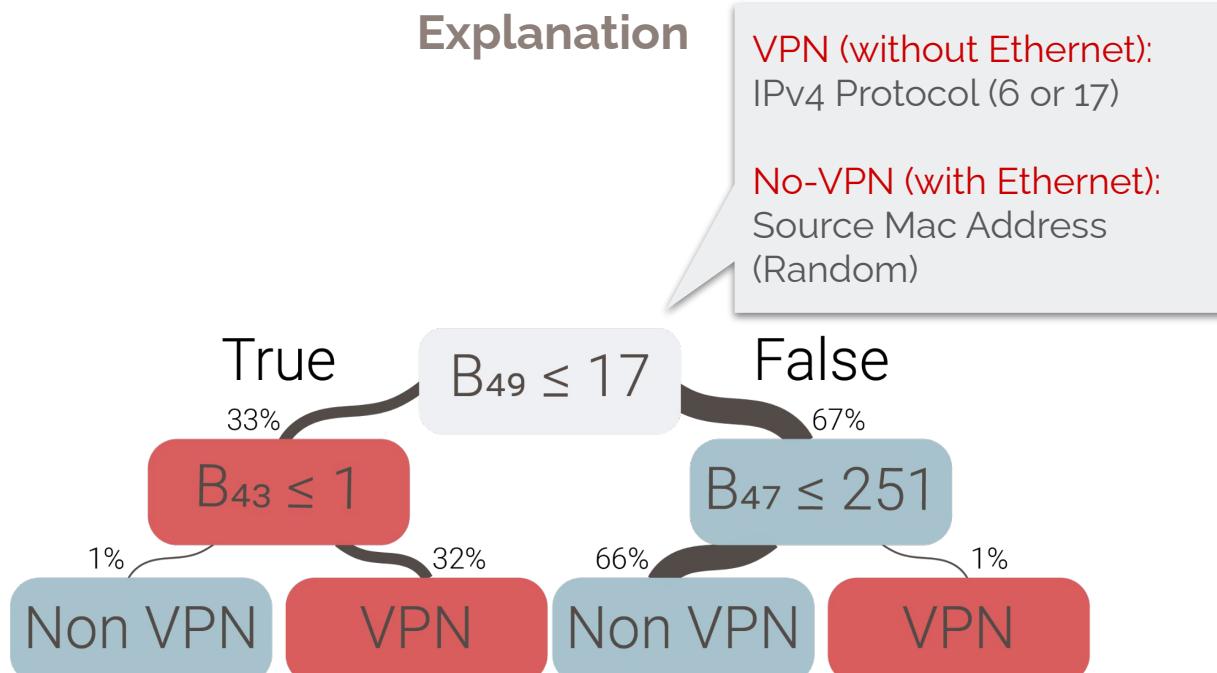
Non VPN

	0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19
Pcap	161	178	195	212	0	2	0	4	0	0	0	0	0	0	0	0	0	255	255	
Meta	0	0	0	1	85	65	10	69	0	5	80	24	0	0	0	64	0	0	0	64
Eth	1	0	94	0	0	252	184	172	111	54	28	162	8	0	69	0	0	50	65	228
IPv4	0	0	1	17	34	185	131	202	240	87	224	0	0	252	201	86	20	235	0	...

