

Protecting the Protector:

Hardening Machine Learning
Defenses Against Adversarial Attacks

Jugal Parikh, Senior Data Scientist
Holly Stewart, Principal Research Manager
Randy Treit, Senior Researcher

A dark-themed world map showing the locations of attacks. Numerous yellow and orange dots are scattered across the map, with a higher concentration in Europe, North America, and East Asia. The map includes country borders and is set against a dark background.

In a single day...

2.6 million people

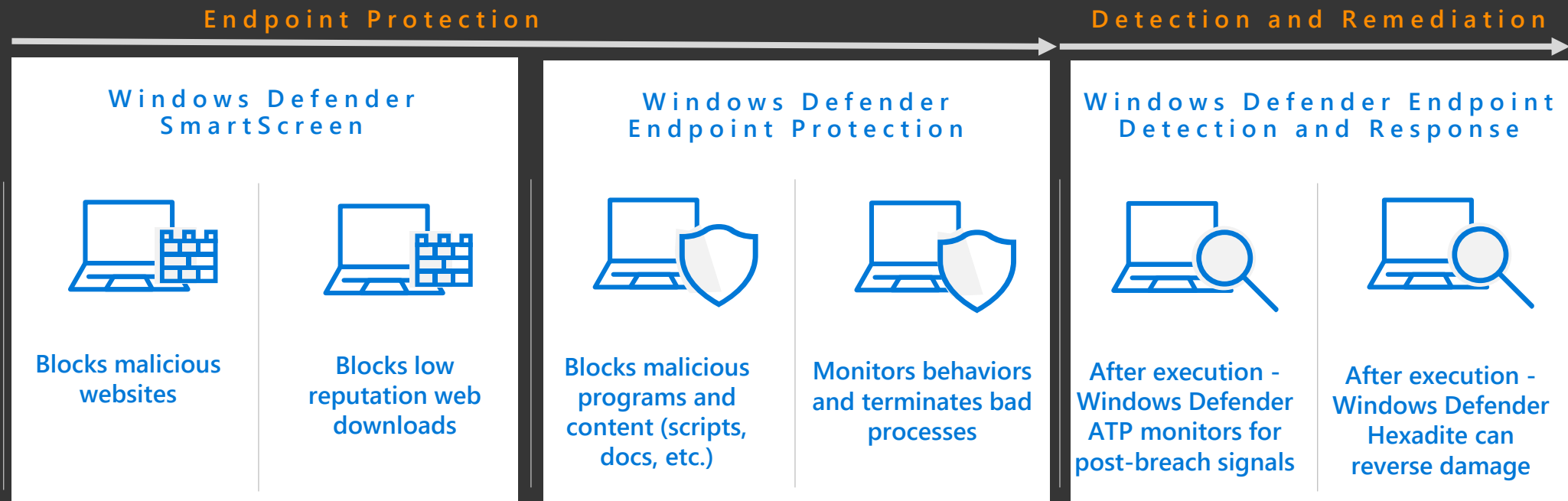
232 country/regions

1.7 million first seen attacks



60% of these attacks
were over within the
hour

Windows Defender Advanced Threat Protection



Holly!



Jugal!

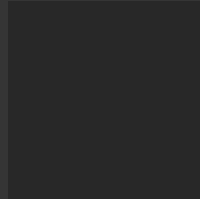
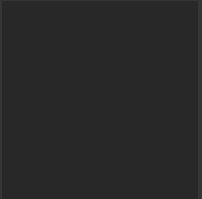
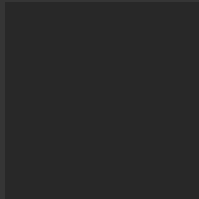
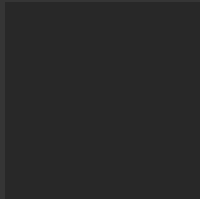
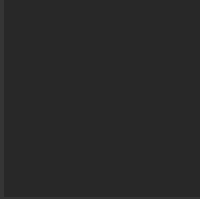


Randy!

Windows Defender ATP Research



Threat Predict Research Team



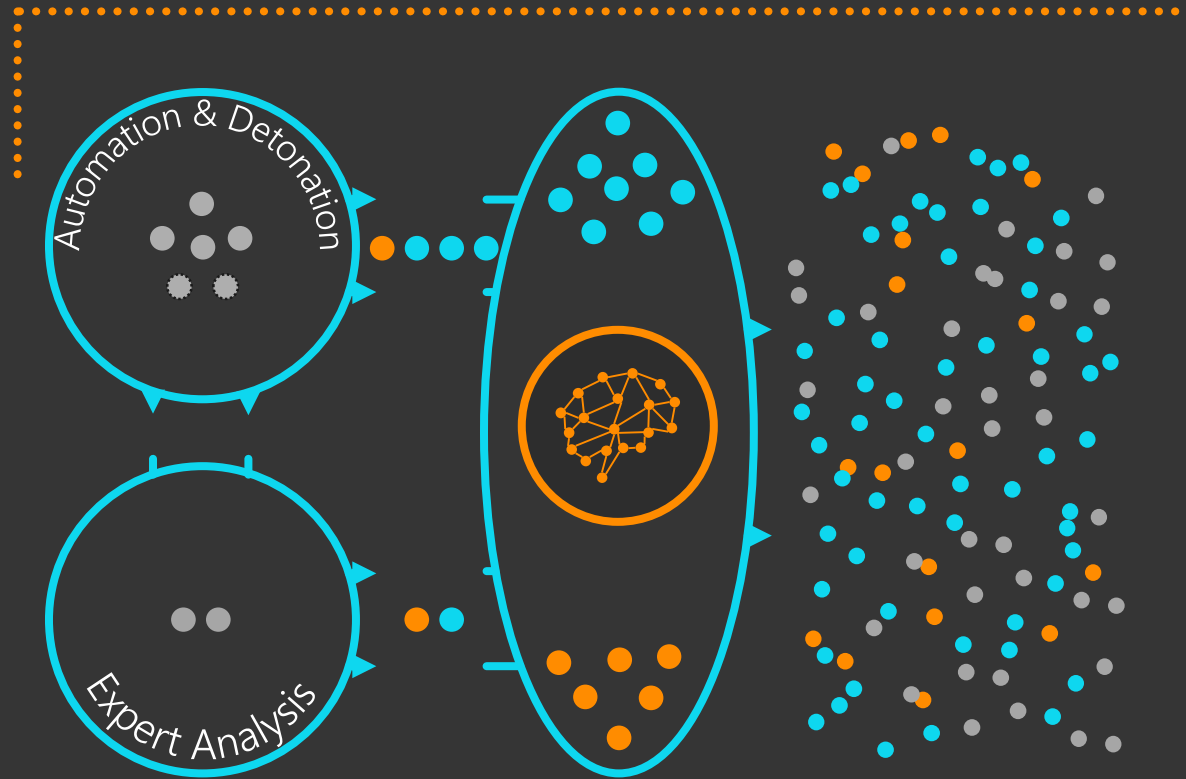
Our focus: Use machine learning
to block attacks for Windows
Defender ATP



