



AI/ML for Network Security: The Emperor has no Clothes

Arthur Selle Jacobs¹

Roman Beltiukov⁴

Walter Willinger²

Ronaldo A. Ferreira³

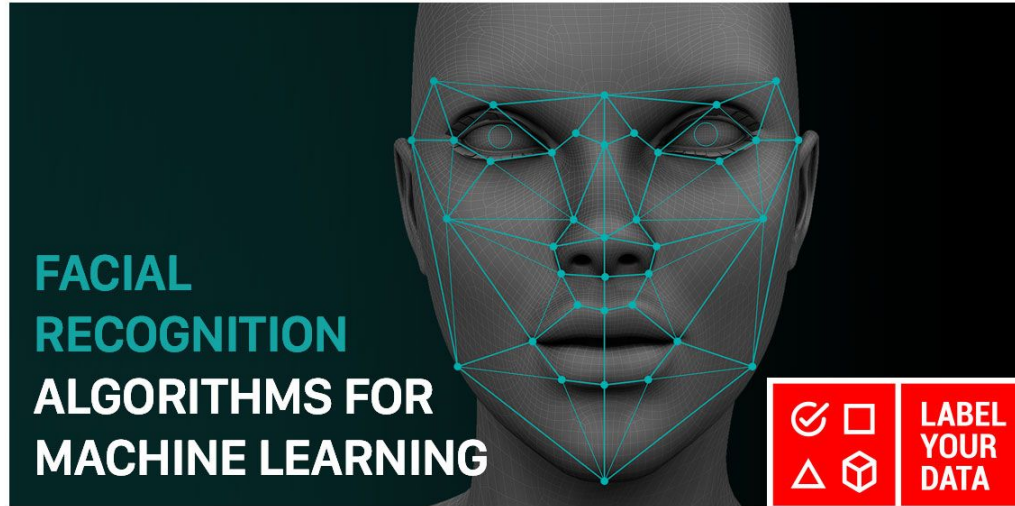
Arpit Gupta⁴

Lisandro Z. Granville¹

March 27th, 2023



The Rise of AI



The Rise of AI

FACIAL
RECOGNITION
ALGORITHMS FOR
MACHINE LEARNING



The Rise of AI

FACIAL
RECOGNITION
ALGORITHMS F
MACHINE LEAF

AI & MACHINE LEARNING

How Kaggle solved a spam problem in 8 days using AutoML

Will Cukierski
Staff Developer Advocate and
Head of Competitions, Kaggle

May 27, 2020

Try Google Cloud

Start building on Google Cloud
with \$300 in free credits and
20+ always free products.

FREE TRIAL

Kaggle is a data science community of nearly 5 million users. In September of 2019, we found ourselves under a sudden siege of spam traffic that threatened to overwhelm visitors to our site. We had to come up with an effective solution, fast. Using AutoML Natural Language on Google Cloud, Kaggle was able to train, test, and deploy a spam detection model to production in just eight days. In this post, we'll detail our success story about using machine learning to rapidly solve an urgent business dilemma.

A spam dilemma

Malicious users were suddenly creating large numbers of Kaggle accounts in order to leave spammy search engine optimization (SEO) content in the user bio section. Search engines were indexing these bios, and our existing spam detection heuristics were failing to flag them. In short, we faced a growing and embarrassing predicament.

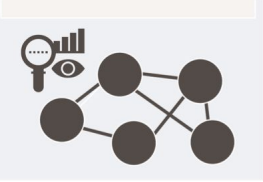
Our problem was context. Kaggle is a community focused on data science and machine learning. As a result of our topical data-science focus, a user bio that seems harmless in isolation may be the work of a spammer. Here is a real example of one such bio:



How does it work?

Traditional AI/ML Development Pipeline

Collect Data

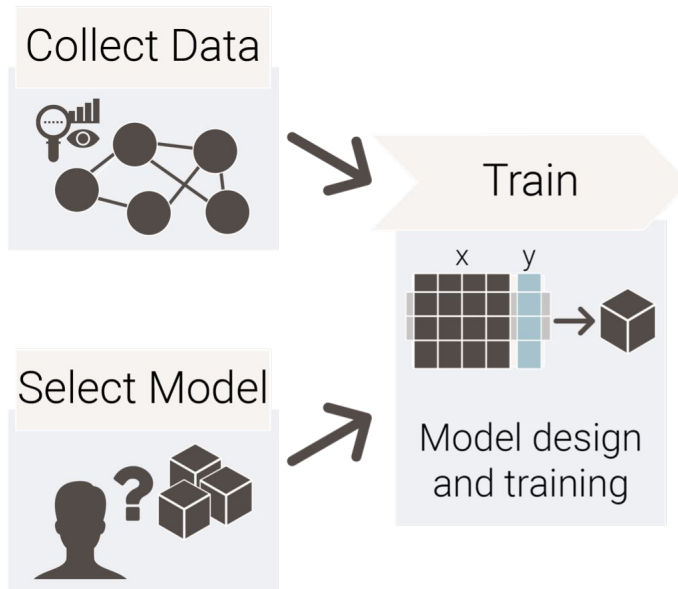


Select Model



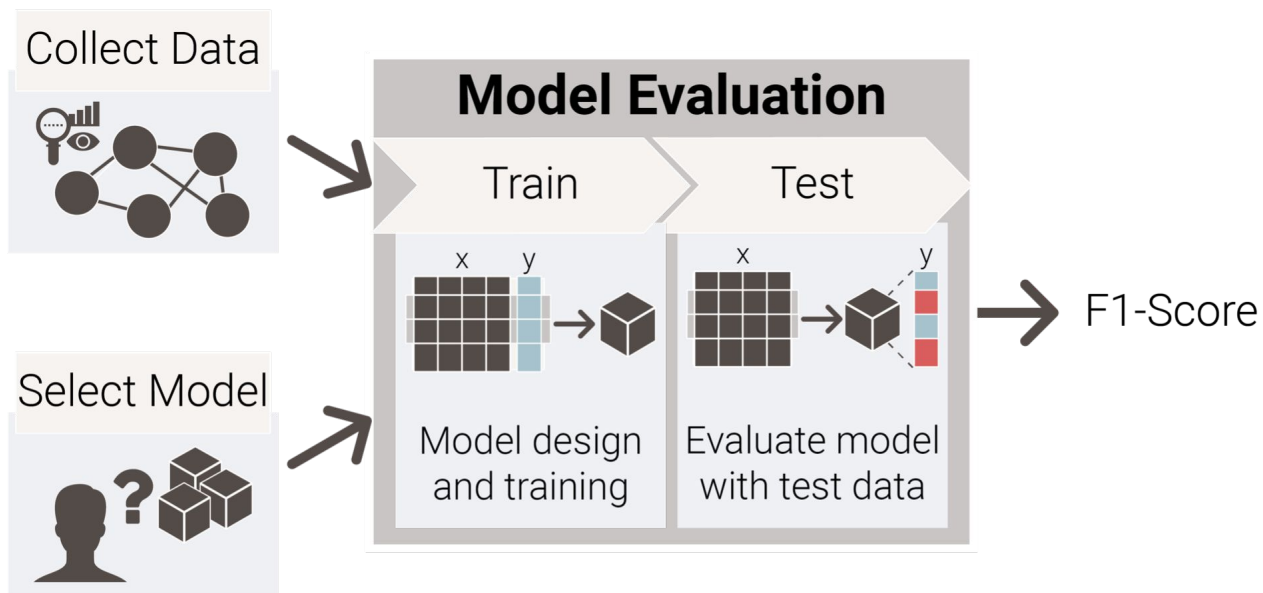
How does it work?

Traditional AI/ML Development Pipeline



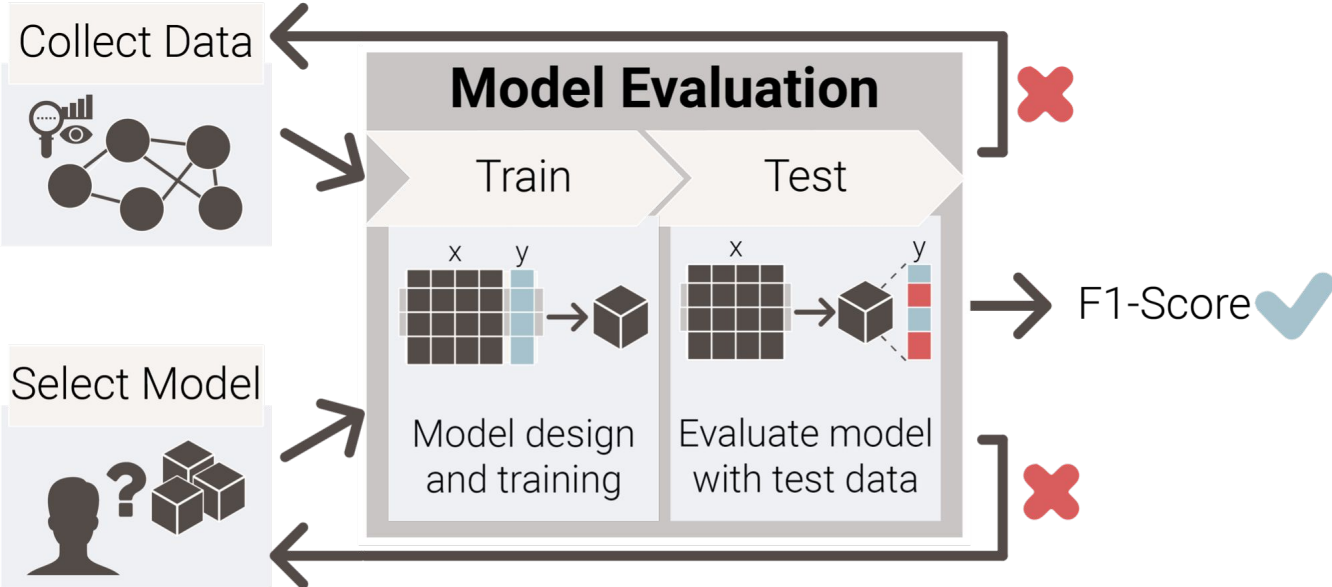
How does it work?

Traditional AI/ML Development Pipeline



How does it work?

Traditional AI/ML Development Pipeline



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!

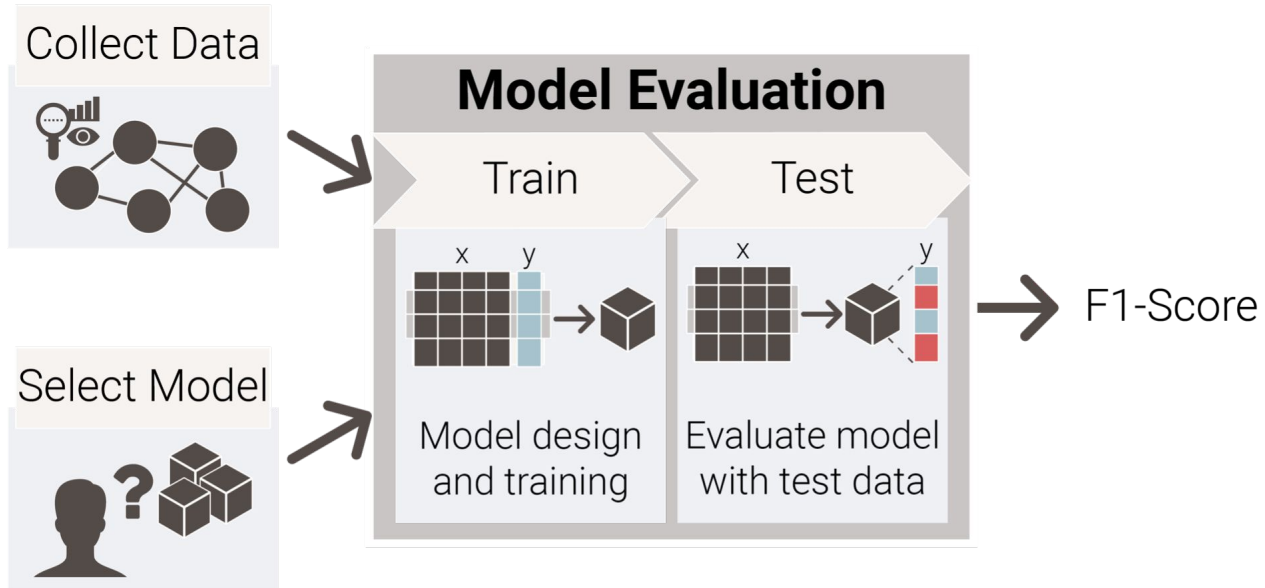
O.O.D. Samples

Model does not generalize!

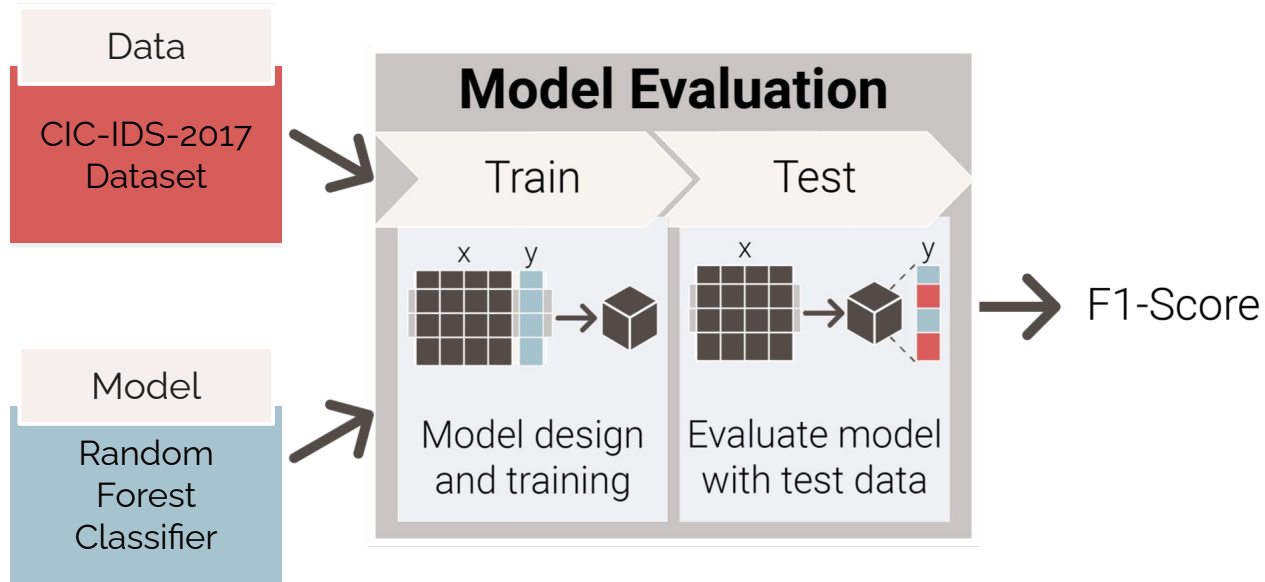
Spurious Correlations

Model picks up wrong correlations in the data!

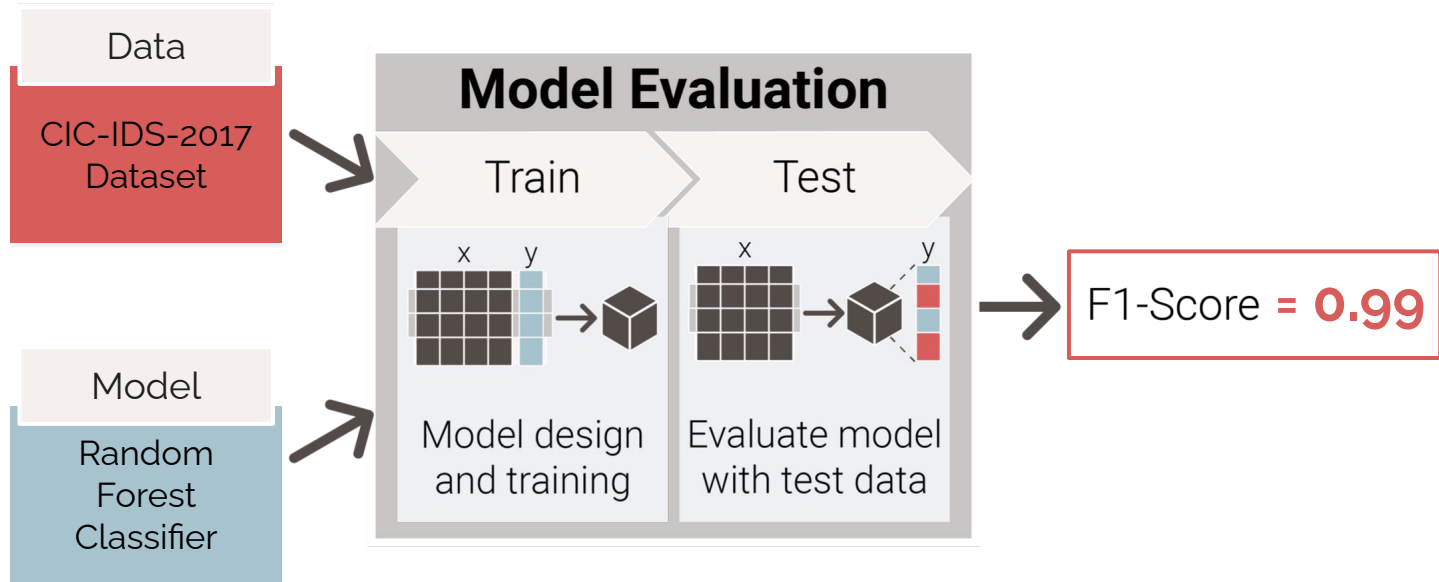
Consider this example...