VPN

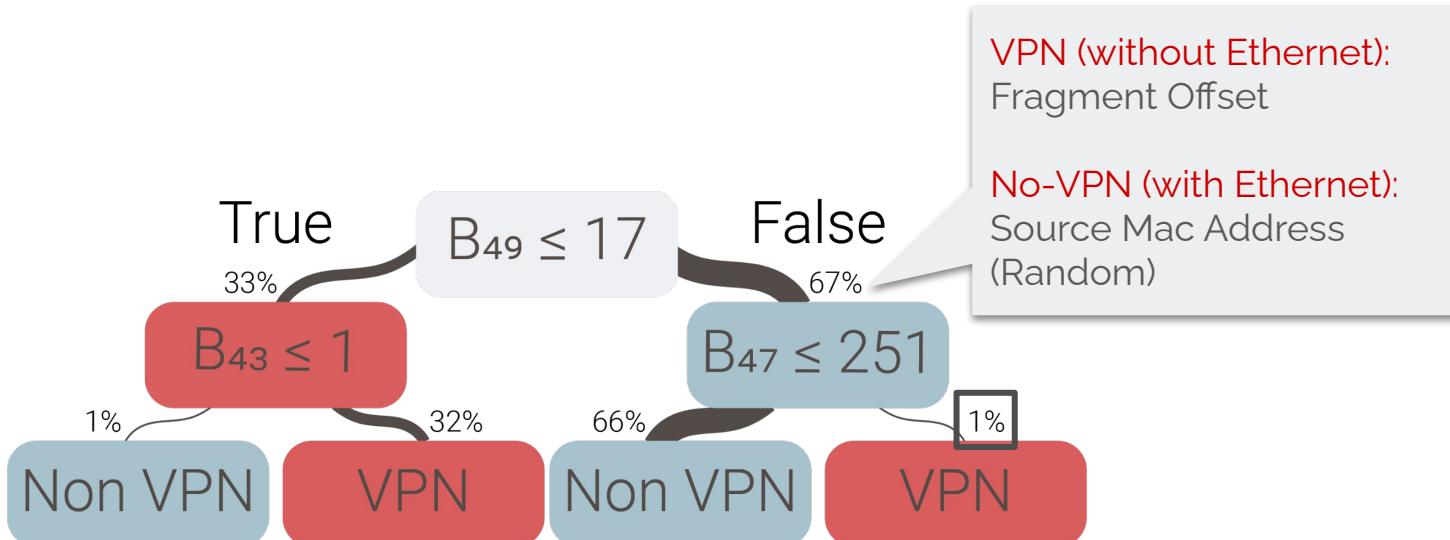
	0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19
Pcap	161	178	195	212	0	2	0	4	0	0	0	0	0	0	0	0	0	0	255	255
Meta	0	0	0	101	85	45	101	91	0	0	111	11	0	0	0	56	0	0	0	56
IPv4	69	0	0	56	99	213	64	0	14	17	5	254	10	8	0	10	69	171	255	36
UDP	146	214	13	150	0	36	120	43	0	1	0	8	33	18	164	66	52	167	9	...

Use Case #1: Detecting VPN vs Non-VPN Traffic



Use Case #1: Detecting VPN vs Non-VPN Traffic

Explanation

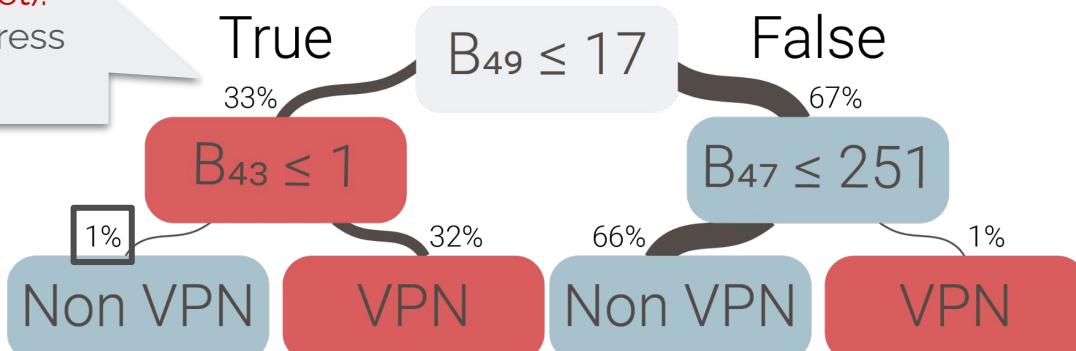


Use Case #1: Detecting VPN vs Non-VPN Traffic

Explanation

VPN (without Ethernet):
IP Total Length

No-VPN (with Ethernet):
Destination Mac Address
(Always 0)



Use Case #1: Detecting VPN vs Non-VPN Traffic

Validation

- Validation dataset:
 - Tampering with packet headers from original PCAPs

Validation Dataset	Avg. Precision	Avg. Recall	Avg. F1
<i>Untampered</i>	0.959	0.956	0.955
<i>Tampered-43-47-49</i>	0.959	0.956	0.955

Use Case #1: Detecting VPN vs Non-VPN Traffic

No VPN

		Byte 23: PCAP Link Type																		
		No-VPN (With Ethernet): 1																		
Pcap Meta		161	178	195	2	0	2	0	4	0	0	0	0	0	0	0	0	0	255	255
		0	0	0	1	85	65	10	69	0	5	80	24	0	0	0	64	0	0	64
Ethernet		Destination MAC Address		Source MAC Address																
		1	0	94	0	0	252	184	172	111	54	28	162	8	0	69	0	0	50	65
IPv4		0	0	1	17	34	185	131	202	240	87	224	0	0	252	201	86	20	235	0
		Byte 23: PCAP Link Type																		