Machine Learning Primer

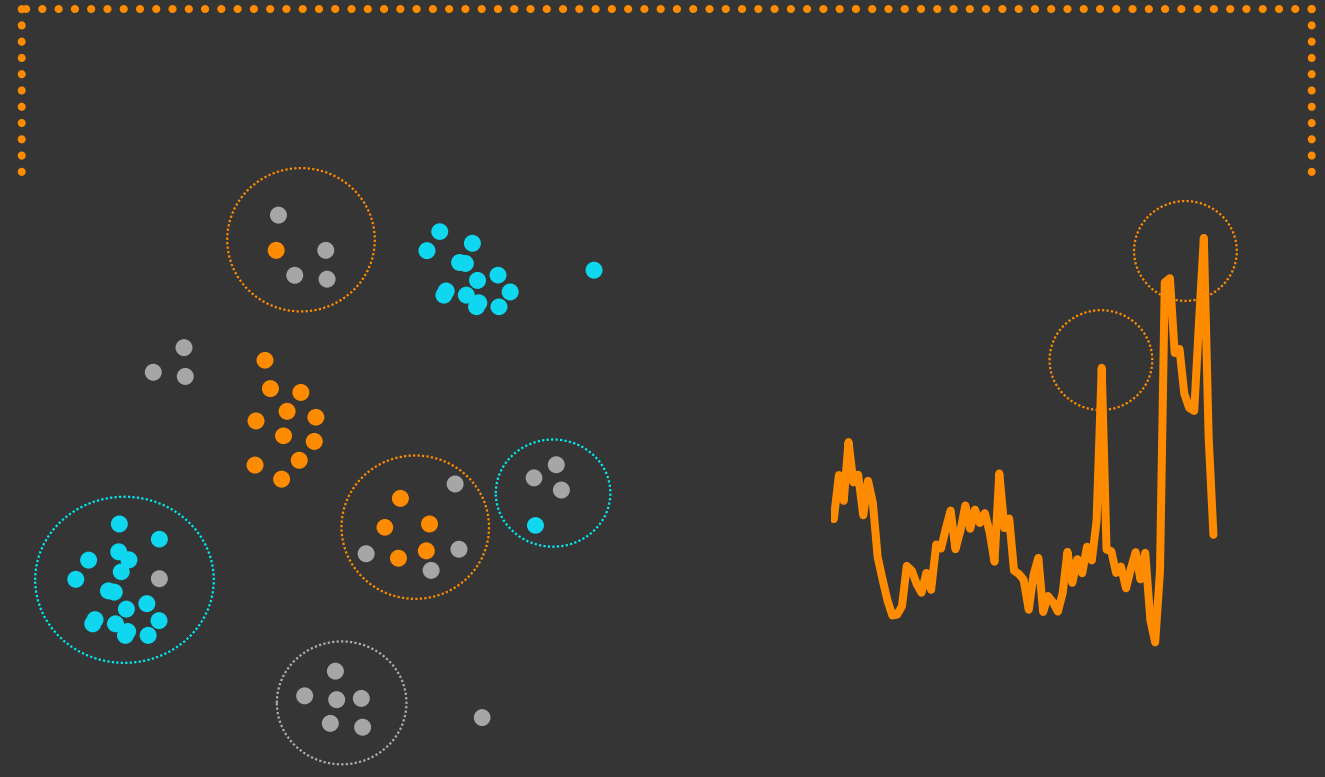
Types of Machine Learning

SUPERVISED



EXPERTS → LABELS → ML → PREDICTIONS

UNSUPERVISED

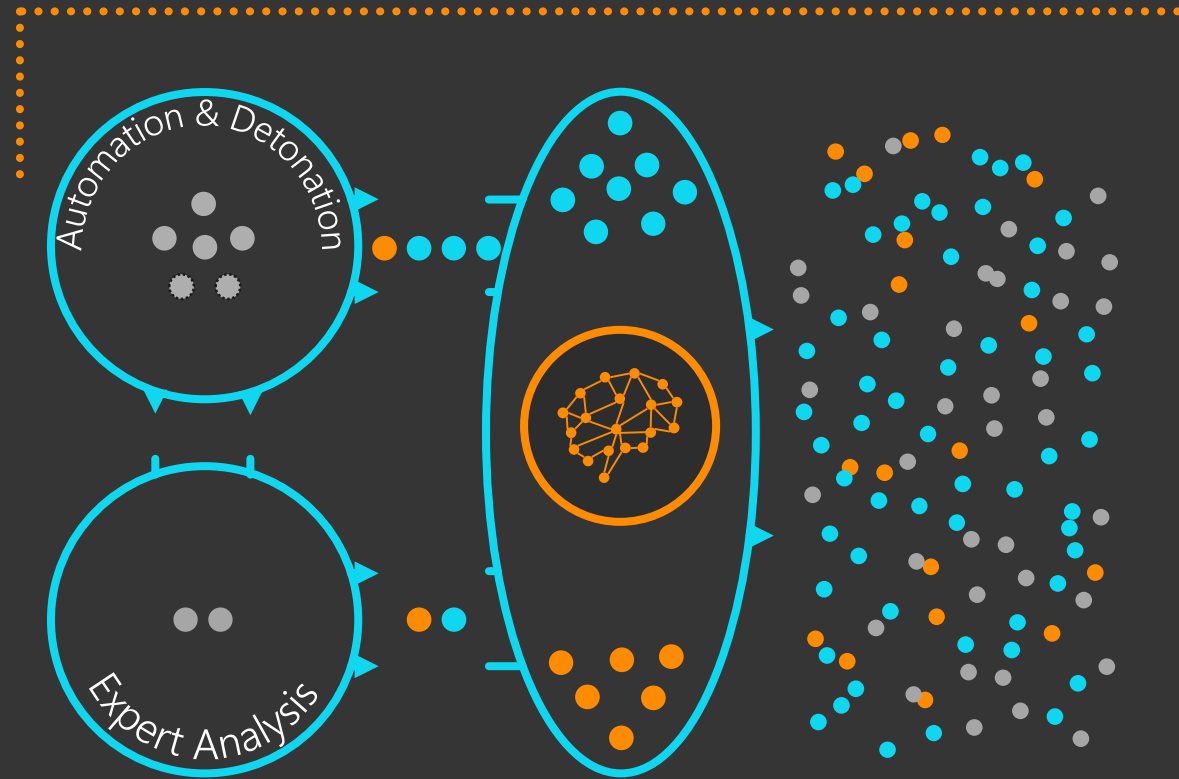


UNKNOWN UNKNOWNNS

ANOMALY DETECTION

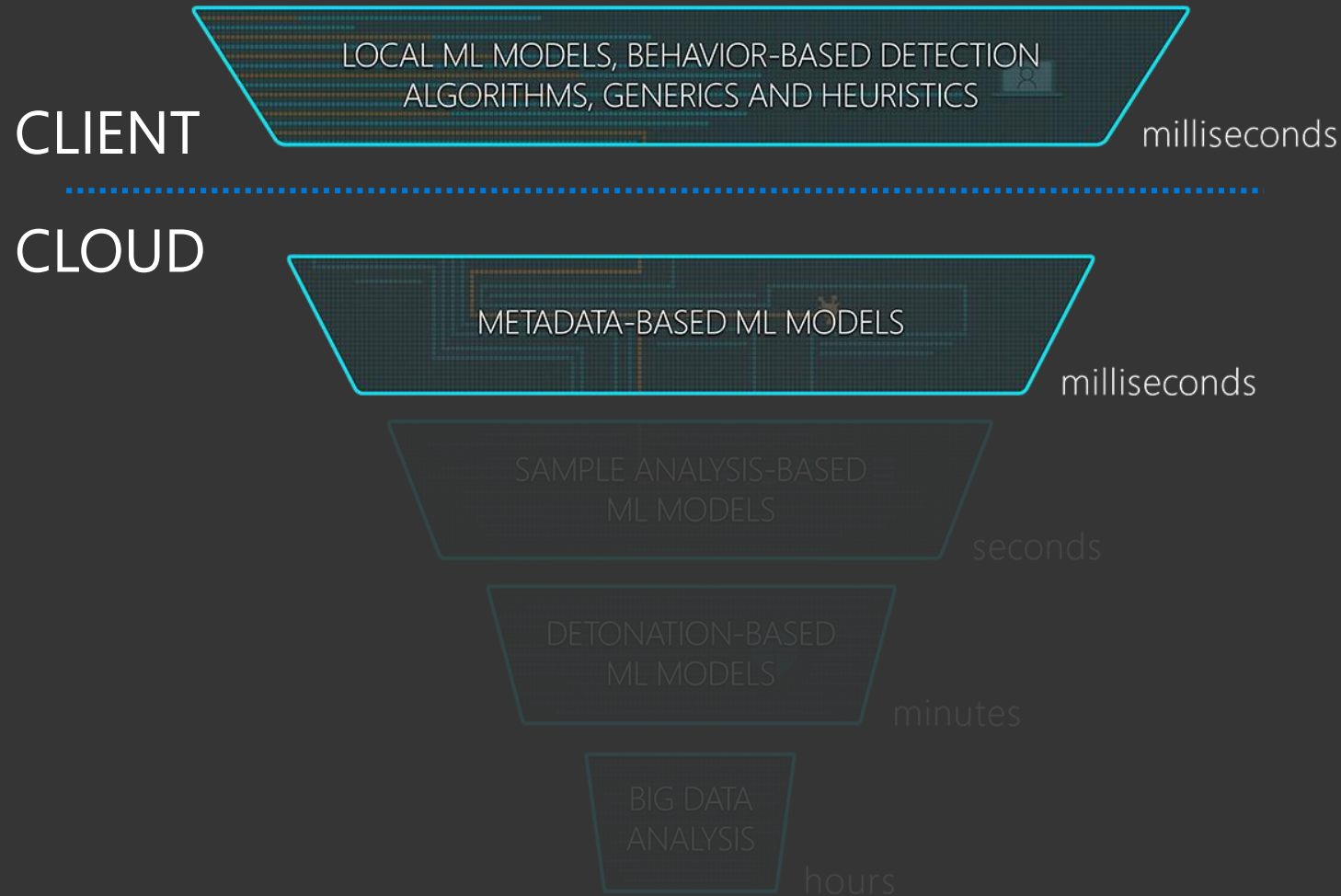
Types of Machine Learning

SUPERVISED

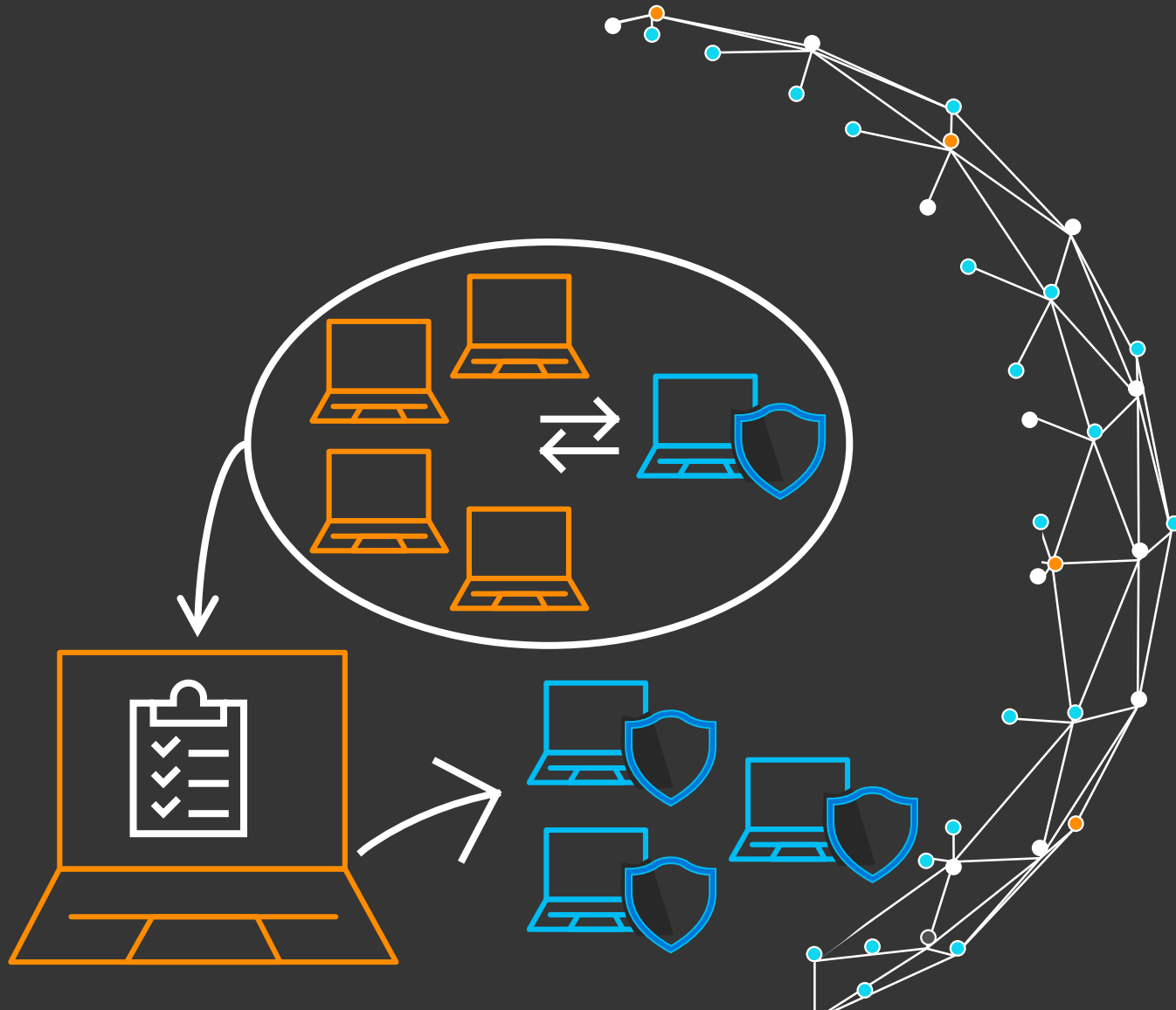


EXPERTS → LABELS → ML → PREDICTIONS

Machine Learning for Endpoint Protection