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?

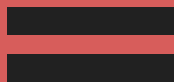




Can you **trust** this model?

Can you **trust** this model?

Trust in AI/ML model

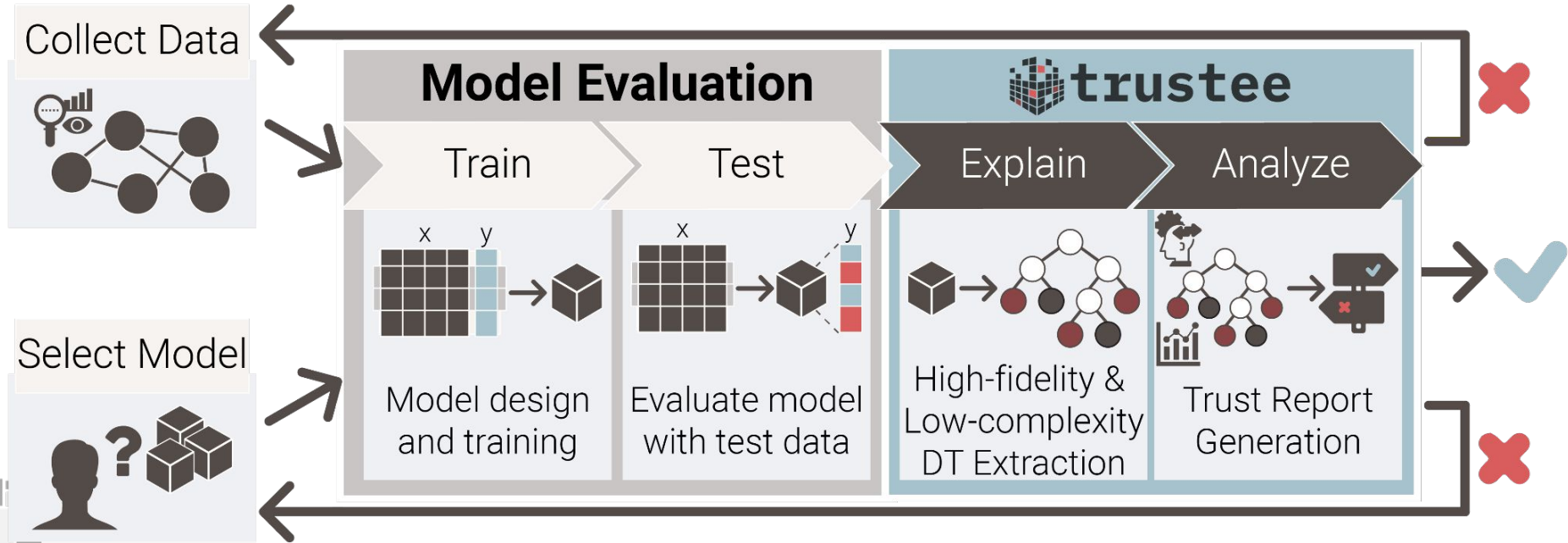


Hand over control to the AI/ML model

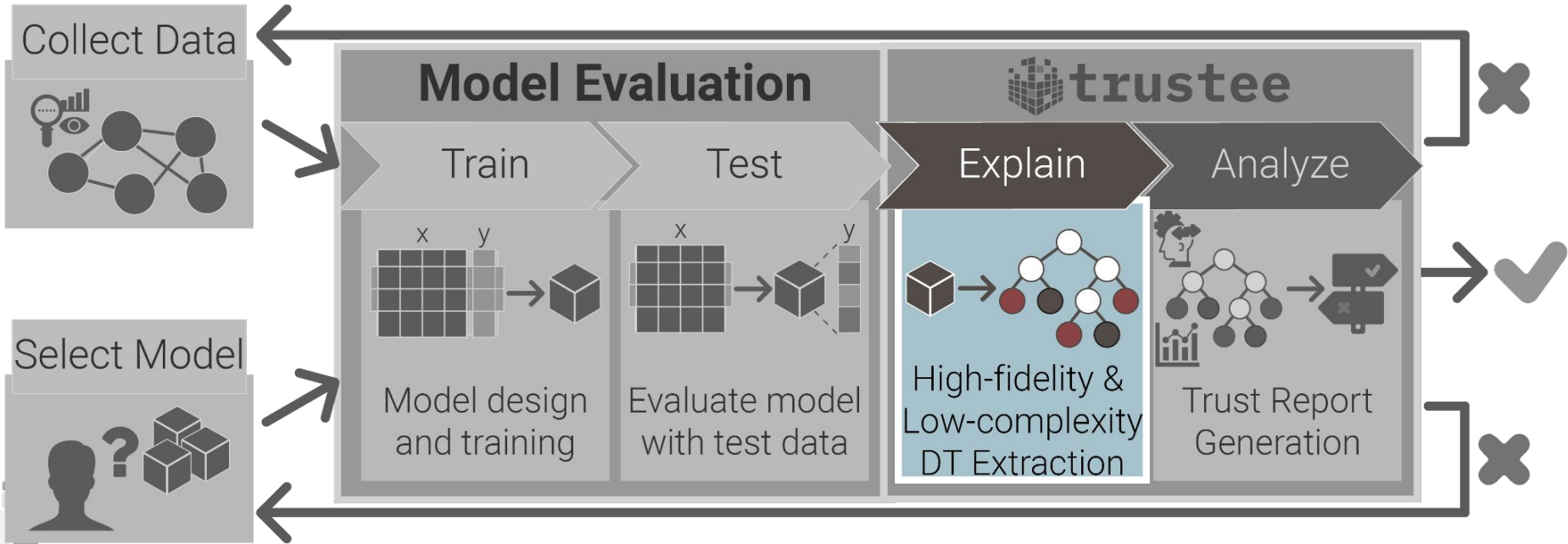


trustee

Augmented AI/ML Development Pipeline

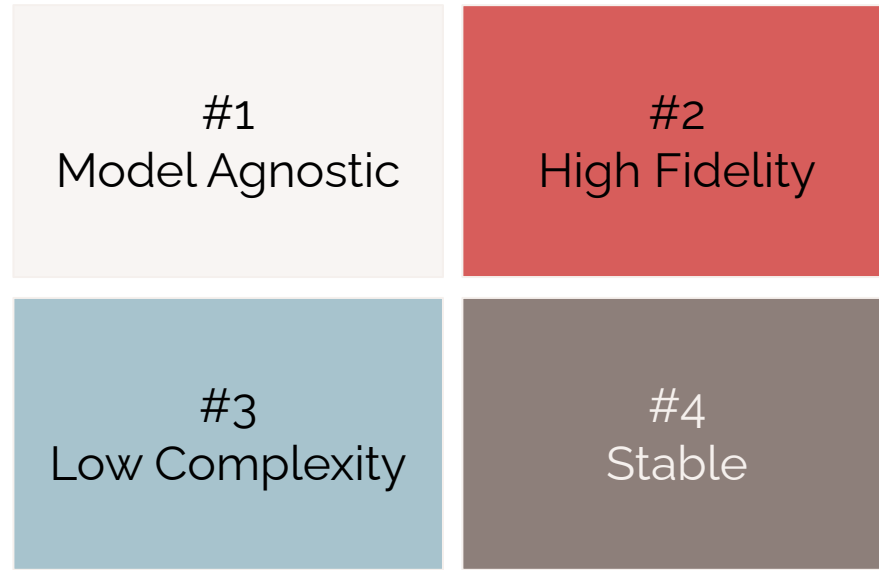


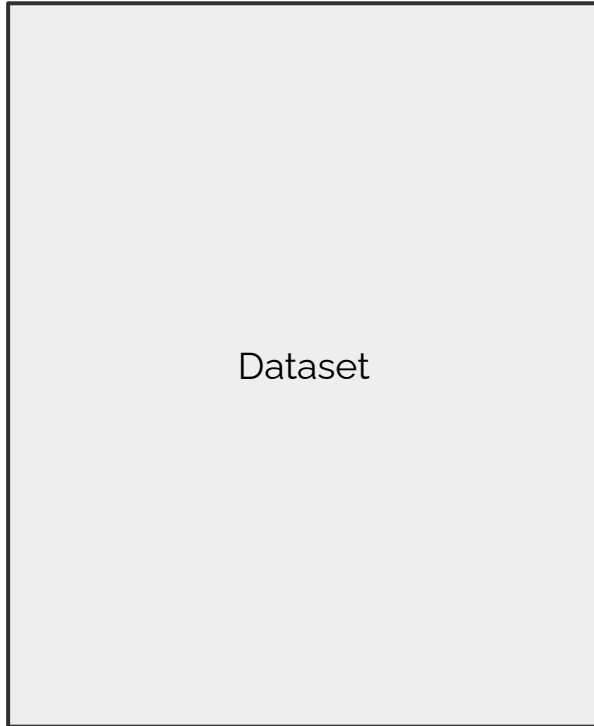
Augmented AI/ML Development Pipeline



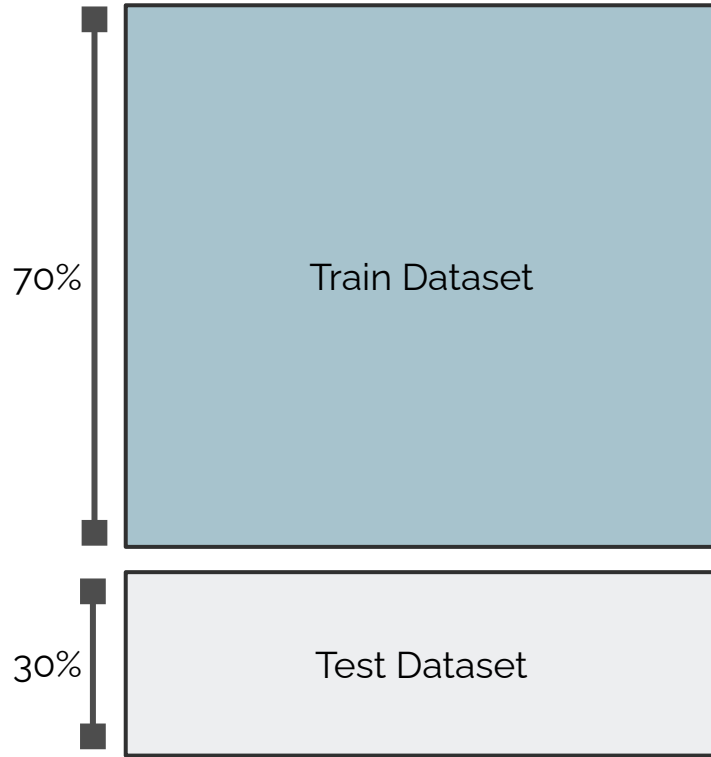