VPN

		Byte 23: PCAP Link Type																		
		VPN (Without Ethernet): 101																		
Pcap Meta		161	178	195	2	0	2	0	4	0	0	0	0	0	0	0	0	0	255	255
		0	0	0	101	85	45	101	91	0	0	111	11	0	0	0	56	0	0	56
IPv4		Total Length		Frag. Off.		Protocol														
		69	0	0	56	199	213	64	0	64	17	35	254	10	8	0	10	69	171	255
UDP		146	214	13	150	0	36	120	43	0	1	0	8	33	18	164	66	52	167	9
		Byte 23: PCAP Link Type																		

Use Case #1: Detecting VPN vs Non-VPN Traffic

Validation

- Validation dataset:
 - Tampering with packet headers from original PCAPs

Validation Dataset	Avg. Precision	Avg. Recall	Avg. F1
<i>Untampered</i>	0.959	0.956	0.955
<i>Tampered-43-47-49</i>	0.959	0.956	0.955
<i>Tampered-32-to-63</i>	0.889	0.867	0.856
<i>Tampered-0-to-63</i>	0.831	0.757	0.734
<i>Tampered-0-to-127</i>	0.753	0.555	0.398

Use Case #1: Detecting VPN vs Non-VPN Traffic

Validation

- Validation dataset:
 - Tampering with packet headers from original PCAPs

Validation Dataset	Avg. Precision	Avg. Recall	Avg. F1
<i>Untampered</i>	0.959	0.956	0.955
<i>Tampered-43-47-49</i>	0.959	0.956	0.955
<i>Tampered-32-to-63</i>	0.889	0.867	0.856
<i>Tampered-0-to-63</i>	0.831	0.757	0.734
<i>Tampered-0-to-127</i>	0.753	0.555	0.398

Takeaway: the model suffers from shortcut learning!

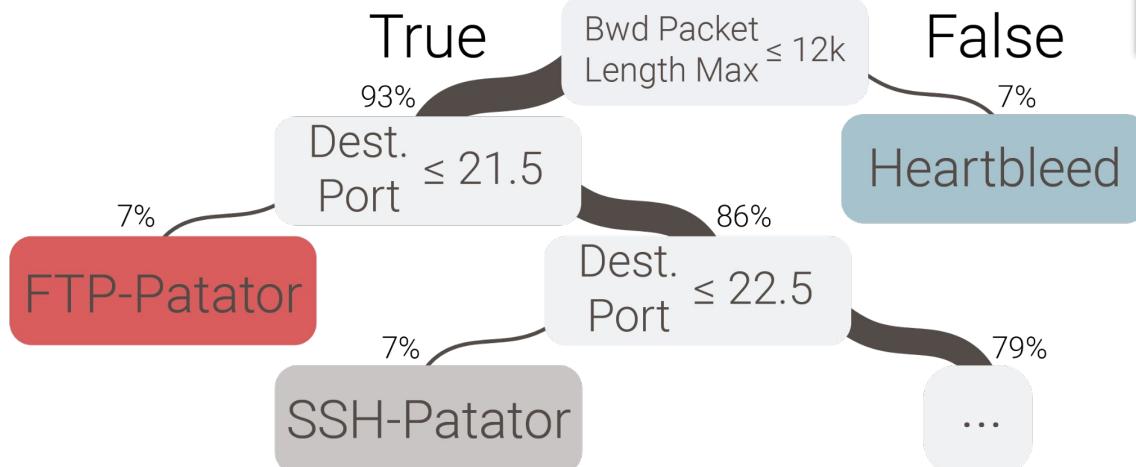
Use Case #2: Detecting Heartbleed Traffic

Problem Setup

- **Selected publications:**
 - Many papers that rely on the CIC-IDS-2017 dataset
 - "Toward Generating a New Intrusion Detection Dataset and Intrusion Traffic Characterization" – Sharafaldin et al., 2018
- **Proposal:**
 - **Model:** Random Forest to classify traffic between benign traffic and 13 different attacks (e.g. PortScan, DDoS, Heartbleed)
 - **Features:** 78 pre-computed features, from flow statistics (e.g. flow duration, mean IAT)
 - **Dataset:** CIC-IDS-2017 [<https://www.unb.ca/cic/datasets/ids-2017.html>]
- **Results:**
 - Reported F1-score: 0.99
 - Reproduced F1-score: 0.99