Client Machine Learning



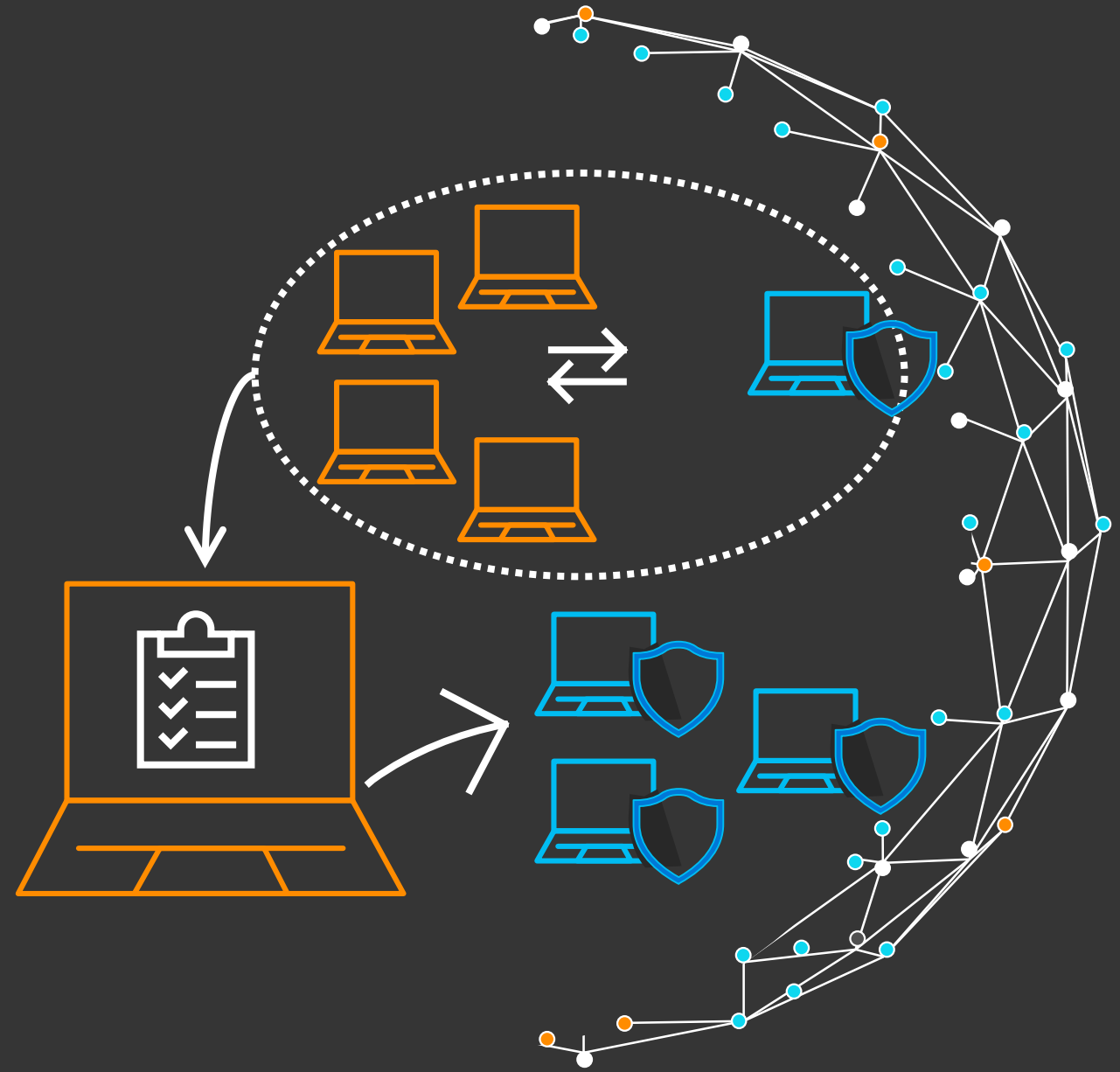
Pro:
Disconnected
protection

Con: Silent
adversarial brute
force attacks

Cloud Machine Learning

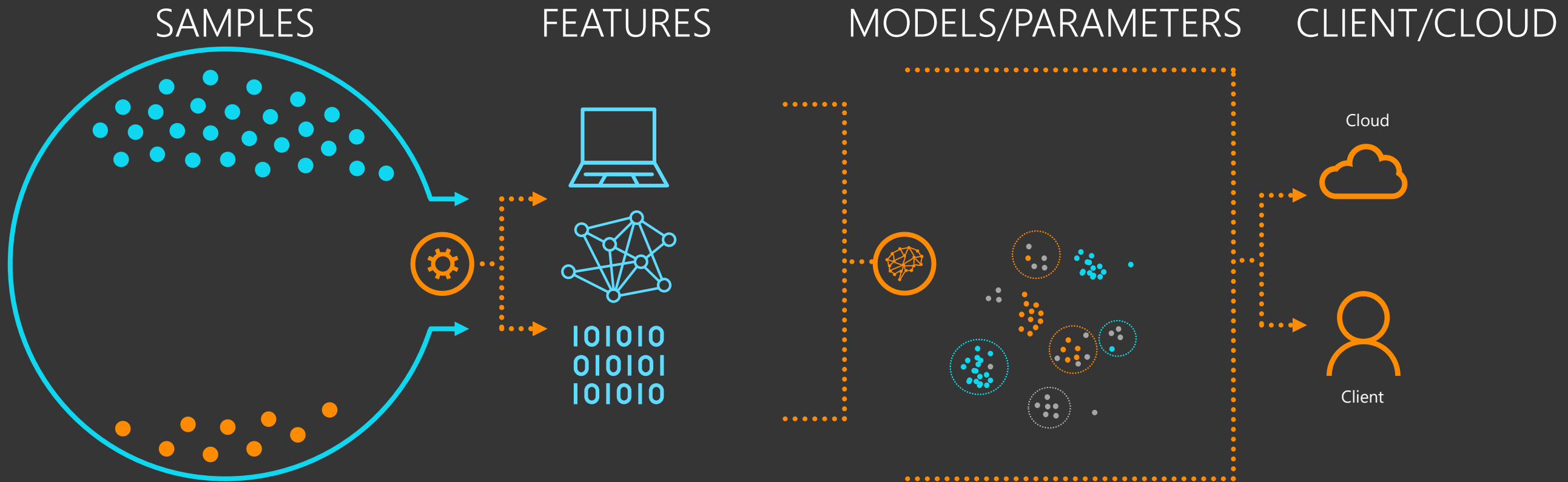
No private brute forcing

Minimal client performance impact



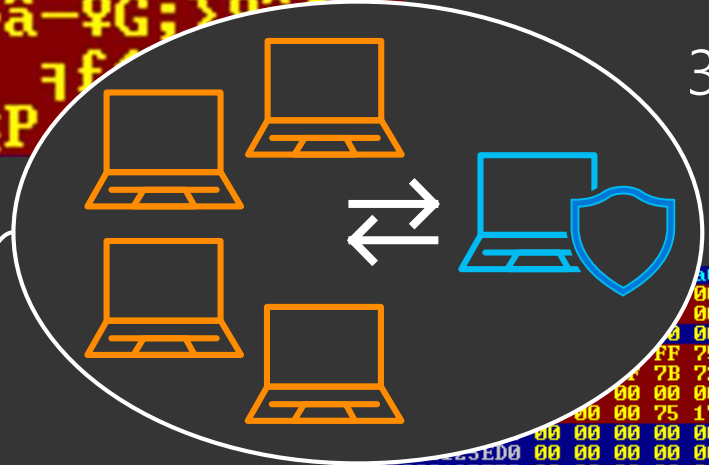
Adversarial ML Examples

Theoretical Attack Vectors: Supervised Model



Specially crafted files → AV industry FPs

1. Identify signature fragments detected as malicious