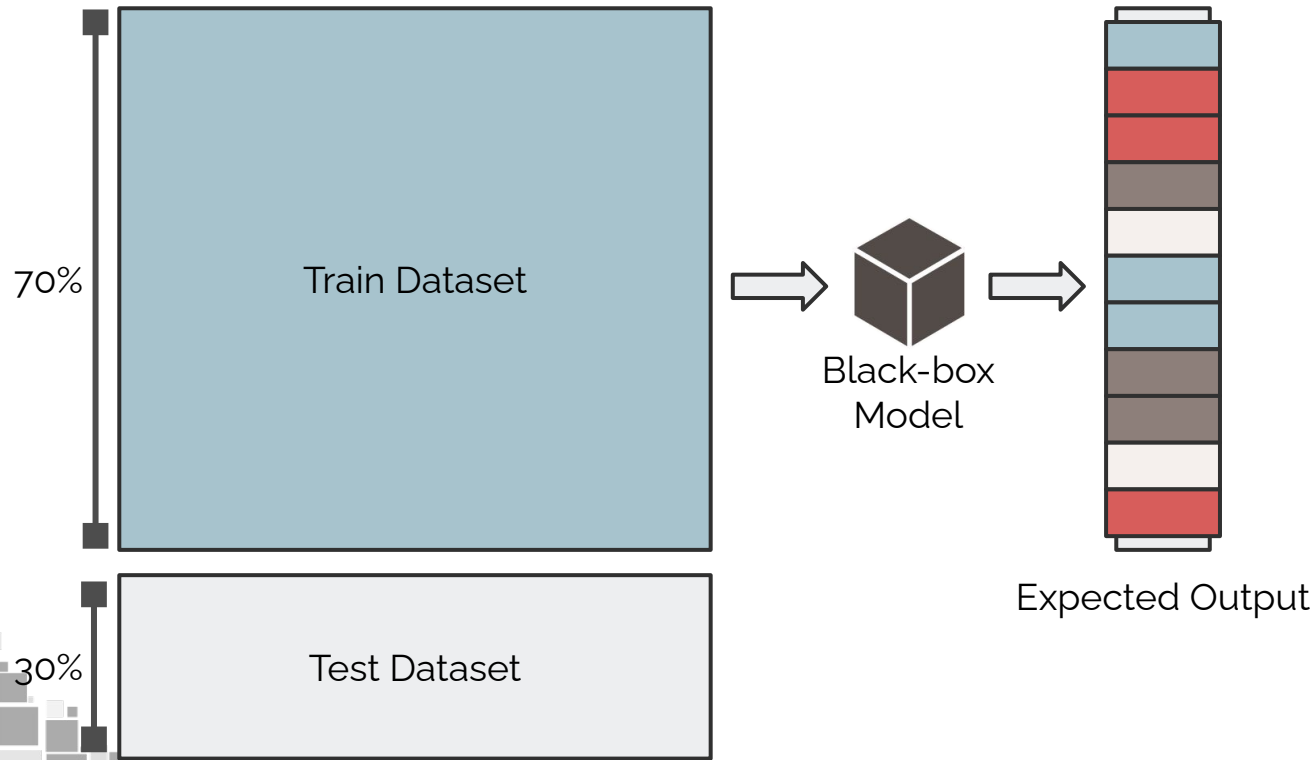
Explanation Requirements



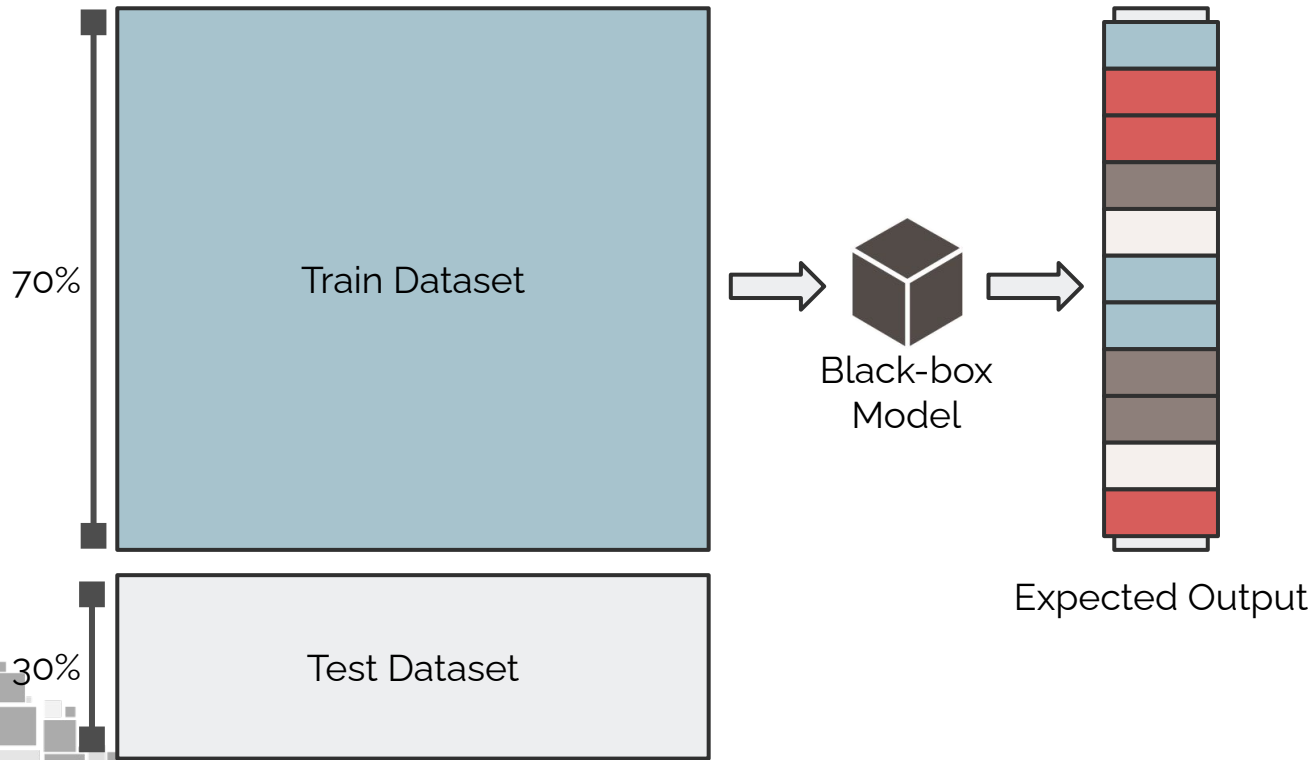


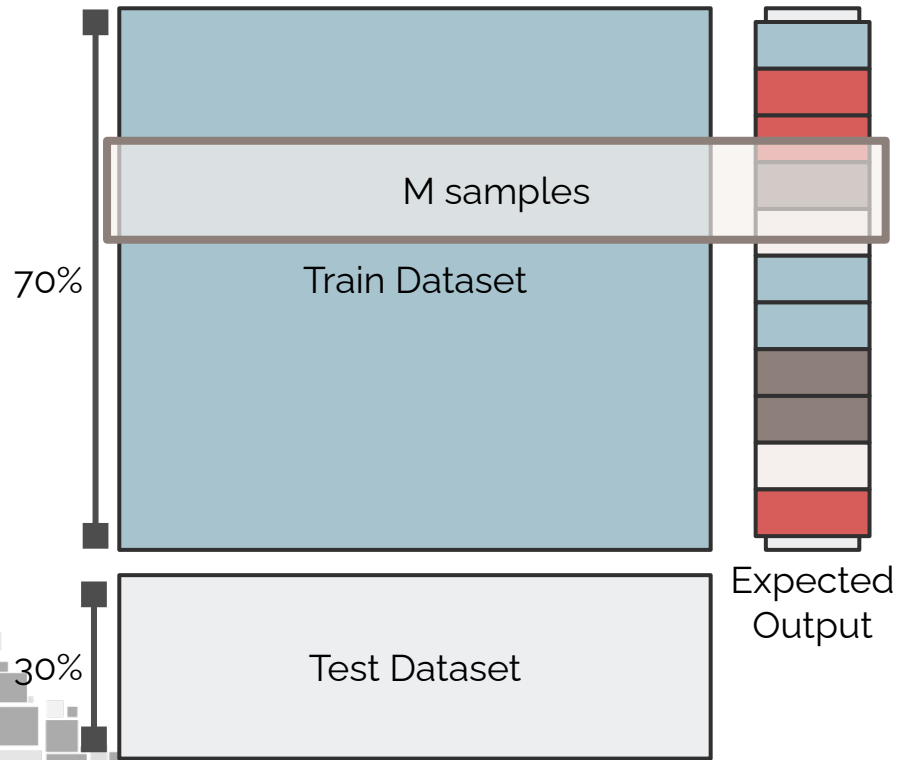
Black-box
Model

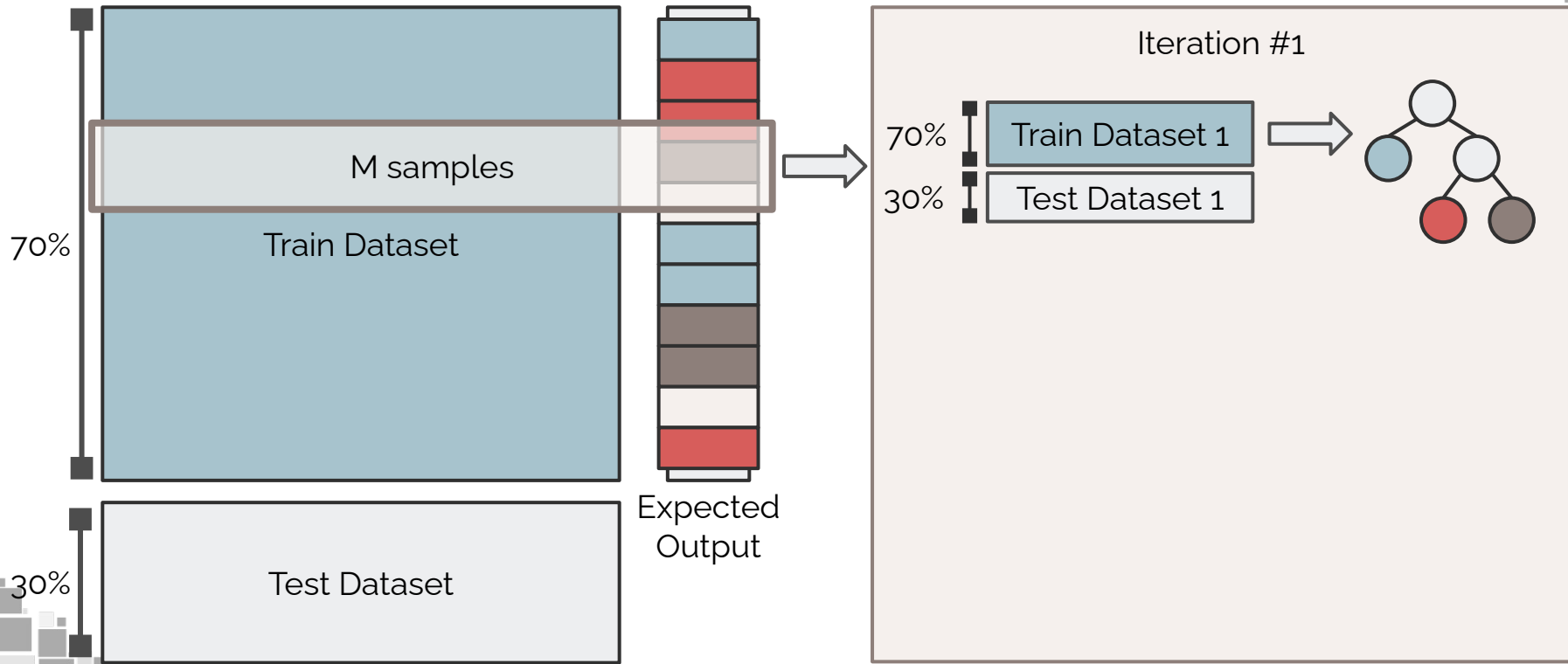


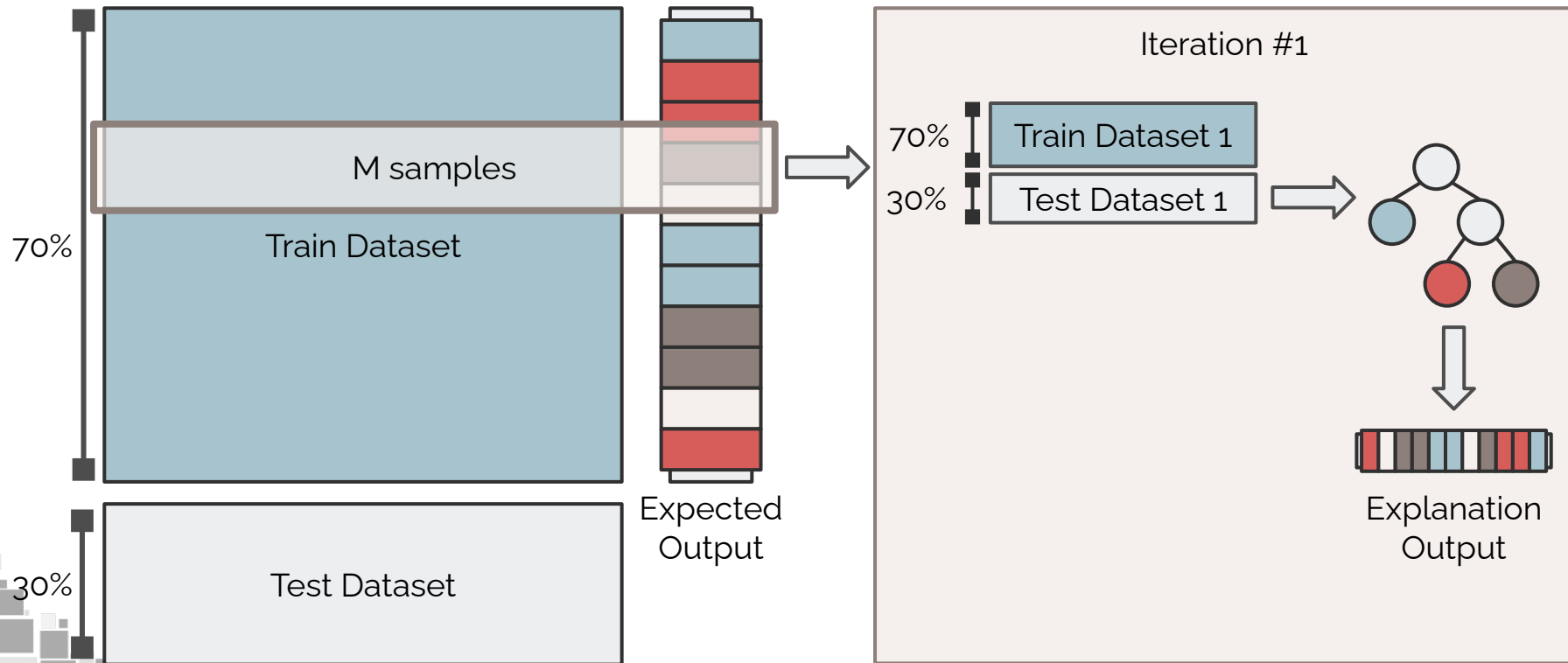


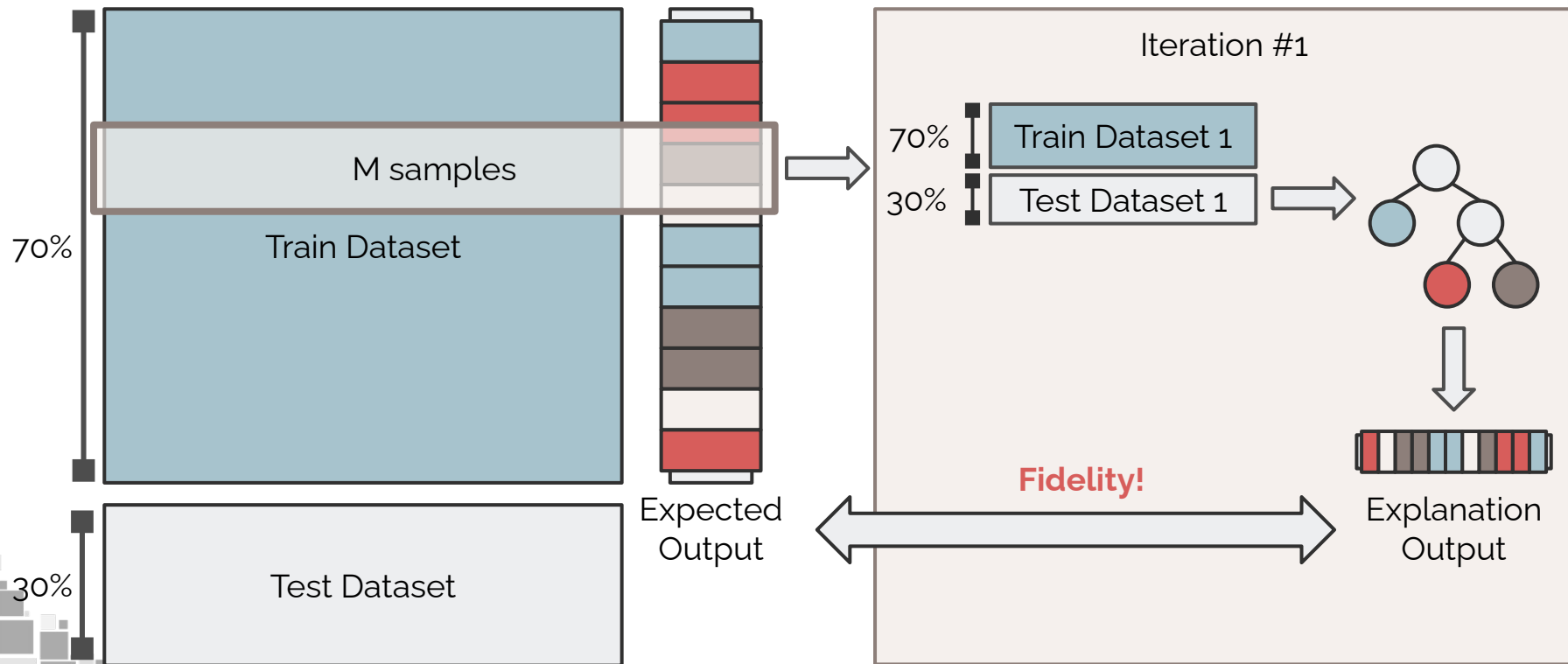
#1
Model
Agnostic

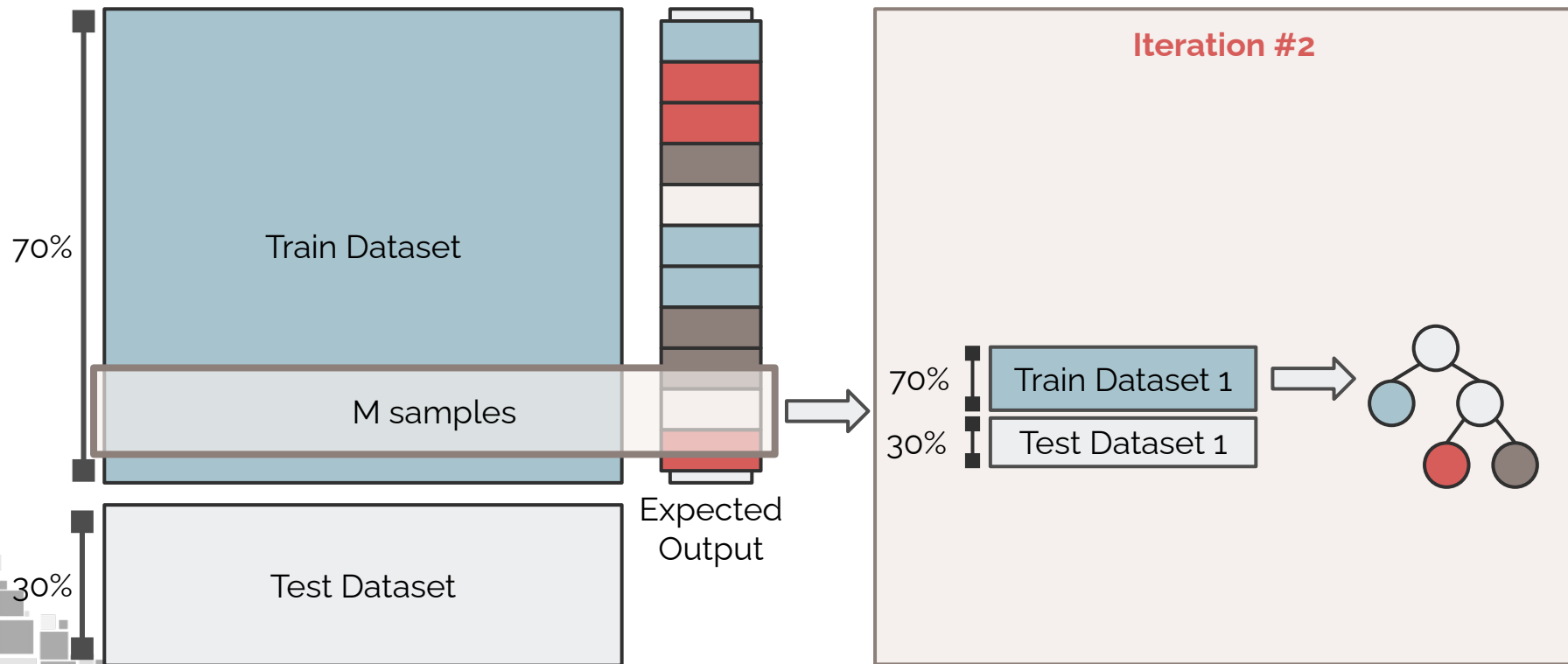


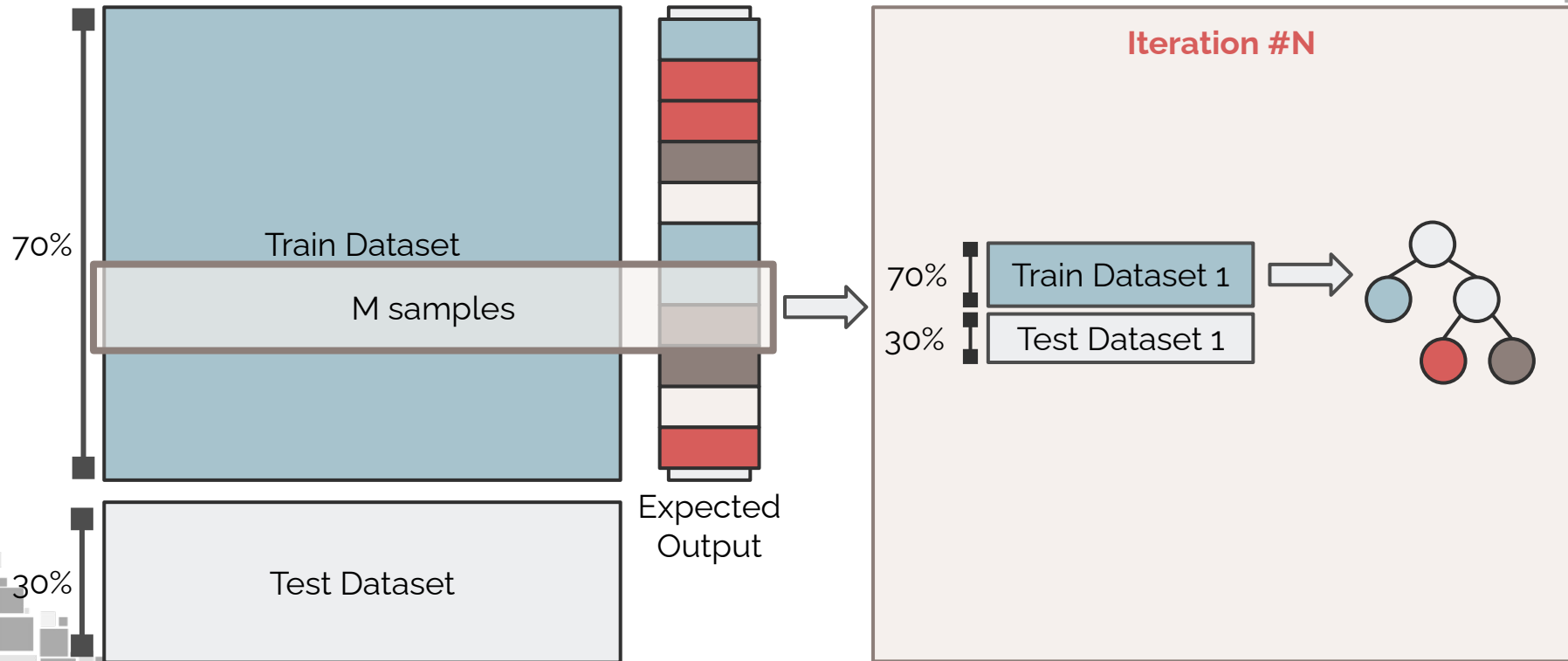


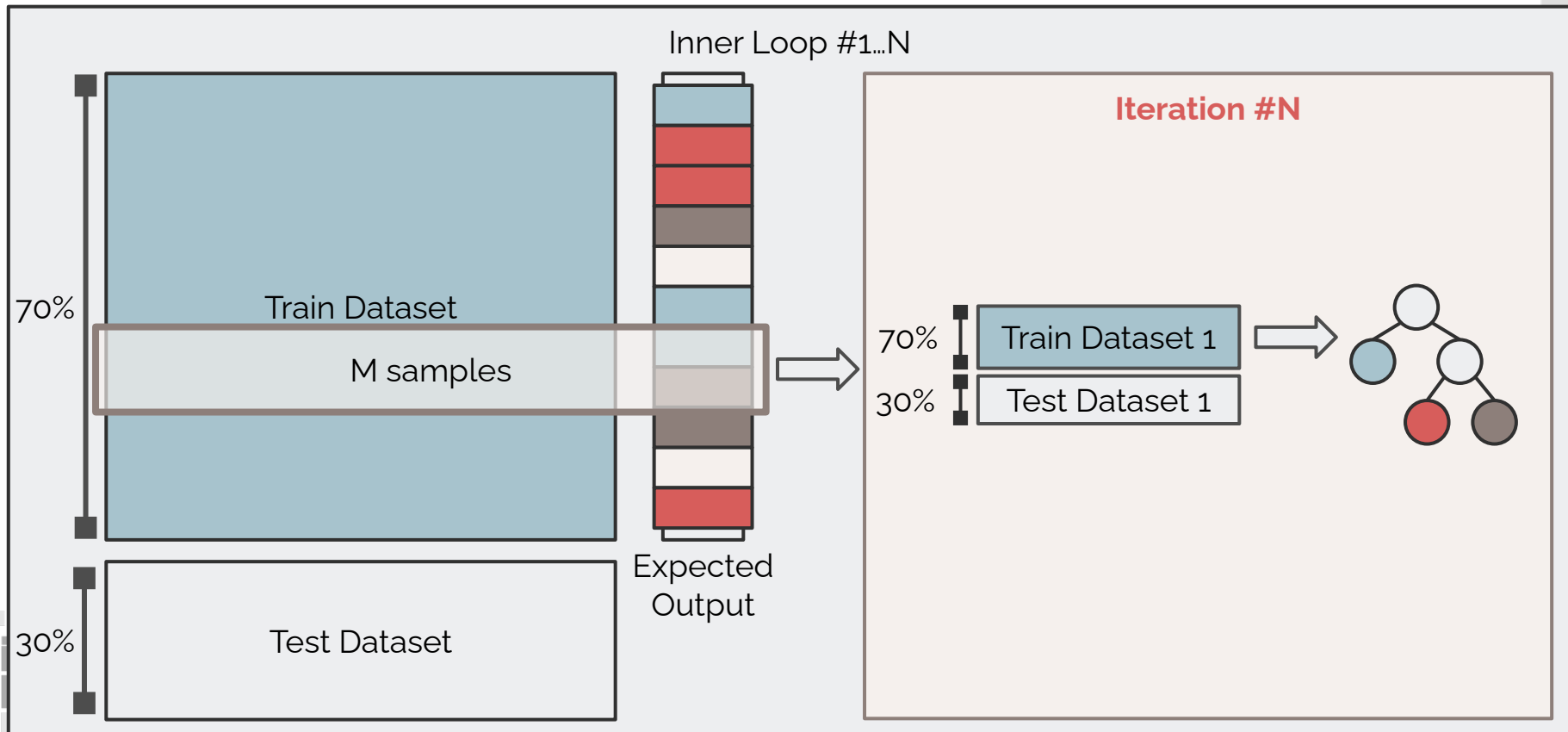


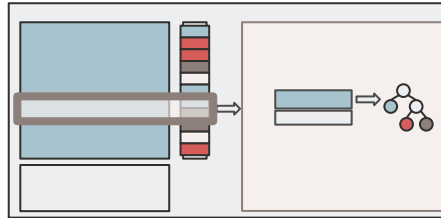


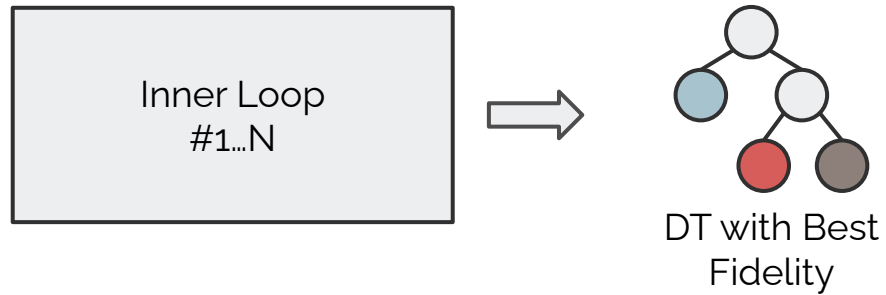




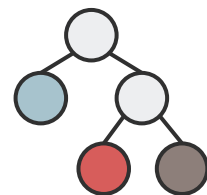