Use Case #2: Detecting Heartbleed Traffic

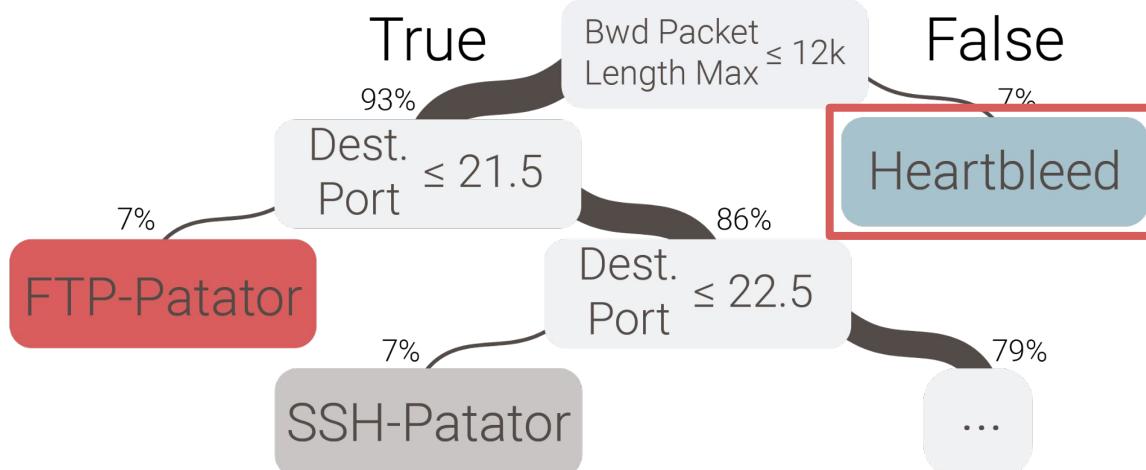
Explanation



Fidelity: 0.99
Top-3 pruning
6 nodes

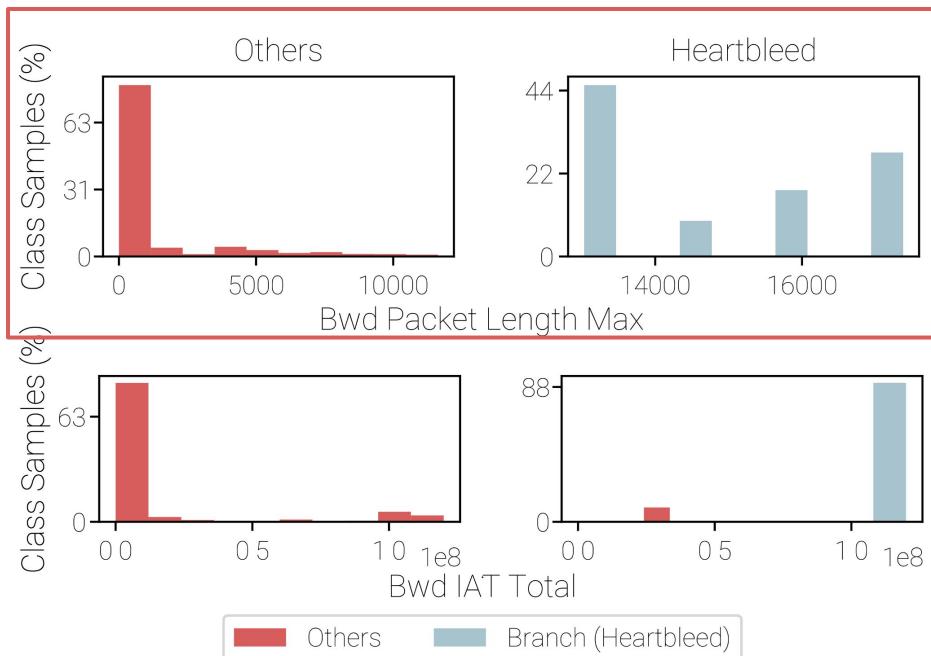
Use Case #2: Detecting Heartbleed Traffic