2. Identify automated detection techniques

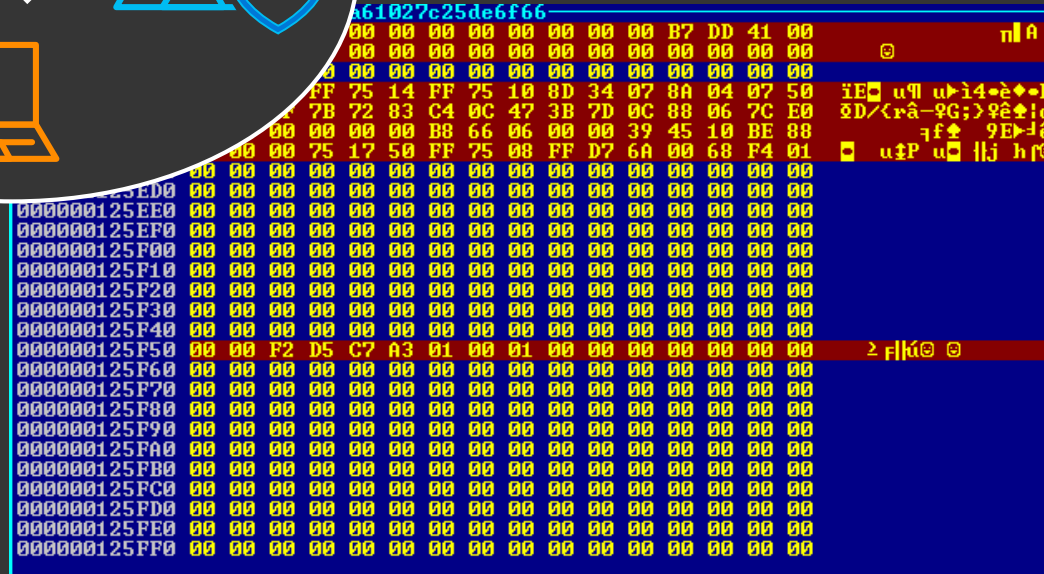


3. Inject signature fragments into clean files

4. Add crafted files to VirusTotal using TOR

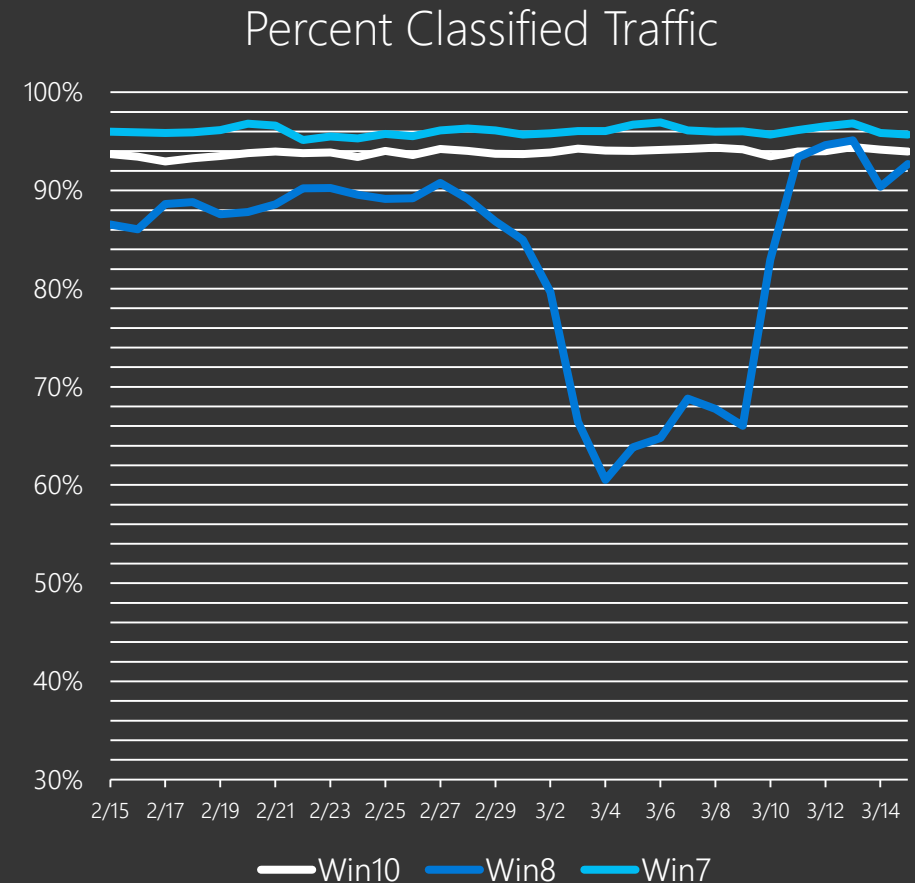


5. Automation signs targeted clean files, multiple vendors have FPs



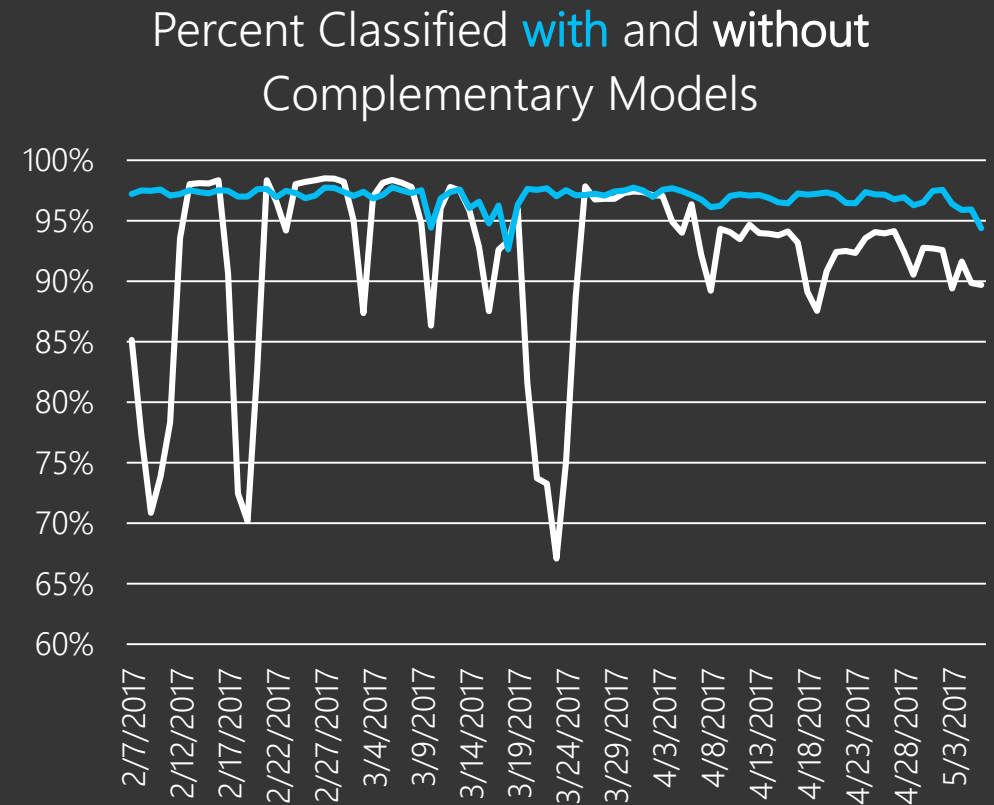
Attacks on Certificate Reputation (Early 2017)

- Synthetic traffic designed to quickly gain reputation on a digital certificate
- Targeted Windows 8
- Originally surfaced as a high percentage of traffic that wasn't classified
- **Low-volume** and **unsigned** file attacks were also identified during investigation



Attacks on Certificate Reputation (cont.)

- Attackers guessed major features (time, traffic, digital certificate)
- Team developed complementary models with additional features that filtered fake traffic out of telemetry
- Combination of models removed attack traffic from training data



Previous research pointed to ensembles

- Research on adversarial attacks against deep learning classifiers
- Showed that an ensemble of classifiers helped defend against the attacks tested in the paper
- See more at:

Attack and Defense of Dynamic Analysis-Based, Adversarial Neural Malware Classification Models

[Jack W. Stokes](#), [De Wang](#), [Mady Marinescu](#), [Marc Marino](#), [Brian Bussone](#)