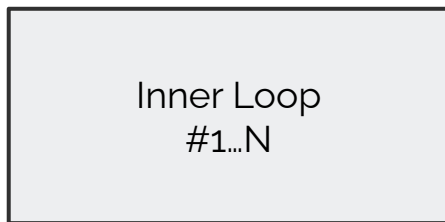




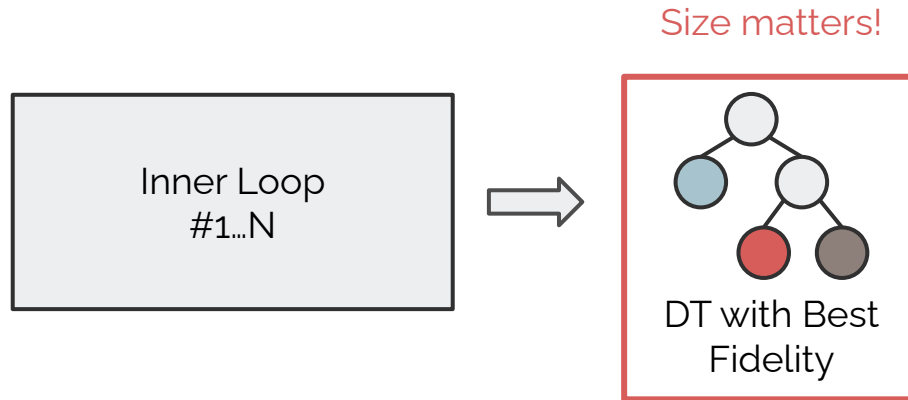


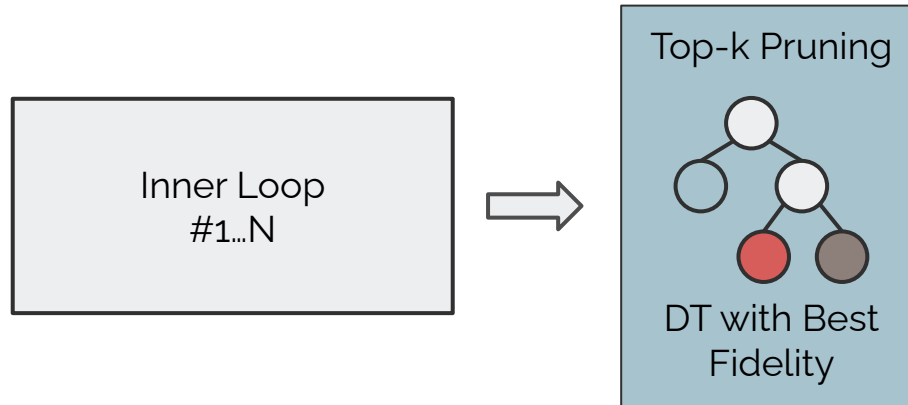


#2
High
Fidelity



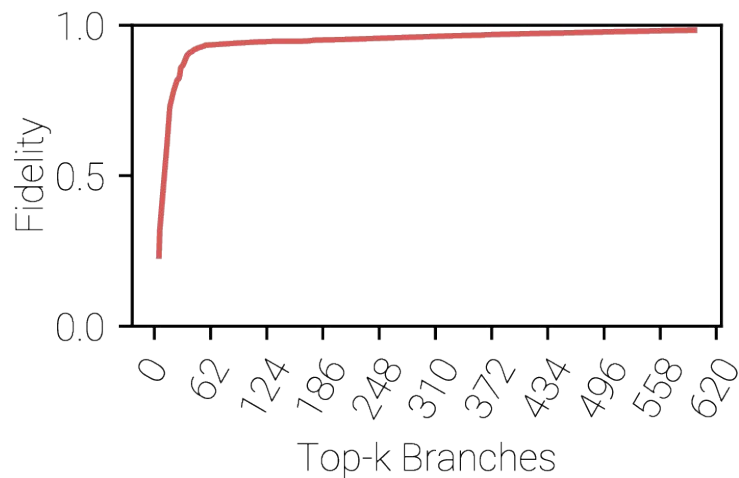
DT with Best
Fidelity



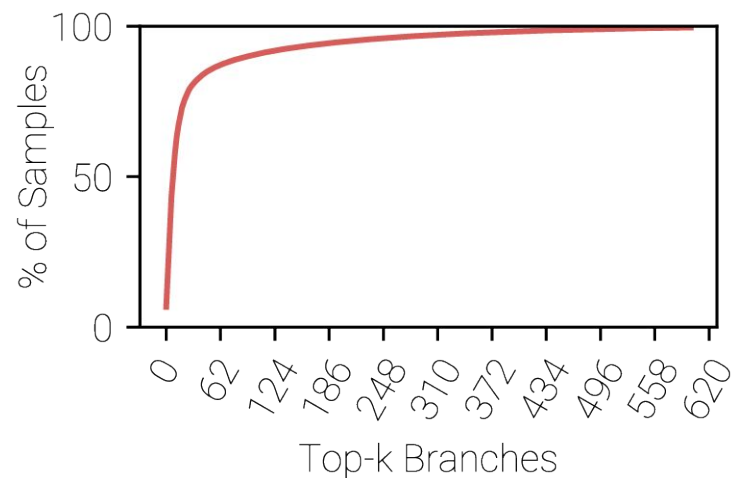


Top-k Pruning

Fidelity

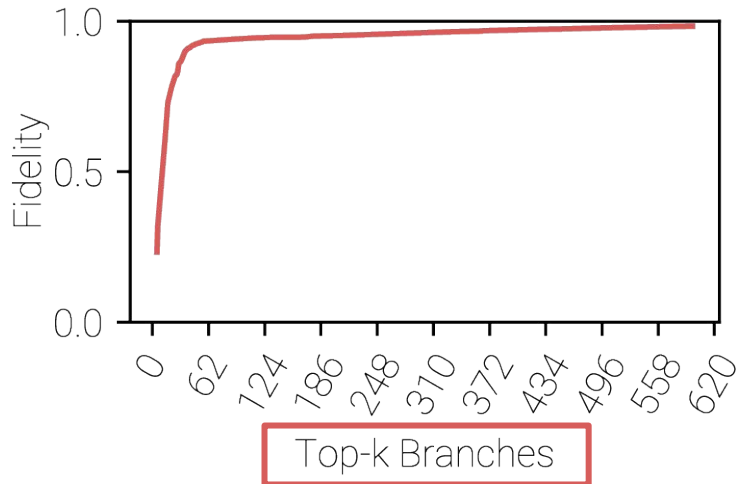


Samples

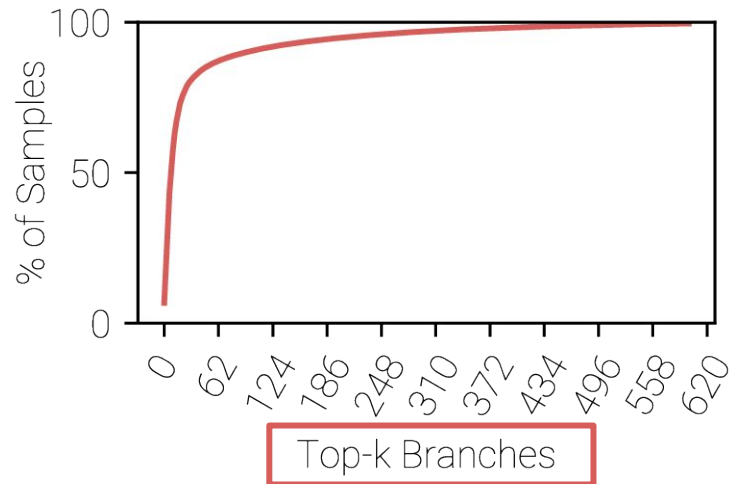


Top-k Pruning

Fidelity



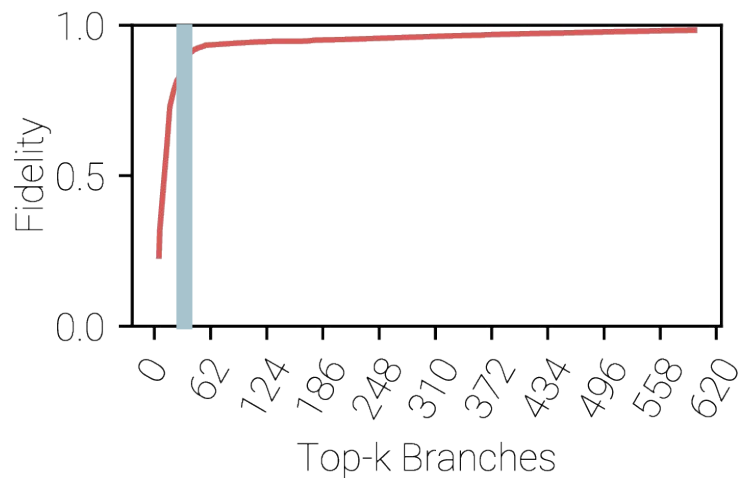
Samples



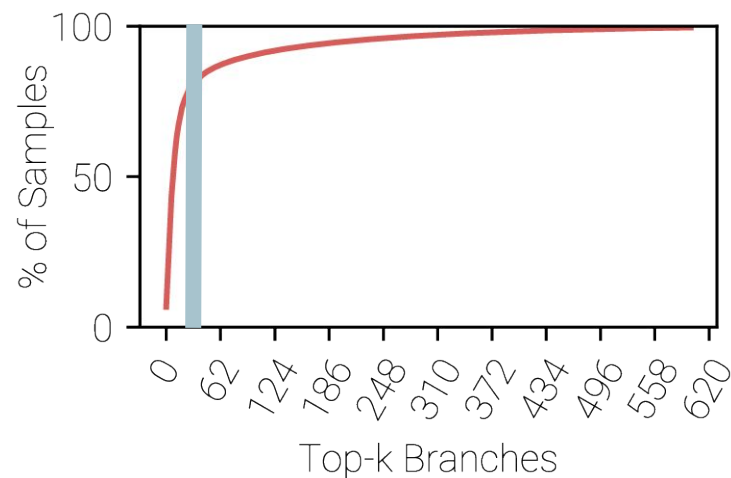
Diminishing returns!