Explanation



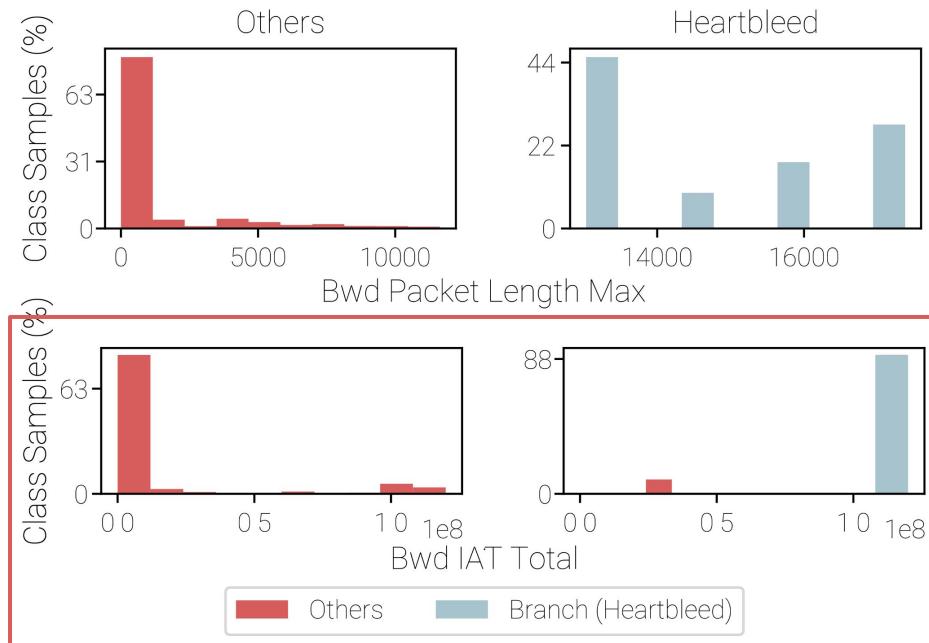
Use Case #2: Detecting Heartbleed Traffic

Explanation



Use Case #2: Detecting Heartbleed Traffic

Explanation



Use Case #2: Detecting Heartbleed Traffic

- Heartbleed attack:
 - An attacker sends an HTTPS heartbeat message with a value in the size field bigger than the message
 - e.g., 16k bytes packet with 64k bytes size value
 - A vulnerable server responds with a message with the size equal to the value specified in the size field and reveals information stored locally in its memory
 - e.g. server returns 64k bytes (16k from packet and 48k from memory)
- In the CIC-IDS-2017 dataset:
 - HTTPS connection was never closed during the duration of the attack
 - Huge number of backward bytes and very high IAT in the flow!

Use Case #2: Detecting Heartbleed Traffic

Validation

- Validation dataset:
 - 1000 new heartbleed flows **closing connection after every heartbeat**
 - **Backward bytes** and **IAT** similar to benign traffic

Class	Precision	Recall	F1
<i>Heartbleed (i.i.d.)</i>	1.000	1.000	1.000
<i>Heartbleed (o.o.d)</i>	0.000	0.000	0.000

Use Case #2: Detecting Heartbleed Traffic

Validation

- Validation dataset:
 - 1000 new heartbleed flows closing connection after every heartbeat
 - Backward bytes and IAT similar to benign traffic

Class	Precision	Recall	F1
<i>Heartbleed (i.i.d.)</i>	1.000	1.000	1.000
<i>Heartbleed (o.o.d.)</i>	0.000	0.000	0.000

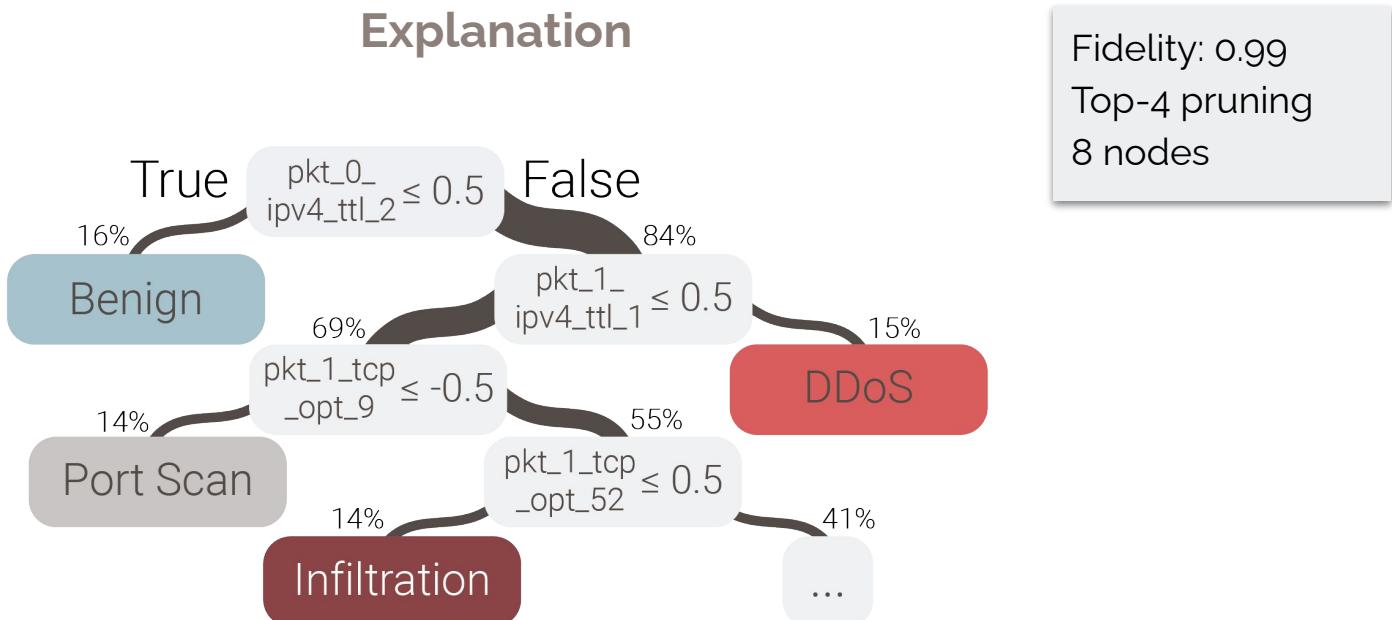
Takeaway: the model is overfitted to training data and fails to identify o.o.d. samples!