Top-k Pruning

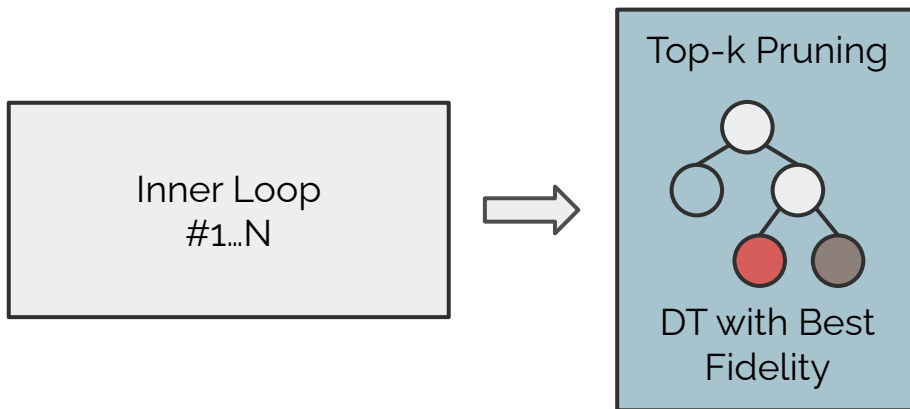
Fidelity

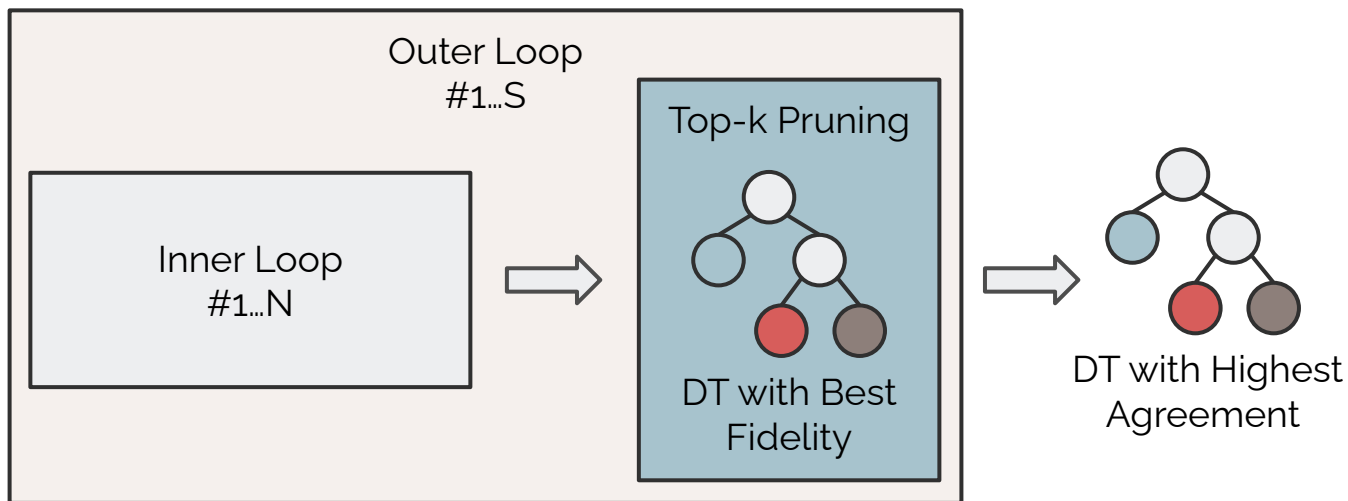


Samples

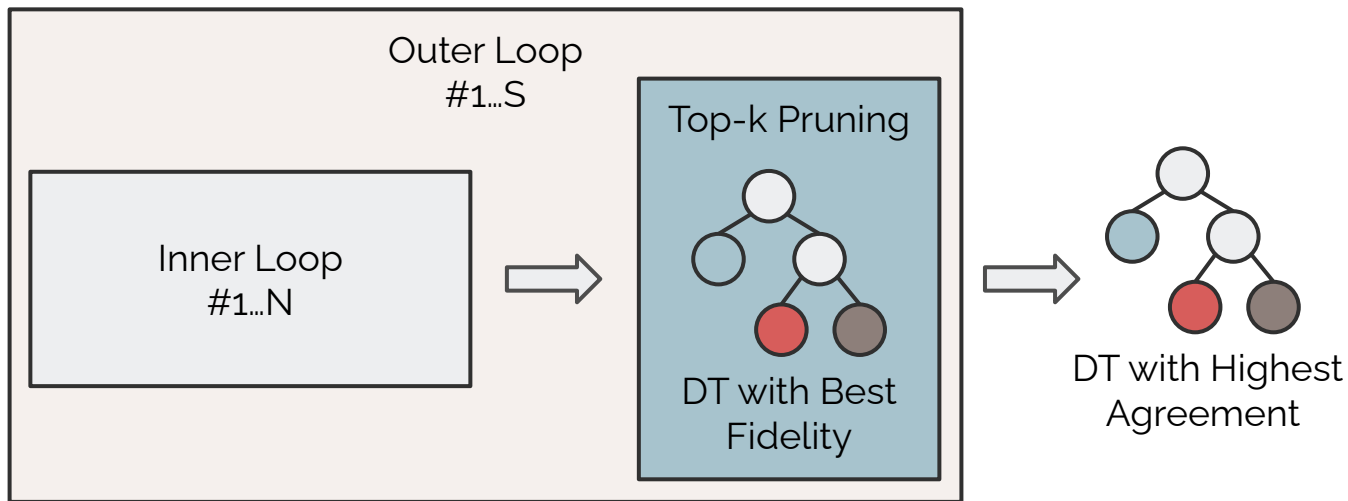


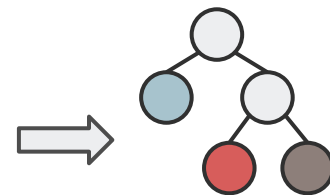
#3
Low
Complexity





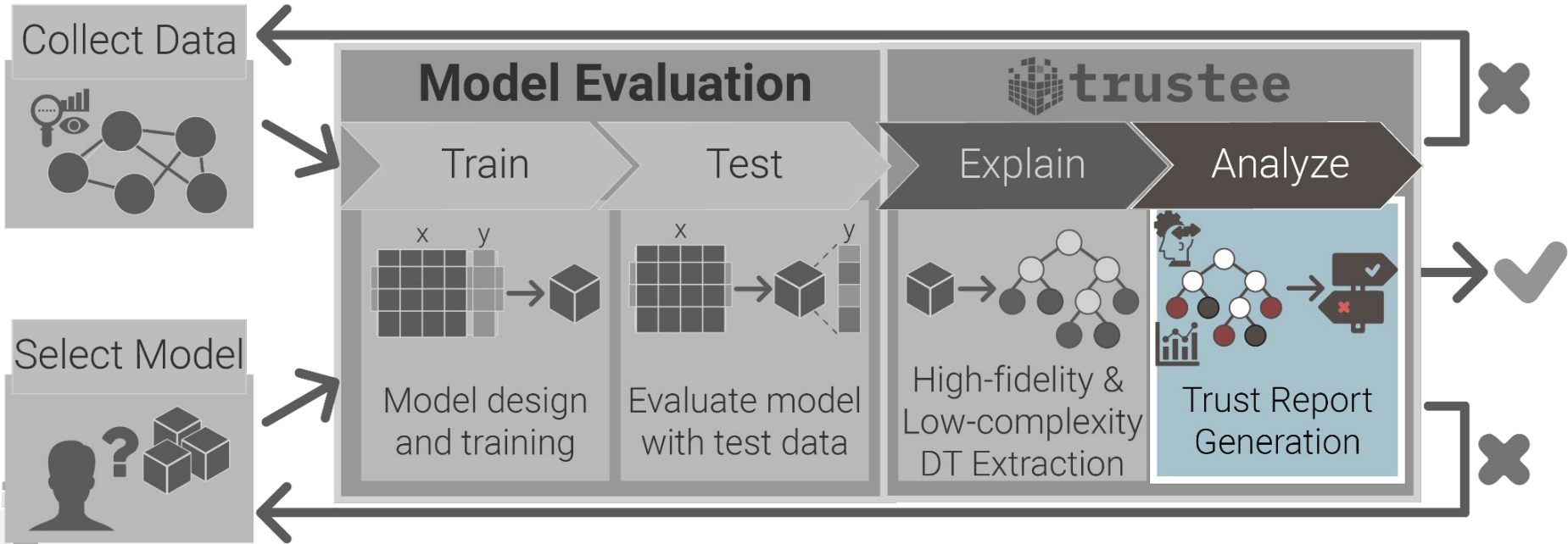
#4
Stable





DT with Highest
Agreement

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!

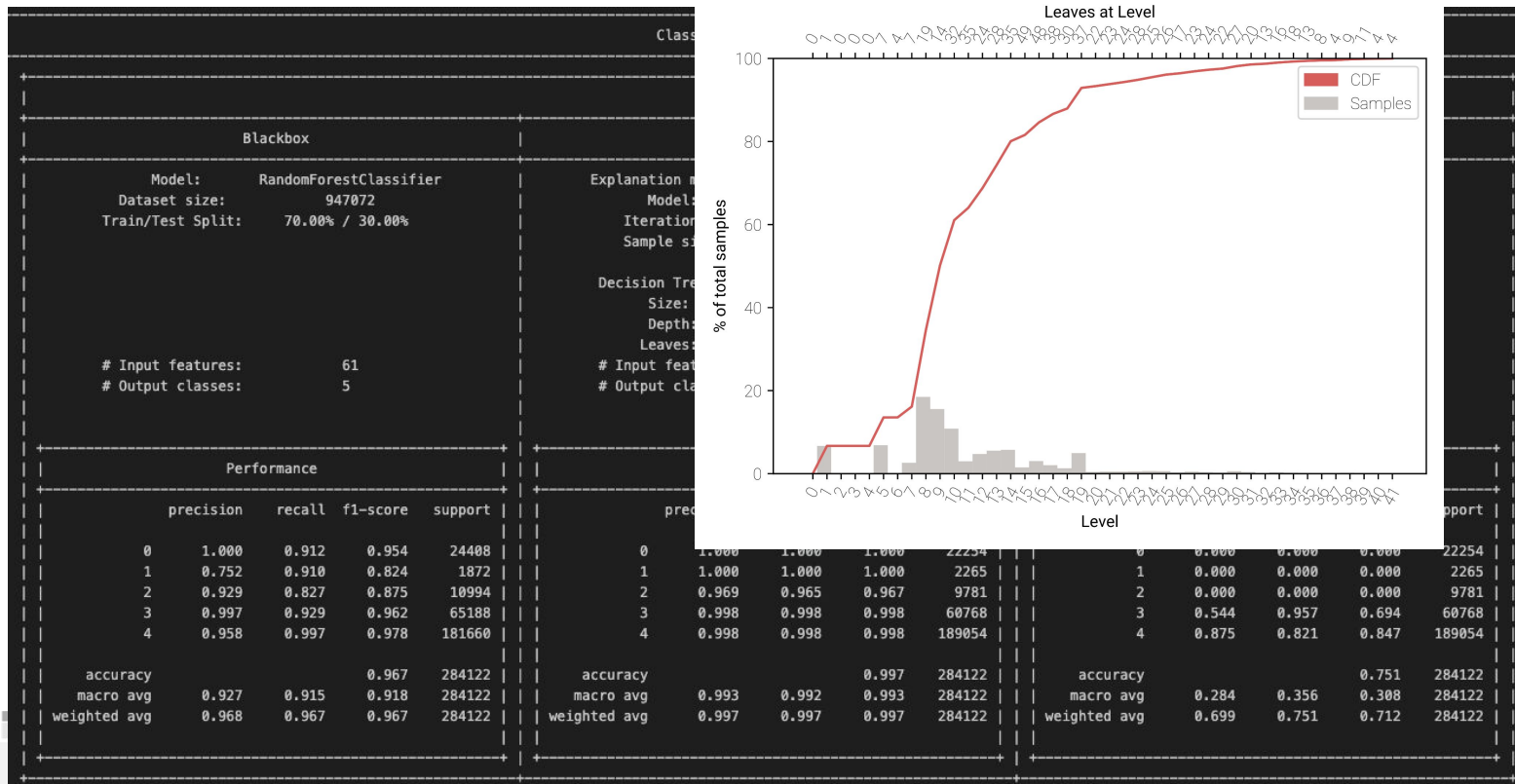
Generating Trust Reports

Classification Trust Report															
Summary															
Blackbox				Whitebox				Top-k Whitebox							
Model:		RandomForestClassifier		Explanation method:		Trustee		Explanation method:		Trustee					
Dataset size:		947072		Model:		DecisionTreeClassifier		Model:		DecisionTreeClassifier					
Train/Test Split:		70.00% / 30.00%		Iterations:		1		Iterations:		1					
				Sample size:		50.00%		Sample size:		50.00%					
				Decision Tree Info				Decision Tree Info							
				Size:		2437		Size:		9					
				Depth:		31		Depth:		4					
				Leaves:		1219		Leaves:		5					
# Input features:		61		# Input features:		18 (29.51%)		# Input features:		-					
# Output classes:		5		# 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	22254
	1	0.752	0.910	0.824	1872	1	1.000	1.000	1.000	2265	1	0.000	0.000	0.000	2265
	2	0.929	0.827	0.875	10994	2	0.969	0.965	0.967	9781	2	0.000	0.000	0.000	9781
	3	0.997	0.929	0.962	65188	3	0.998	0.998	0.998	60768	3	0.544	0.957	0.694	60768
	4	0.958	0.997	0.978	181660	4	0.998	0.998	0.998	189054	4	0.875	0.821	0.847	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

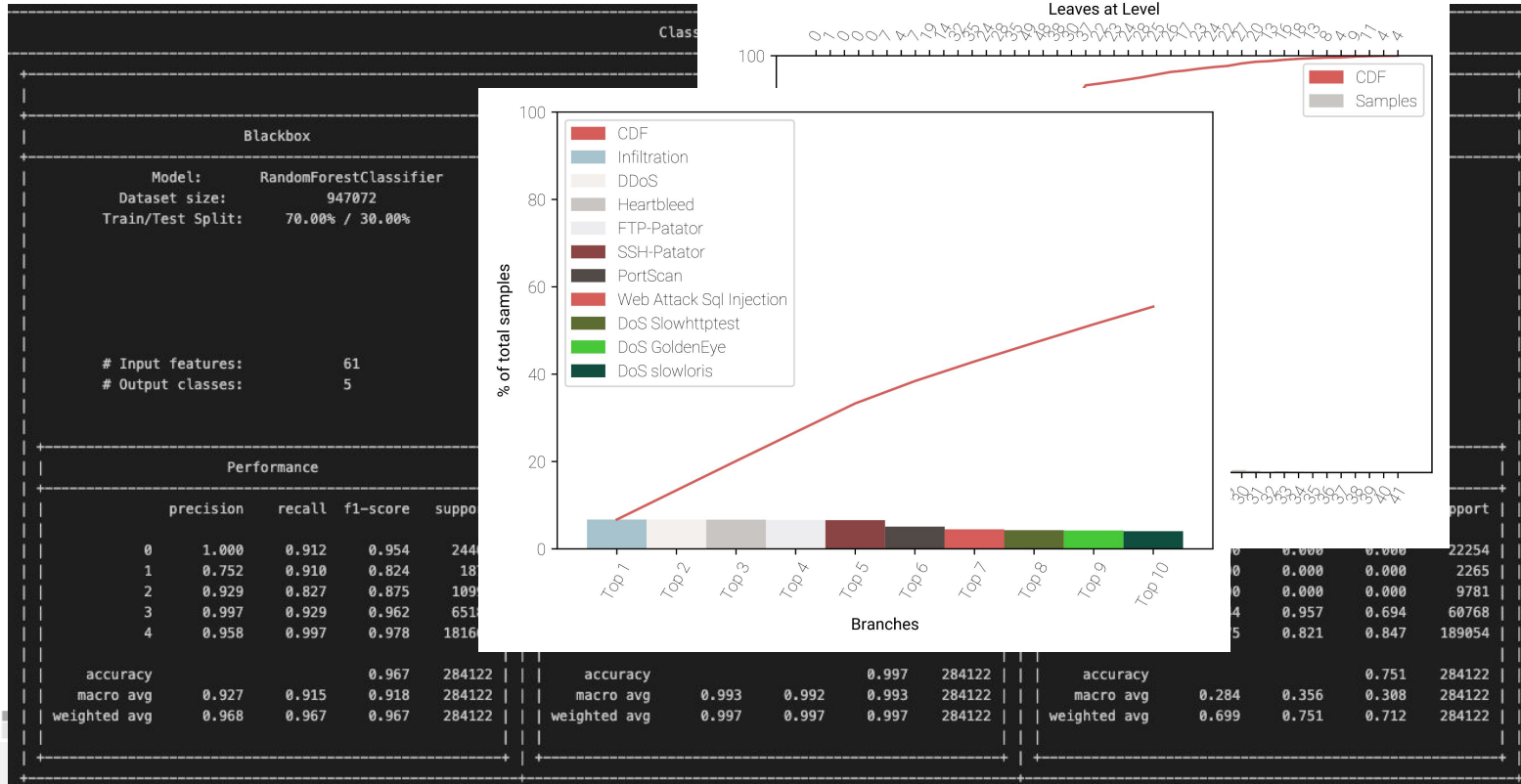
Generating Trust Reports

Classification Trust Report																	
Summary																	
Blackbox				Whitebox				Top-k Whitebox									
Model:		RandomForestClassifier		Explanation method:		Trustee		Explanation method:		Trustee							
Dataset size:		947072		Model:		DecisionTreeClassifier		Model:		DecisionTreeClassifier							
Train/Test Split:		70.00% / 30.00%		Iterations:		1		Iterations:		1							
				Sample size:		50.00%		Sample size:		50.00%							
				Decision Tree Info				Decision Tree Info									
				Size:		2437		Size:		9							
				Depth:		31		Depth:		4							
				Leaves:		1219		Leaves:		5							
# Input features:		61		# Input features:		18 (29.51%)		# Input features:		-							
# Output classes:		5		# 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	22254			
1	0.752	0.910	0.824	1872	1	1.000	1.000	1.000	2265	1	0.000	0.000	0.000	2265			
2	0.929	0.827	0.875	10994	2	0.969	0.965	0.967	9781	2	0.000	0.000	0.000	9781			
3	0.997	0.929	0.962	65188	3	0.998	0.998	0.998	60768	3	0.544	0.957	0.694	60768			
4	0.958	0.997	0.978	181660	4	0.998	0.998	0.998	189054	4	0.875	0.821	0.847	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

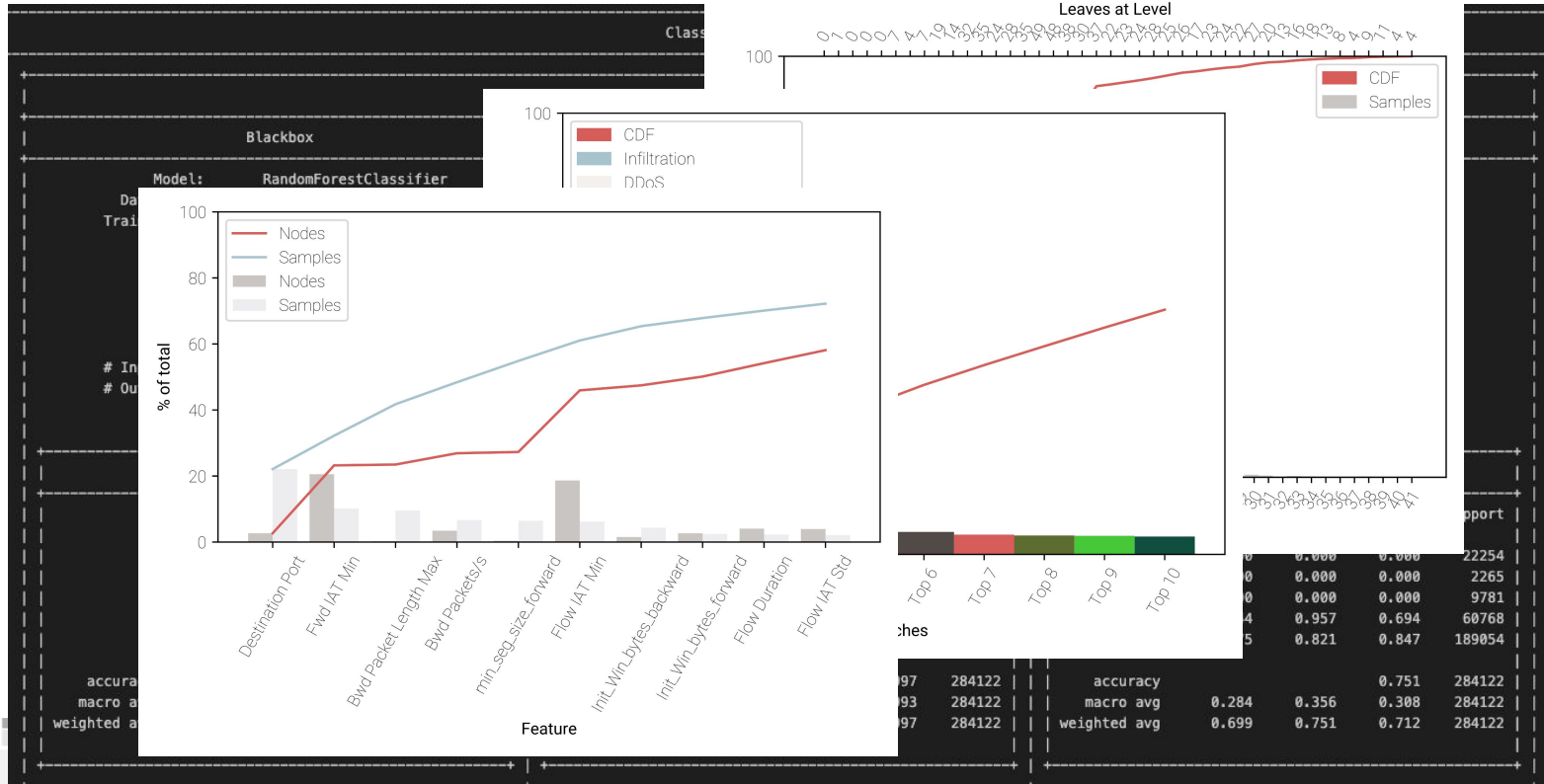
Generating Trust Reports



Generating Trust Reports



Generating Trust Reports



Use Cases

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

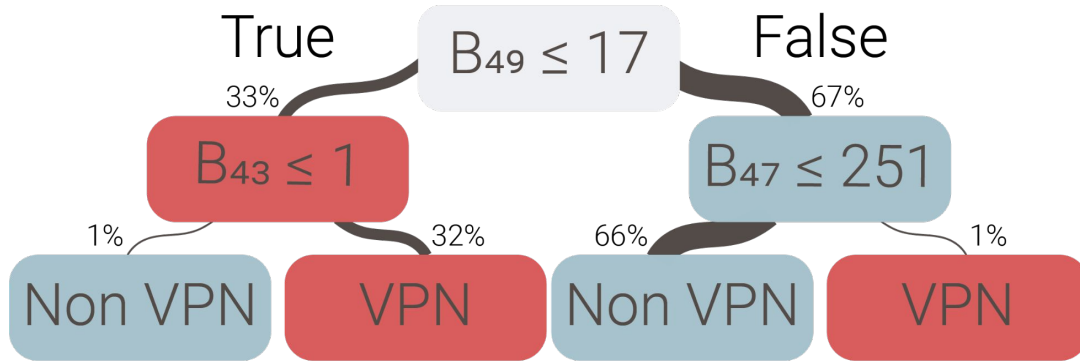
Problem Setup

- **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

Explanation

Non VPN

	0																9	10						19
Pcap	0	161	178	195	212	0	2	0	4	0	0	0	0	0	0	0	0	0	0	255	255			
Meta	20	0	0	0	1	85	65	10	69	0	5	80	24	0	0	0	64	0	0	0	64			
Eth	40	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	228			
IPv4	60	0	0	1	17	34	185	131	202	240	87	224	0	0	252	201	86	20	235	0	...			

VPN

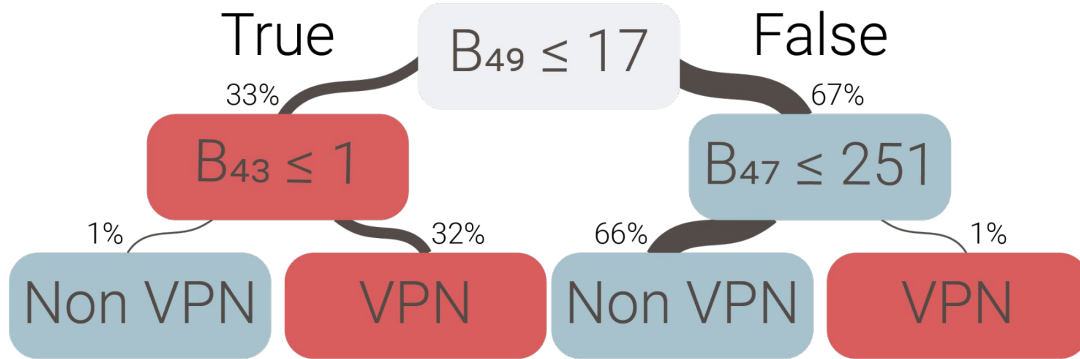
	0																9	10						19
Pcap	0	161	178	195	212	0	2	0	4	0	0	0	0	0	0	0	0	0	0	255	255			
Meta	20	0	0	0	101	85	45	101	91	0	0	111	11	0	0	0	56	0	0	0	56			
IPv4	40	Total Length			Frag. Off.			Protocol																
		69	0	0	56	99	213	64	0	17	5	254	10	8	0	10	69	171	255	36				
UDP	60	146	214	13	150	0	36	120	43	0	1	0	8	33	18	164	66	52	167	9	...			

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

Explanation

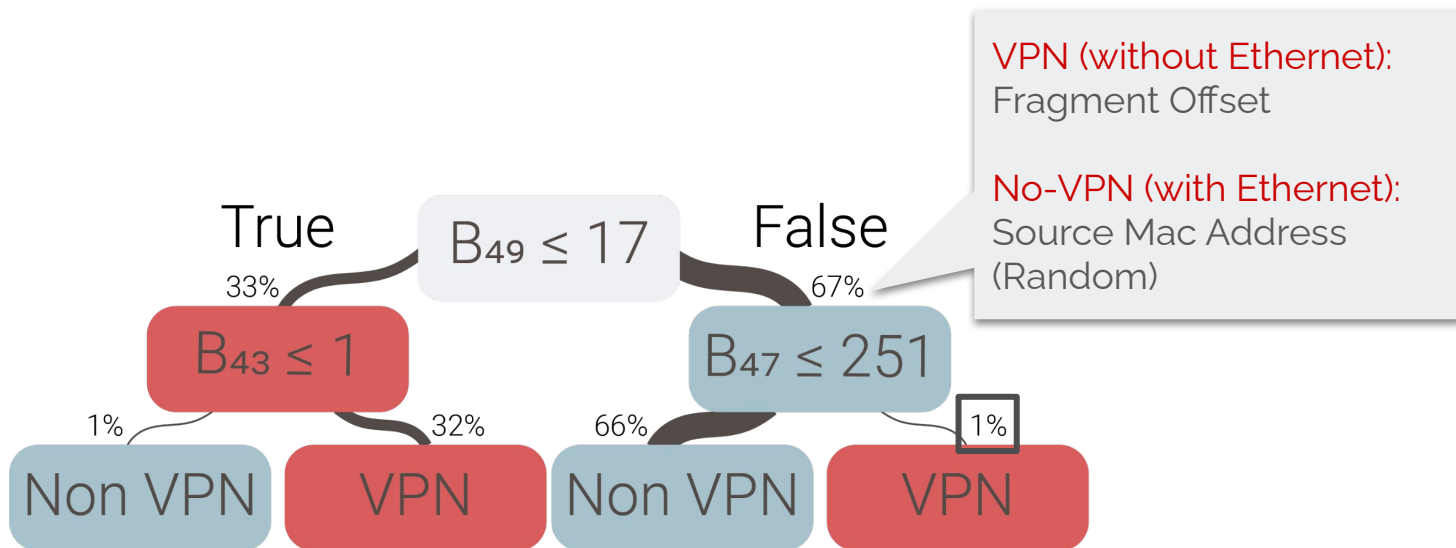
VPN (without Ethernet):
IPv4 Protocol (6 or 17)

No-VPN (with Ethernet):
Source Mac Address
(Random)



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

Explanation

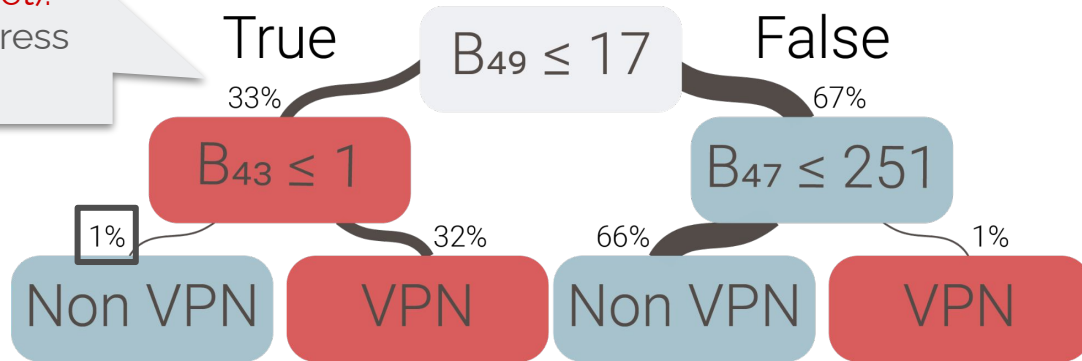


Use Case #1: Detecting VPN vs Non-VPN Traffic

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

	0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19
Pcap Meta	161	178	195	212	229	246	263	280	297	314	331	348	365	382	399	416	433	450	467	484
Ethernet	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

Byte 23: PCAP Link Type

VPN (Without Ethernet): 101

	0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19
Pcap Meta	161	178	195	212	229	246	263	280	297	314	331	348	365	382	399	416	433	450	467	484
IPv4	69	0	0	56	199	213	64	0	64	17	35	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

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!

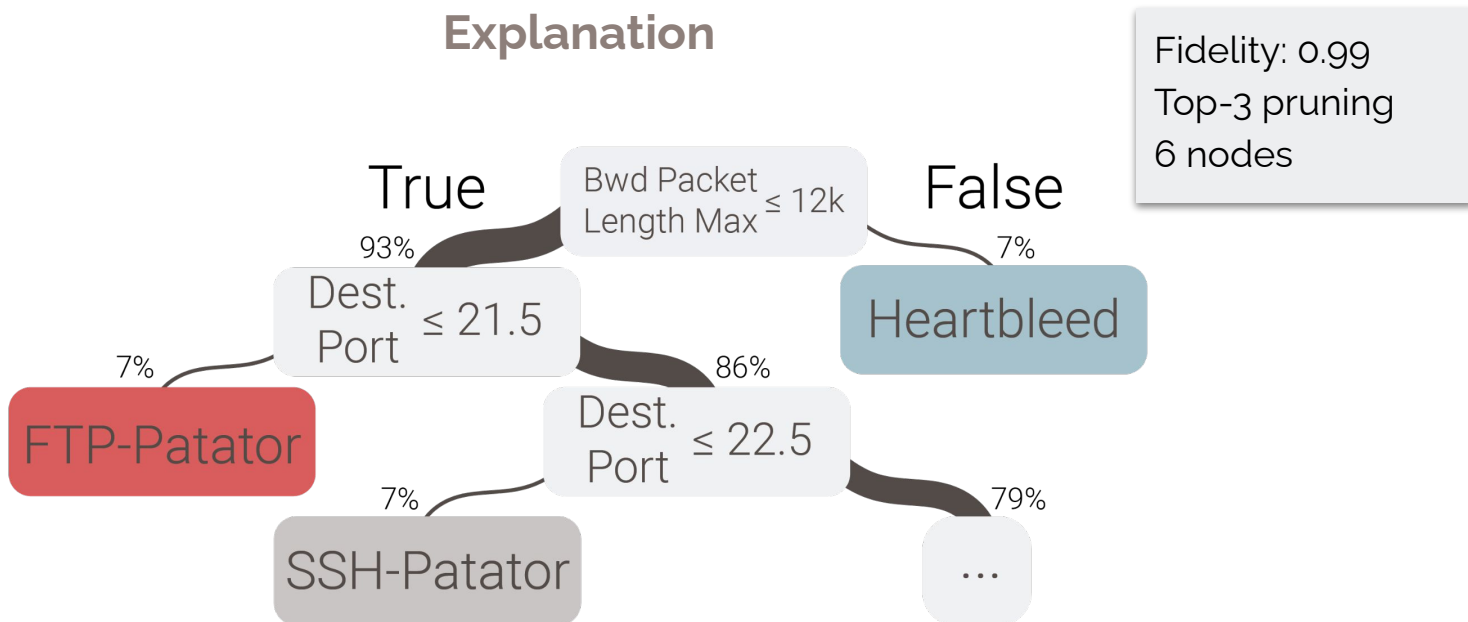
Use Case #2: Detecting Heartbleed Traffic

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

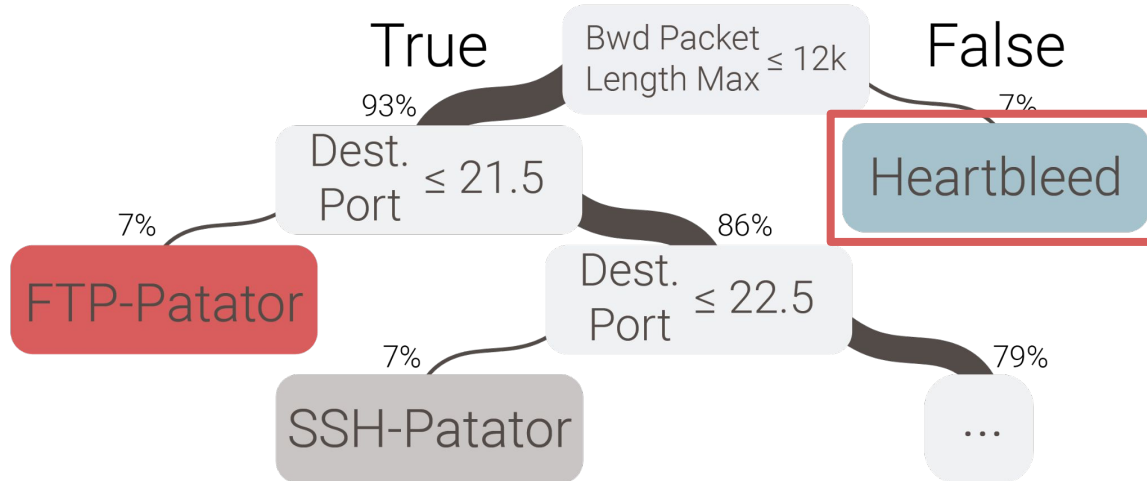
Use Case #2: Detecting Heartbleed Traffic