Use Case #3: Inferring Malicious Traffic for IDS

Problem Setup

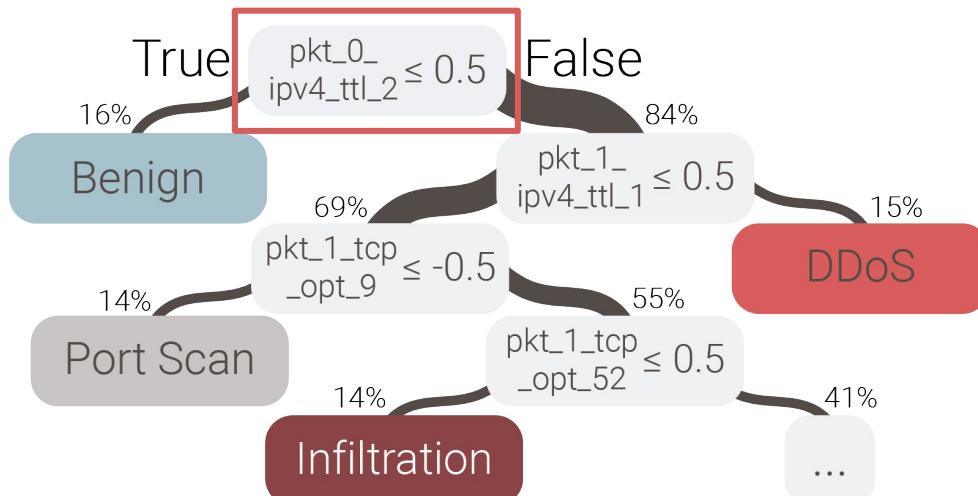
- **Selected publications:**
 - “New Directions in Automated Traffic Analysis” — Holland et al., 2020
- **Proposal:**
 - **Model:** nPrintML, an AutoML model for an Intrusion Detection System (IDS). It uses AutoGluton to ensemble multiple single models.
 - **Features:** 4,480 features with values -1, 0, or 1, each feature represents a bit of a set of pre-established protocol headers.
 - **Dataset:** CIC-IDS-2017 [<https://www.unb.ca/cic/datasets/ids-2017.html>]
- **Results:**
 - Reported F1-score: 0.99
 - Reproduced F1-score: 0.99