Explanation



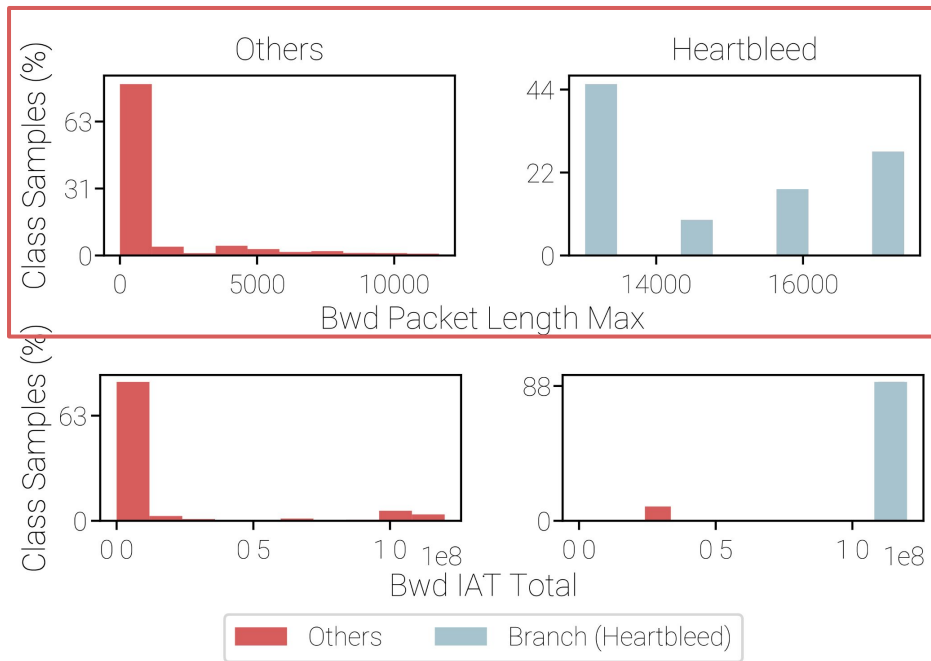
Use Case #2: Detecting Heartbleed Traffic

Explanation



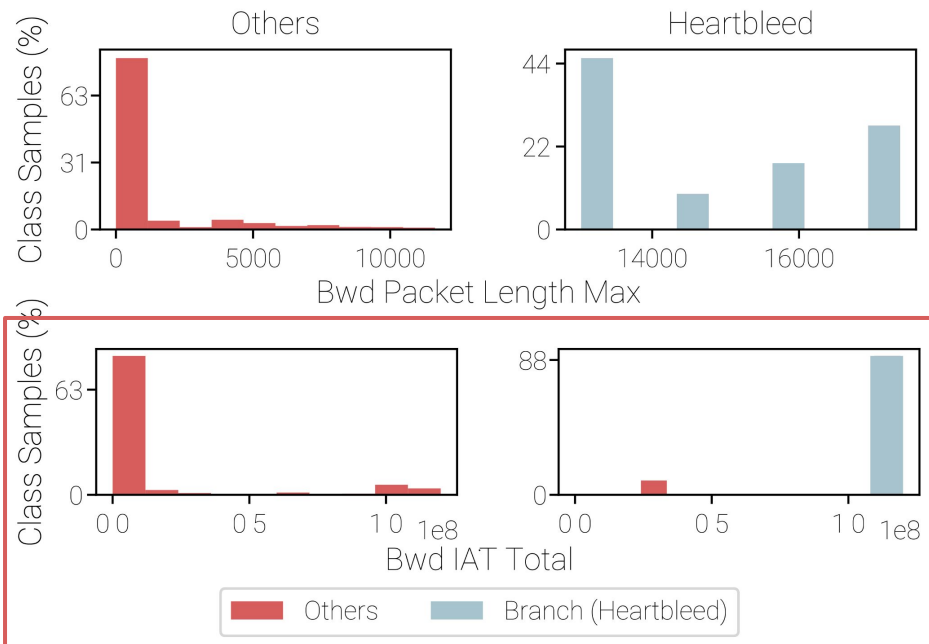
Use Case #2: Detecting Heartbleed Traffic

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!

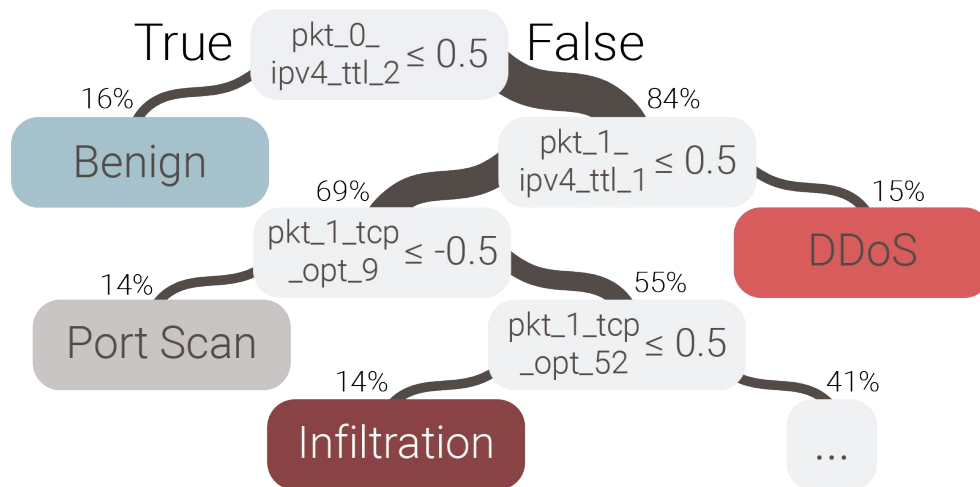
Use Case #3: Inferring Malicious Traffic for IDS

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)
 - **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

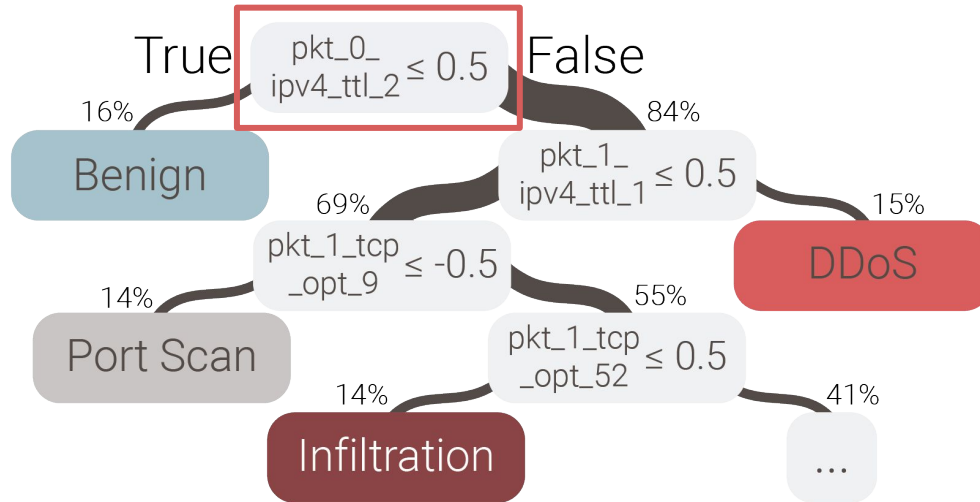
Explanation



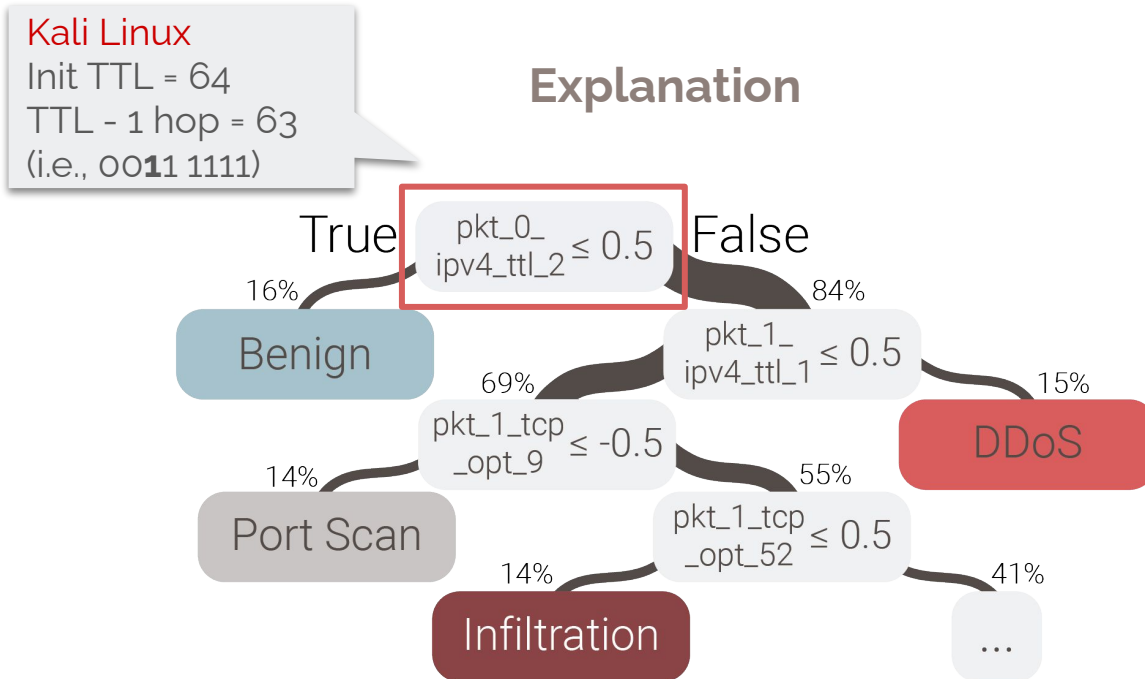
Fidelity: 0.99
Top-4 pruning
8 nodes

Use Case #3: Inferring Malicious Traffic for IDS

Explanation

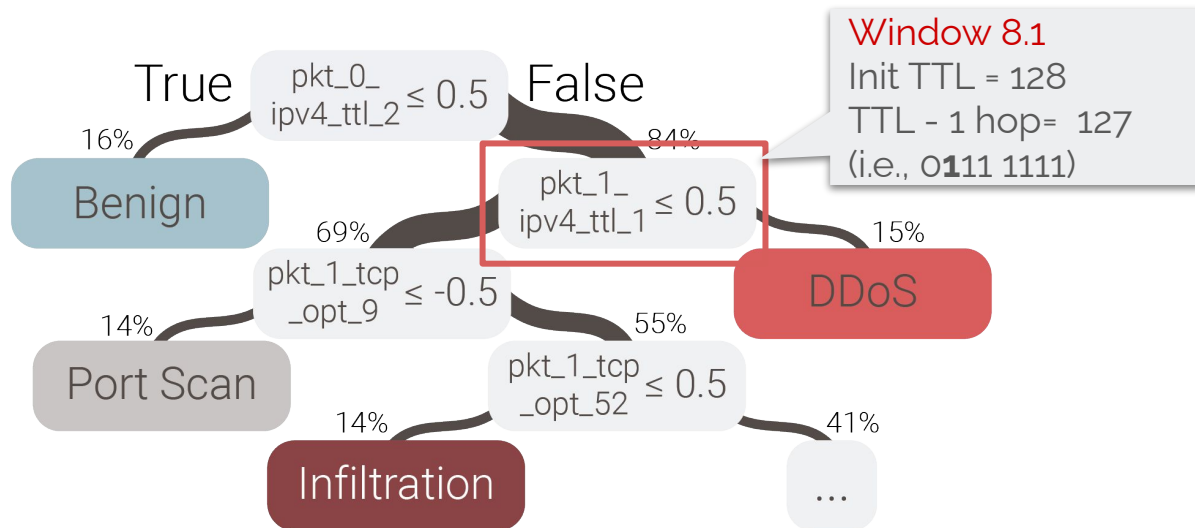


Use Case #3: Inferring Malicious Traffic for IDS



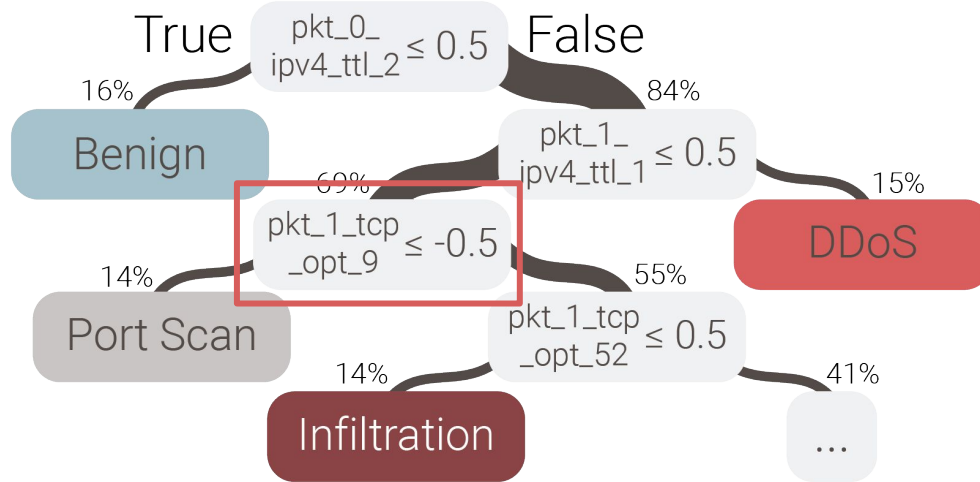
Use Case #3: Inferring Malicious Traffic for IDS

Explanation



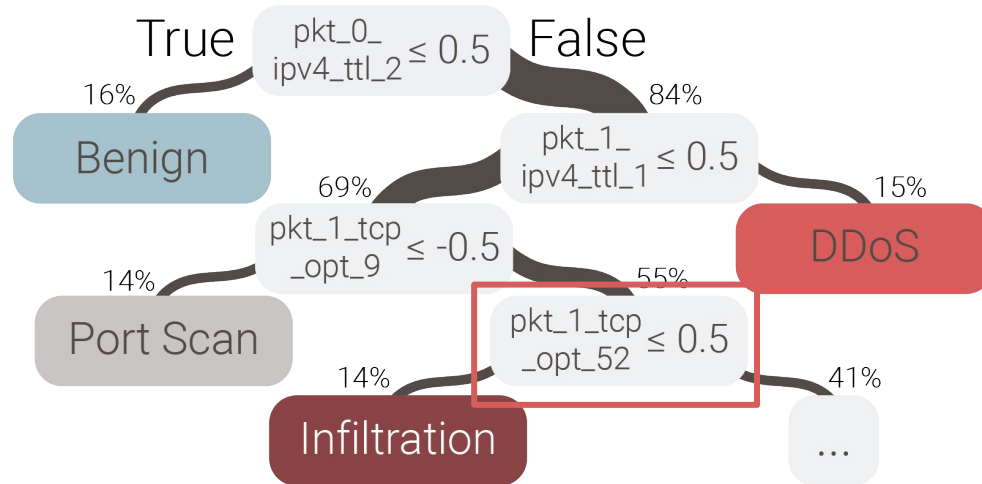
Use Case #3: Inferring Malicious Traffic for IDS

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)	lisy	UNSW-IoT	0.99	Shortcut learning
Adaptive Bit-rate (Mao <i>et al.</i> , SIGCOMM'17)	Pensieve	HSDPA Norway	0.99	O.O.D

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)	lisy	UNSW-IoT	0.99	Shortcut learning
Adaptive Bit-rate (Mao <i>et al.</i> , SIGCOMM'17)	Pensieve	HSDPA Norway	0.99	O.O.D

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)	lisy	UNSW-IoT	0.99	Shortcut learning
Adaptive Bit-rate (Mao <i>et al.</i> , SIGCOMM'17)	Pensieve	HSDPA Norway	0.99	O.O.D

Additional details (see paper)

Algorithmic Description of Trustee

Algorithm 1 Model agnostic decision tree explanation extraction.

```
1: procedure TRUSTEE(  
     $\pi^*$ : Black-box model,  
     $\mathcal{D}_0$ : Initial training dataset,  
     $M$ : Number of samples to train the decision tree,  
     $N$ : Number of iterations of inner loop,  
     $S$ : Number of iterations of outer loop,  
     $k$ : Parameter for Top- $k$  Pruning),  
2: Initialize dataset using black-box  $\mathcal{D} \leftarrow \pi^*(\forall x \in \mathcal{D}_0)$   
3: Initialize stabilization set of DTs  $\mathcal{R} \leftarrow \emptyset$   
4: for  $i \leftarrow 1 \dots S$  do  
5:     for  $j \leftarrow 1 \dots N$  do  
6:         Sample  $M$  training cases uniformly from  $\mathcal{D}$   
            $\mathcal{D}' \leftarrow \{(x, y) \stackrel{\text{i.i.d.}}{\sim} U(\mathcal{D})\}$   
7:         Split sampled dataset for training and testing  
            $\mathcal{D}'_{train}, \mathcal{D}'_{test} \leftarrow \text{TRAINTESTSPLIT}(\mathcal{D}')$   
8:         Train DT  
            $\hat{\pi}_j \leftarrow \text{TRAINDECISIONTREE}(\mathcal{D}'_{train})$   
9:         Test and get samples DT misclassifies  
            $\mathcal{D}'_e \leftarrow \{\forall (x, y) \in \mathcal{D}'_{test} \mid \hat{\pi}_j(x) \neq \pi^*(x)\}$   
10:        Get correct outcome from black-box  
            $\mathcal{D} \leftarrow \pi^*(\forall x \in \mathcal{D})$ 
```