Use Case #3: Inferring Malicious Traffic for IDS



Use Case #3: Inferring Malicious Traffic for IDS

Explanation



Use Case #3: Inferring Malicious Traffic for IDS

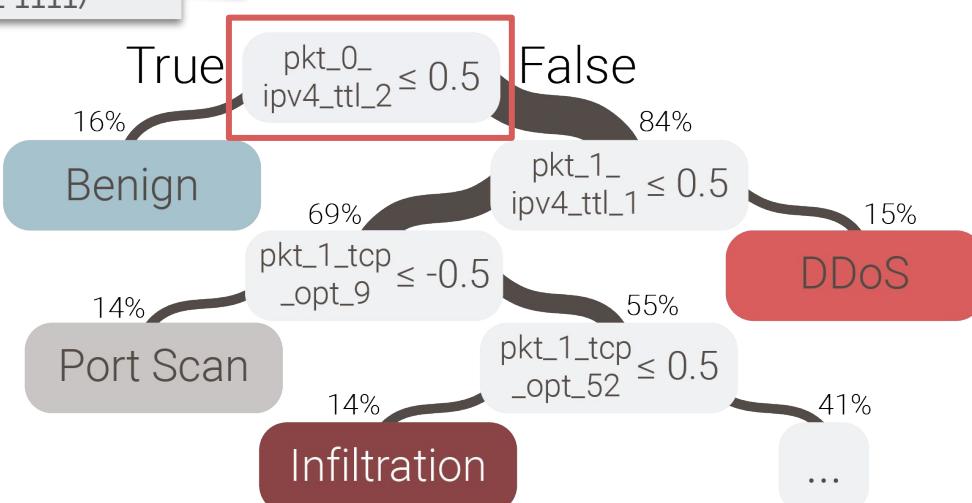
Kali Linux

Init TTL = 64

TTL - 1 hop = 63

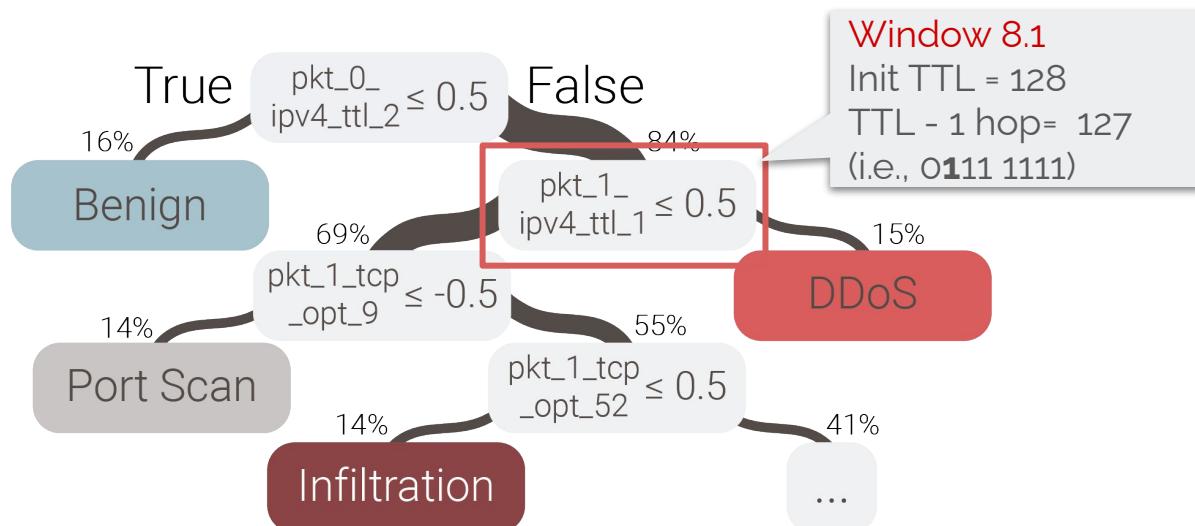
(i.e., 0011 1111)

Explanation



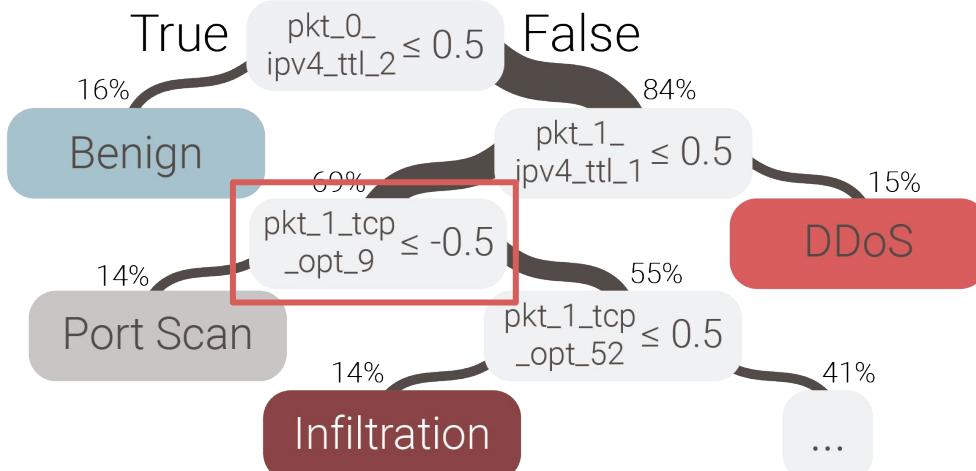
Use Case #3: Inferring Malicious Traffic for IDS

Explanation



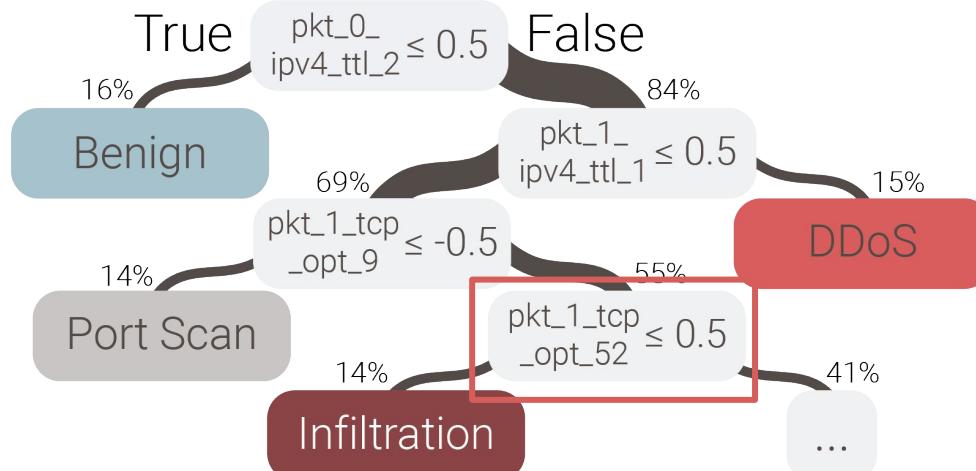
Use Case #3: Inferring Malicious Traffic for IDS

Explanation



Use Case #3: Inferring Malicious Traffic for IDS

Explanation



Use Case #3: Inferring Malicious Traffic for IDS

Validation

- Validation dataset:
 - Curated balanced dataset with **4,047 flows** from **real-world traffic** in UCSB network
 - Used **Suricata-IDS** to generate flow labels

Class	Precision	Recall	F1
<i>Benign</i>	0.653	0.806	0.722
<i>DoS</i>	0.000	0.000	0.000
<i>Port Scan</i>	0.120	0.143	0.130
Average	0.256	0.315	0.282

Use Case #3: Inferring Malicious Traffic for IDS

Validation

- Validation dataset:
 - Curated balanced dataset with **4,047 flows** from **real-world traffic** in UCSB network
 - Used **Suricata-IDS** to generate flow labels

Class	Precision	Recall	F1
<i>Benign</i>	0.653	0.806	0.722
<i>DoS</i>	0.000	0.000	0.000
<i>Port Scan</i>	0.120	0.143	0.130
Average	0.256	0.315	0.282

Takeaway: the model suffers from spurious correlations in the training data!

Other Use Cases

Problem	Model(s)	Dataset(s)	Trustee Fidelity	Inductive Bias
Detect VPN traffic (Wang <i>et al.</i> , ISI'17)	1-D CNN	ISCX VPN-nonVPN	1.00	Shortcut learning
Detect Heartbleed traffic (Sharafaldin <i>et al.</i> , ICISSP'18)	RFC	CIC-IDS-2017	0.99	O.O.D.
Detect Malicious traffic (IDS) (Holland <i>et al.</i> , CCS'21)	nPrintML	CIC-IDS-2017	0.99	Spurious Correlation
Anomaly Detection (Mirsky <i>et al.</i> , NDSS'18)	Kitsune	Mirai dataset	0.99	O.O.D
OS Fingerprinting (Holland <i>et al.</i> , CCS'21)	nPrintML	CIC-IDS-2017	0.99	O.O.D
IoT Device Fingerprinting (Xiong <i>et al.</i> , HotNets'19)	Iisy	UNSW-IoT	0.99	Shortcut learning
Adaptive Bit-rate (Mao <i>et al.</i> , SIGCOMM'17)	Pensieve	HSDPA Norway	0.99	O.O.D

Conclusions

1. ML in high-stakes requires trust
2. Trustee improves trust!
3. Trustee can be used with any existing model
4. Trustee is ready to be used!
 - More in the afternoon session

Thank you!



<https://trusteeml.github.io>

Trustee Python package

- <https://pypi.org/project/trustee/>
Trustee Repository
- <https://github.com/TrusteeML/trustee>
Use Cases Repository
- <https://github.com/TrusteeML/emperor>