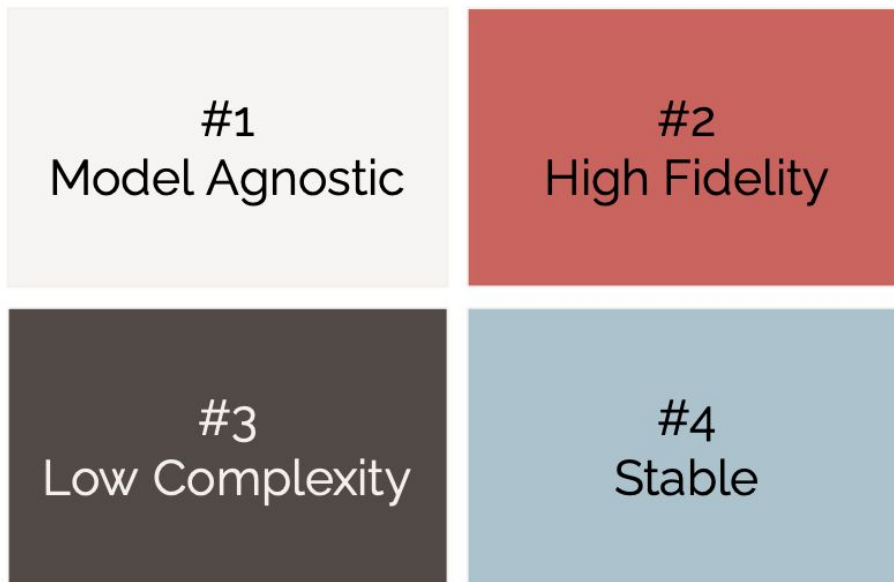
Additional details (see paper)

Algorithmic Description of Trustee

Algorithm 1 Model agnostic decision

```
1: procedure TRUSTEE(  
     $\pi^*$ : Black-box model,  
     $\mathcal{D}_0$ : Initial training dataset,  
     $M$ : Number of samples to train,  
     $N$ : Number of iterations of inner loop,  
     $S$ : Number of iterations of outer loop,  
     $k$ : Parameter for Top- $k$  Pruning  
2: Initialize dataset using black-box model  
3: Initialize stabilization set of iterations  
4: for  $i \leftarrow 1 \dots S$  do  
5:   for  $j \leftarrow 1 \dots N$  do  
6:     Sample  $M$  training samples  
        $\mathcal{D}' \leftarrow \{(x, y)^i\}$   
7:     Split sampled dataset  
        $\mathcal{D}'_{train}, \mathcal{D}'_{test} \leftarrow \dots$   
8:     Train DT  
        $\hat{\pi}_j \leftarrow \text{TRAINDEC}(\dots)$   
9:     Test and get samples  
        $\mathcal{D}'_e \leftarrow \{\forall(x, y) \in \mathcal{D}'_{test}\}$   
10:    Get correct outcome for  $\mathcal{D}'_e$ 
```

Ablation Study on Design Requirements



Additional details (see paper)

Algorithmic Description of Trustee

Algorithm 1 Model agnostic decision

```
1: procedure TRUSTEE(  
    $\pi^*$ : Black-box model,  
    $\mathcal{D}_0$ : Initial training dataset,  
    $M$ : Number of samples to train,  
    $N$ : Number of iterations of initialization,  
    $S$ : Number of iterations of optimization,  
    $k$ : Parameter for Top- $k$  Pruning  
2: Initialize dataset using black-box model  
3: Initialize stabilization set of  $I$   
4: for  $i \leftarrow 1 \dots S$  do  
5:   for  $j \leftarrow 1 \dots N$  do  
6:     Sample  $M$  training cases  $\mathcal{D}' \leftarrow \{(x, y)\}^i$   
7:     Split sampled dataset into  $\mathcal{D}'_{train}, \mathcal{D}'_{test}$   
8:     Train DT  $\hat{\pi}_j \leftarrow \text{TRAINDEC}(\mathcal{D}'_{train}, \mathcal{D}'_{test})$   
9:     Test and get samples  $\mathcal{D}'_e \leftarrow \{\forall(x, y) \in \mathcal{D}'_{test} : \hat{\pi}_j(x) \neq y\}$   
10:    Get correct outcome  $\mathcal{D}'_c \leftarrow \{(x, y) \in \mathcal{D}'_{train} : \hat{\pi}_j(x) = y\}$ 
```

Ablation Study on Design Requirements

#1
Model Agnostic

#3
Low Complexity

Trust Report and User Guide

Classification Trust Report													
Summary													
Model		Whitebox				Blackbox							
DecisionTreeClassifier	947072	0.00%	30.00%	61	5	Explanation method: Trustee							
						Model: DecisionTreeClassifier		Explanation method: Trustee					
						Iterations: 1		Iterations: 1					
						Sample size: 50.00%		Sample size: 50.00%					
						Decision Tree Info							
						Size: 2437		Size: 2437					
						Depth: 31		Depth: 31					
						Leaves: 1219		Leaves: 1219					
						# Input features: 18 (29.51%)		# Input features: 18 (29.51%)					
						# Output classes: 5 (100.00%)		# Output classes: 5 (100.00%)					
Fidelity													
Whitebox			Blackbox				Whitebox			Blackbox			
call	f1-score	support	precision	recall	f1-score	support	precision	recall	f1-score	support	precision		
0.912	0.954	24408	1.000	1.000	1.000	22254	1.000	1.000	1.000	22254	1.000		
0.910	0.824	1872	1.000	1.000	1.000	2265	1.000	1.000	1.000	2265	1.000		

Trustee Python package



trustee 1.1.1

pip install trustee

Released: Aug 28, 2022

This package implements the Trustee framework to extract decision tree explanation from black-box ML models.

Navigation

Project description

Release history

Download files

Project links

Homepage

Repository

Statistics

GitHub statistics:

Stars: 2

Forks: 0

Open issues/PRs: 0

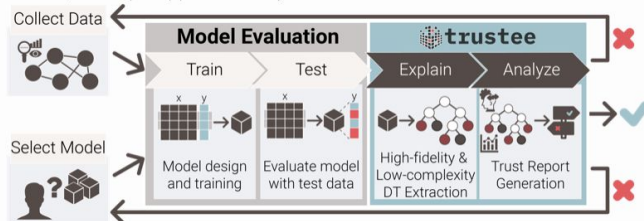
View statistics for this project via

Project description



This package implements the `trustee` framework to extract decision tree explanation from black-box ML models. For more information, please visit the [documentation website](#).

Standard AI/ML development pipeline extended by Trustee.



Downloads 43

Daily

Downloads 175

Weekly

Downloads 6k

Total

GitHub Repo Use Cases Tech Report



Trustee 1.1.1 documentation

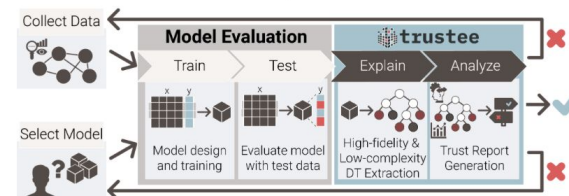
Search

API Examples

TrustReport for Classification
TrustReport for Regression
RegressionTrustee
ClassificationTrustee

Overview

Trustee is a framework to extract decision tree explanation from black-box ML models.



Standard AI/ML development pipeline extended by Trustee.

Getting Started

This section contains basic information and instructions to get started with Trustee.

Python Version

Trustee supports Python ≥ 3.7 .

Install Trustee

Use the following command to install Trustee:

```
$ pip install trustee
```

Sample Code

```
from sklearn import datasets
```

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!
 - Just download our Python package

Thank you!

Arthur Jacobs
asjacobs@inf.ufrgs.br



<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>

Backup

Existing approaches

Method	Model Agnostic	High Fidelity	Domain-specific Pruning
Trepan	✓	—	—
<i>dtextract</i>	✓	—	—
VIPER	—	—	—
Metis	—	—	—
 trustee	✓	✓	✓

Use Case #4: Anomaly Detection for Mirai Attacks

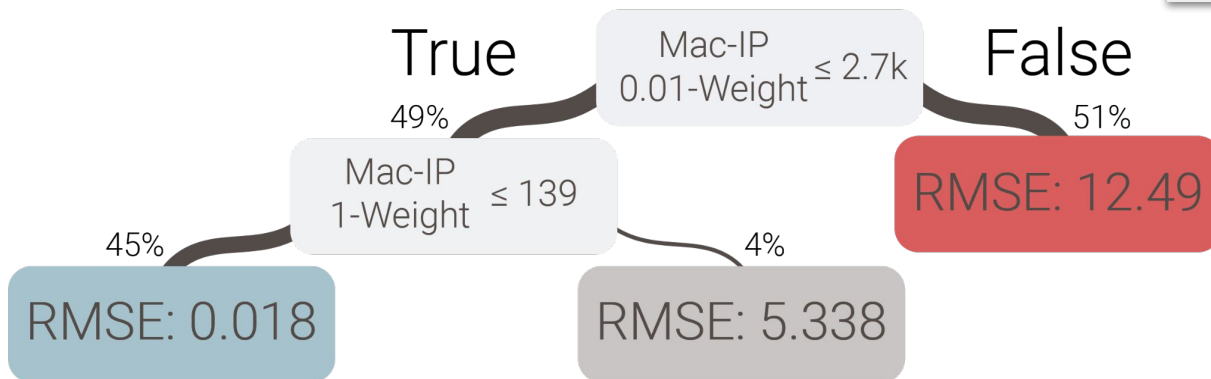
Problem Setup

- **Selected publications:**
 - *"Kitsune: An Ensemble of Autoencoders for Online Network Intrusion Detection"* — Mirsky et al., 2018
- **Proposal:**
 - **Model:** Kitsune, an ensemble of neural networks, trained with unsupervised learning, for anomaly detection
 - **Features:** 110 features based on traffic statistics (e.g., number of packets per **time window**).
 - **Dataset:** synthetic Mirai attack trace.
- **Results:**
 - Reported R-squared: 0.99
 - Reproduced R-squared: 0.99

Use Case #4: Anomaly Detection for Mirai Attacks

Explanation

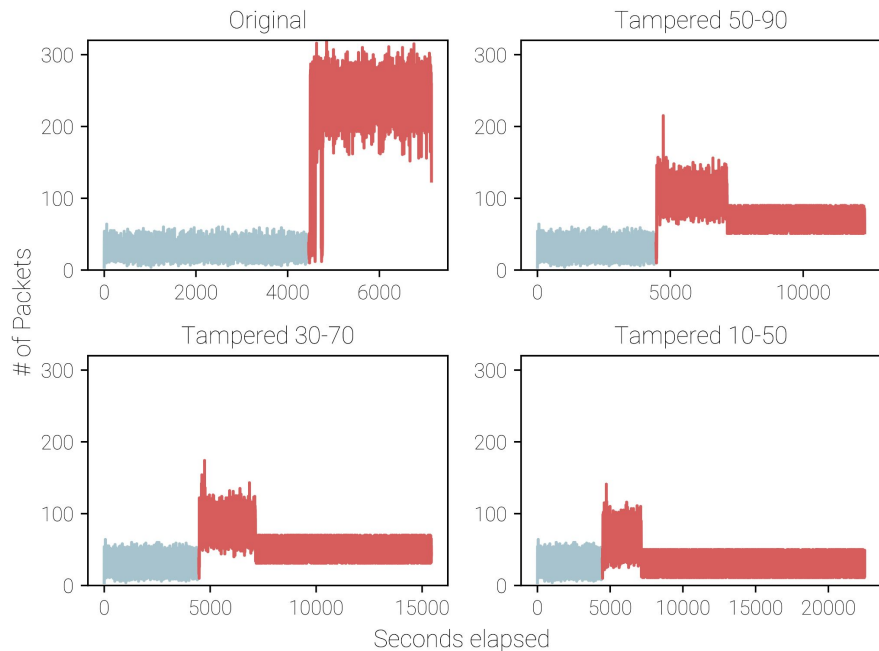
Fidelity: 0.99
Top-3 pruning
5 nodes



Use Case #4: Anomaly Detection for Mirai Attacks

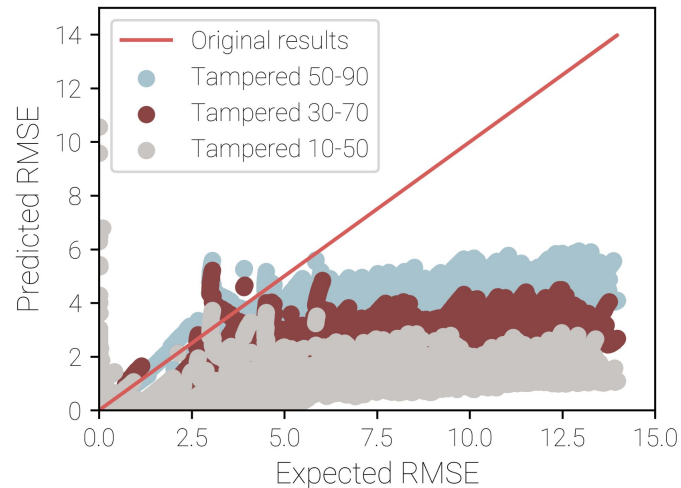
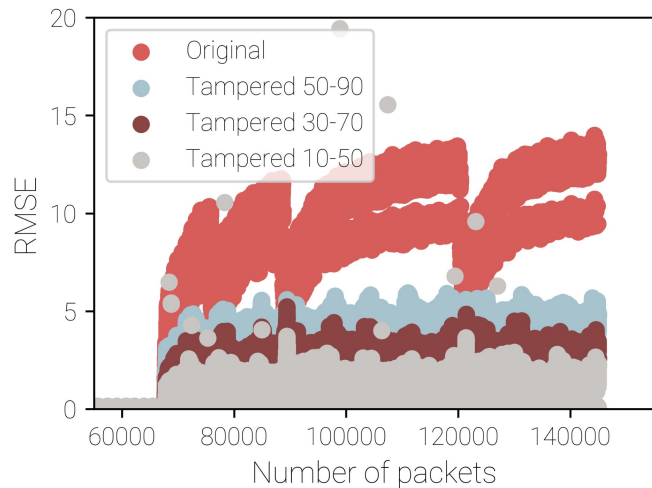
Validation

- Validation datasets:



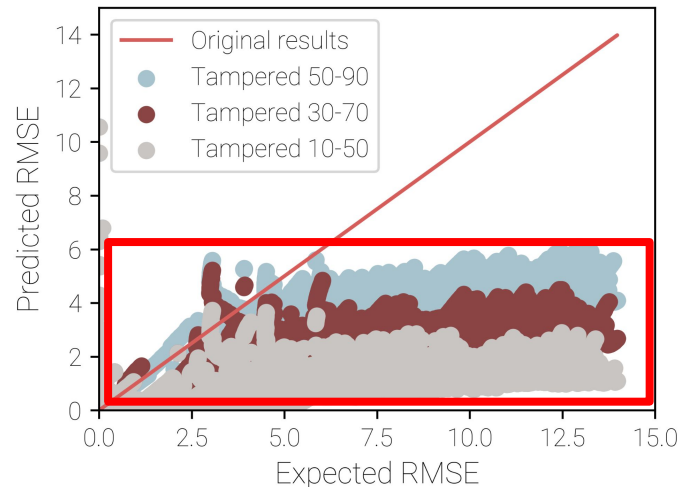
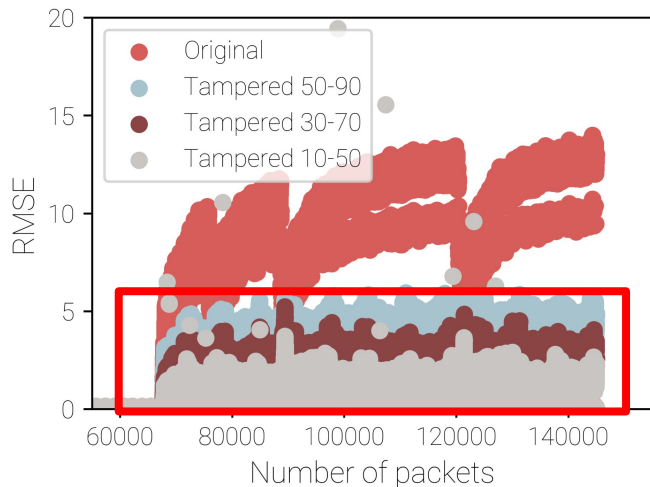
Use Case #4: Anomaly Detection for Mirai Attacks

Validation



Use Case #4: Anomaly Detection for Mirai Attacks

Validation



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