<https://arxiv.org/abs/1712.05919>

Ensemble Machine Learning Primer

Ensemble Model Development and Testing

Ensemble ML Primer

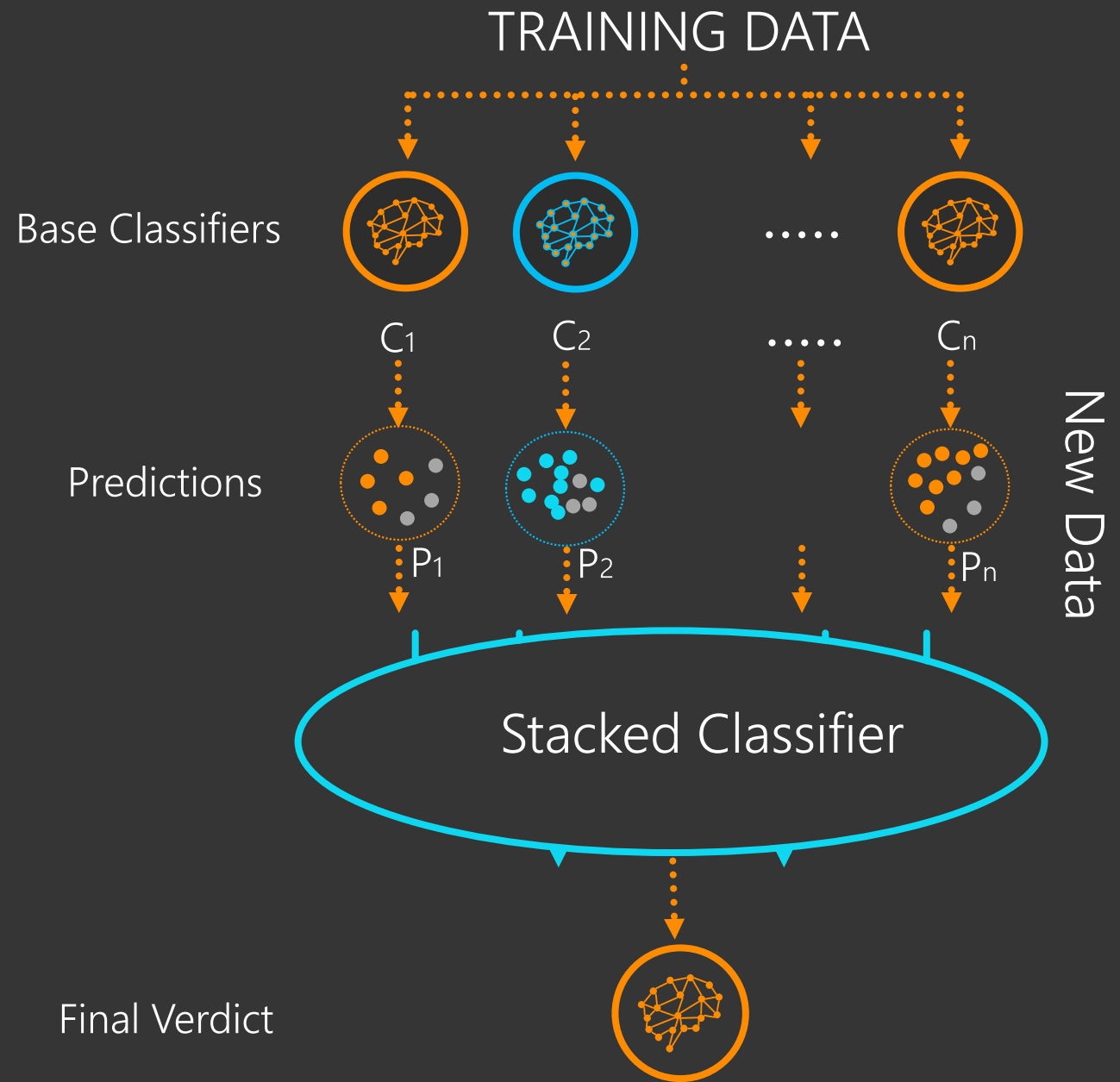


Diversity Requirements

Developing the Model

Testing the Model

Overview



Challenges

Dealing with:

- Active adversarial
- Volatility/ Covariate Shift
- Noisy environment

Scale:

- Petabytes of threat Intelligence daily

Evaluate:

- ~2.3 Billion global queries everyday

Diversity

Diverse Models

1. Different feature sets
2. Different training algorithms
3. Different training data sets
4. Different optimization settings

Features - Highly dimensional data



Machine attributes

OS version
Processor
Security settings



Researcher Expertise

10k+ researcher attributes

100k+ static attributes

10k+ behavioral attributes



Behavioral and contextual attributes

Process and installation

ProcessName
ParentProcess
TriggeringSignal
TriggeringFile
Download IP and URL
Parent/child relationships

Behavioral

Connection IP
System changes
API calls
Process injection

Locale

Locale setting
Geographical location

Static file attributes

Partial and Fuzzy hashes

ClusterHash
ImportHash
Fuzzy hashes

Full File Content

Header
Footer
Raw file content

File properties

File Geometry
FileSize
FileMetaData
....

Signer info

Issuer
Publisher
Signer

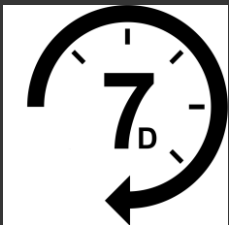
Diverse Set of Classifiers

Client Models	Feature Set	Learner	# of Features	PE
	PE Properties	Fast Tree Ensemble	10K+ features	
Cloud Models	Researcher Expertise	Boosted Tree Ensemble	190K+ features	JavaScript
	Behavioral	Boosted Tree Ensemble	6M+ features	
Full File Content Models	Fuzzy Hash 1	Random Forest	512+ features	VBS
	Fuzzy Hash 2	SDCA	10M+ features	PDF
	Static, Dynamic and Contextual	Averaged Perceptron	16M+ features	
	Researcher Expertise, Fuzzy Hash	Averaged Perceptron	12M+ features	Macro
	File Emulation	DNN	150K+ features	
	File Detonation	DNN	10M+ features	

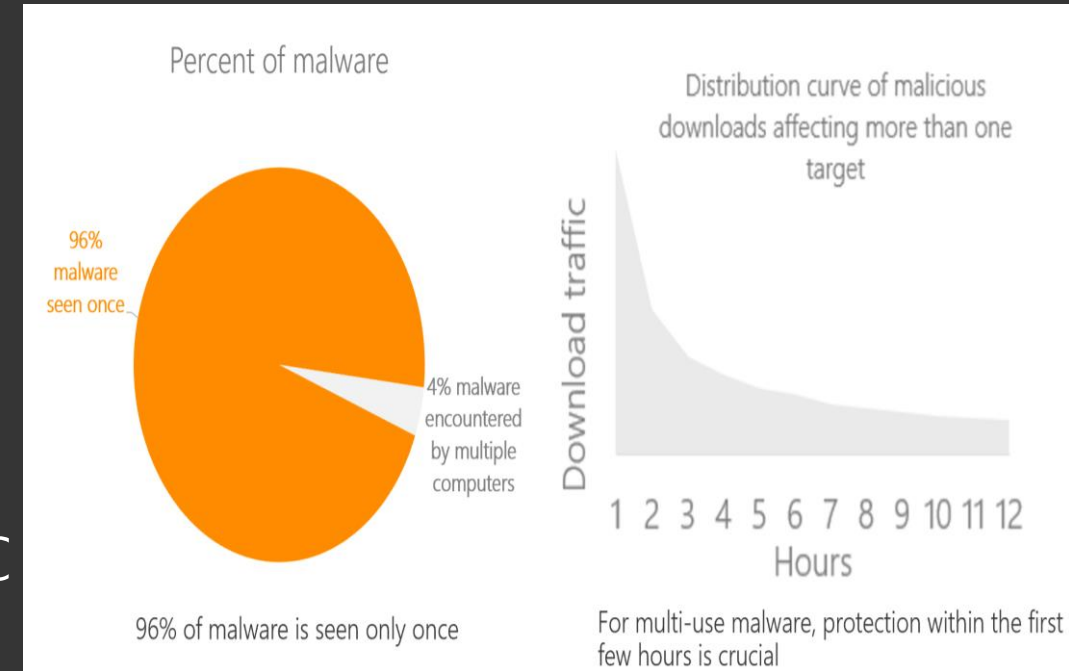


Optimizing for Different Threat Scenarios

Training Cadence: Classifiers:



- Malware
- Clean
- PUA
- Enterprise specific
- File Type specific



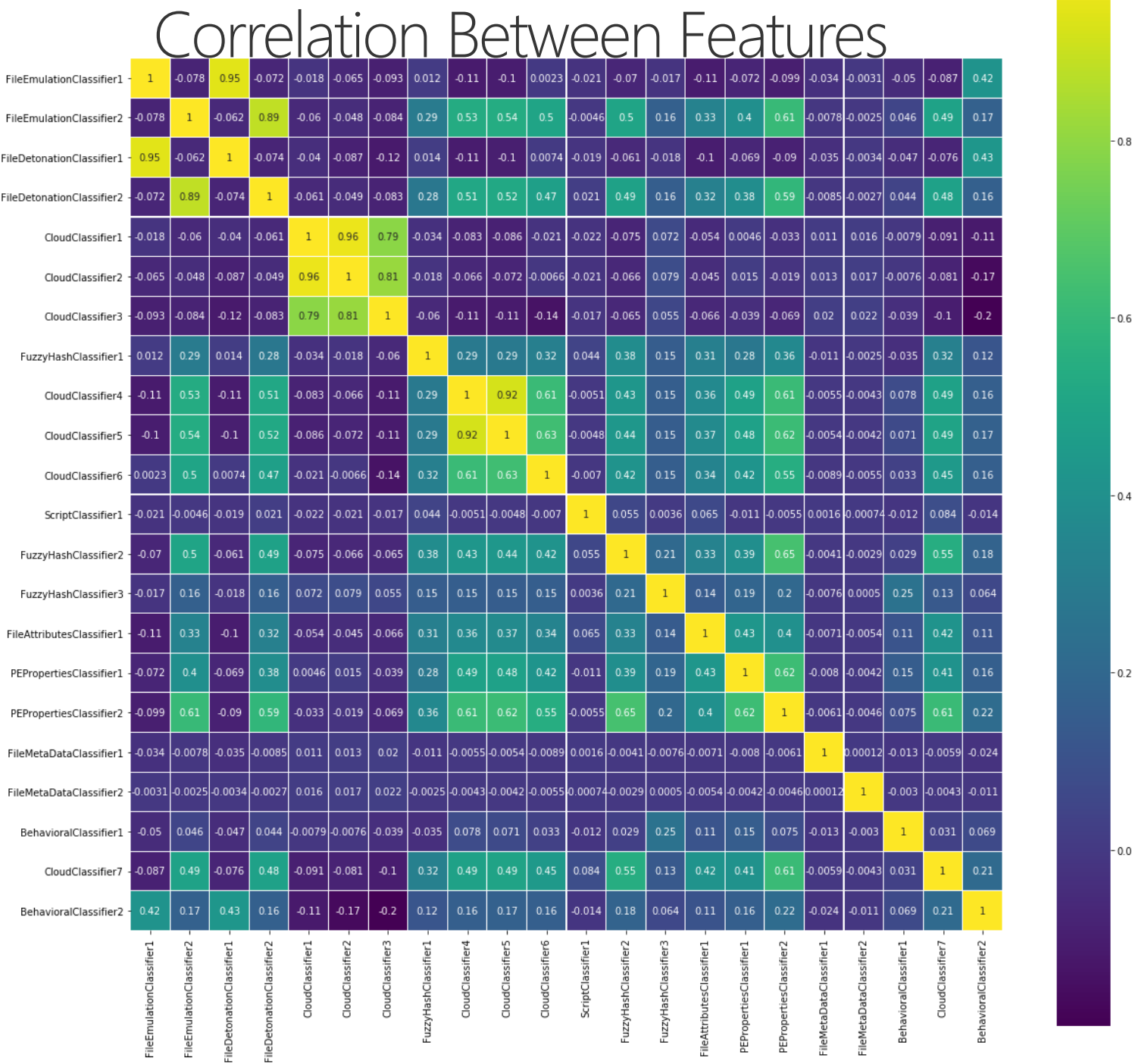
Feature Set

Training Algorithms

Training Data Sets

Optimization Settings

Feature Diversity

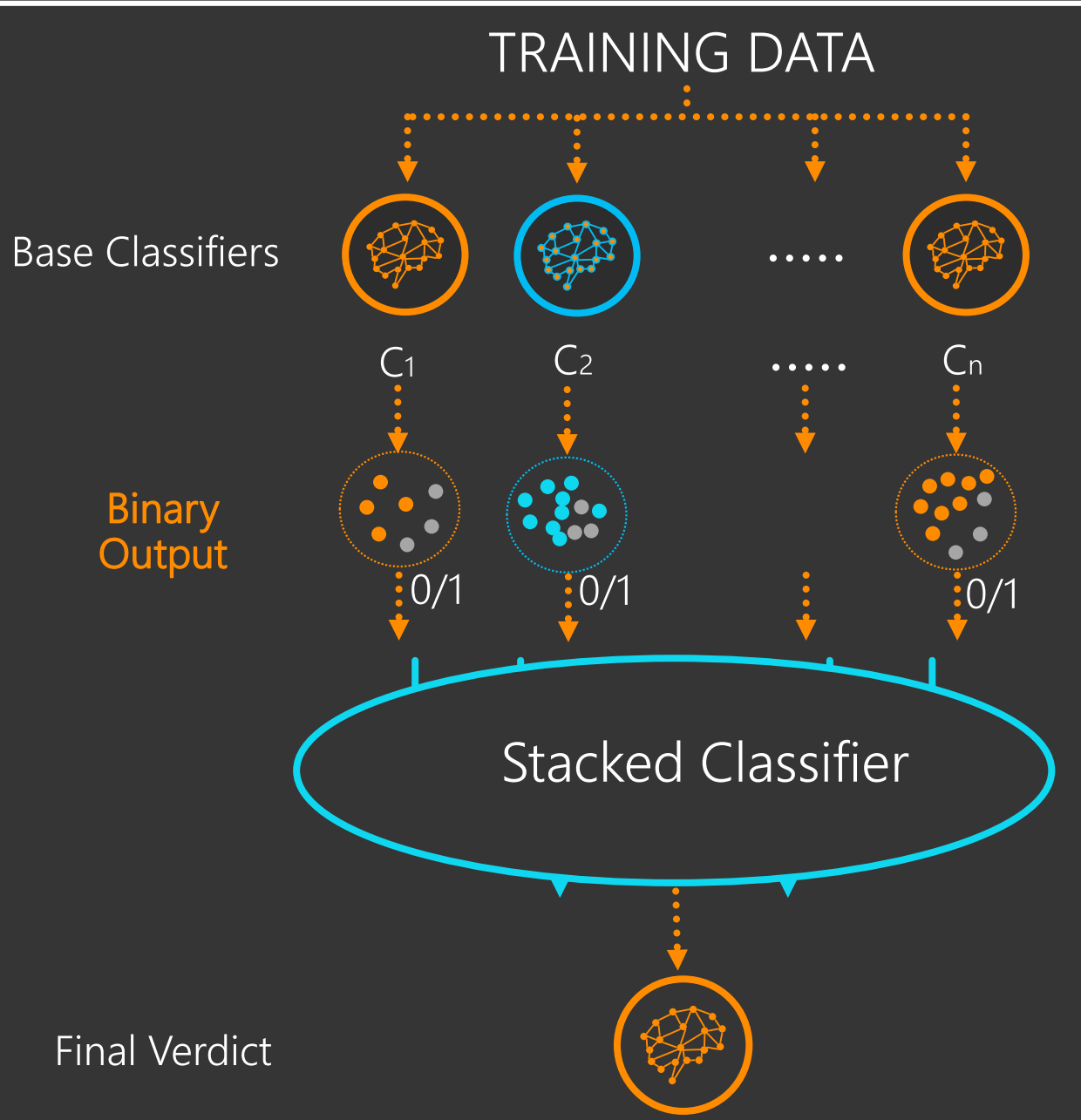


Developing the Stacked Model

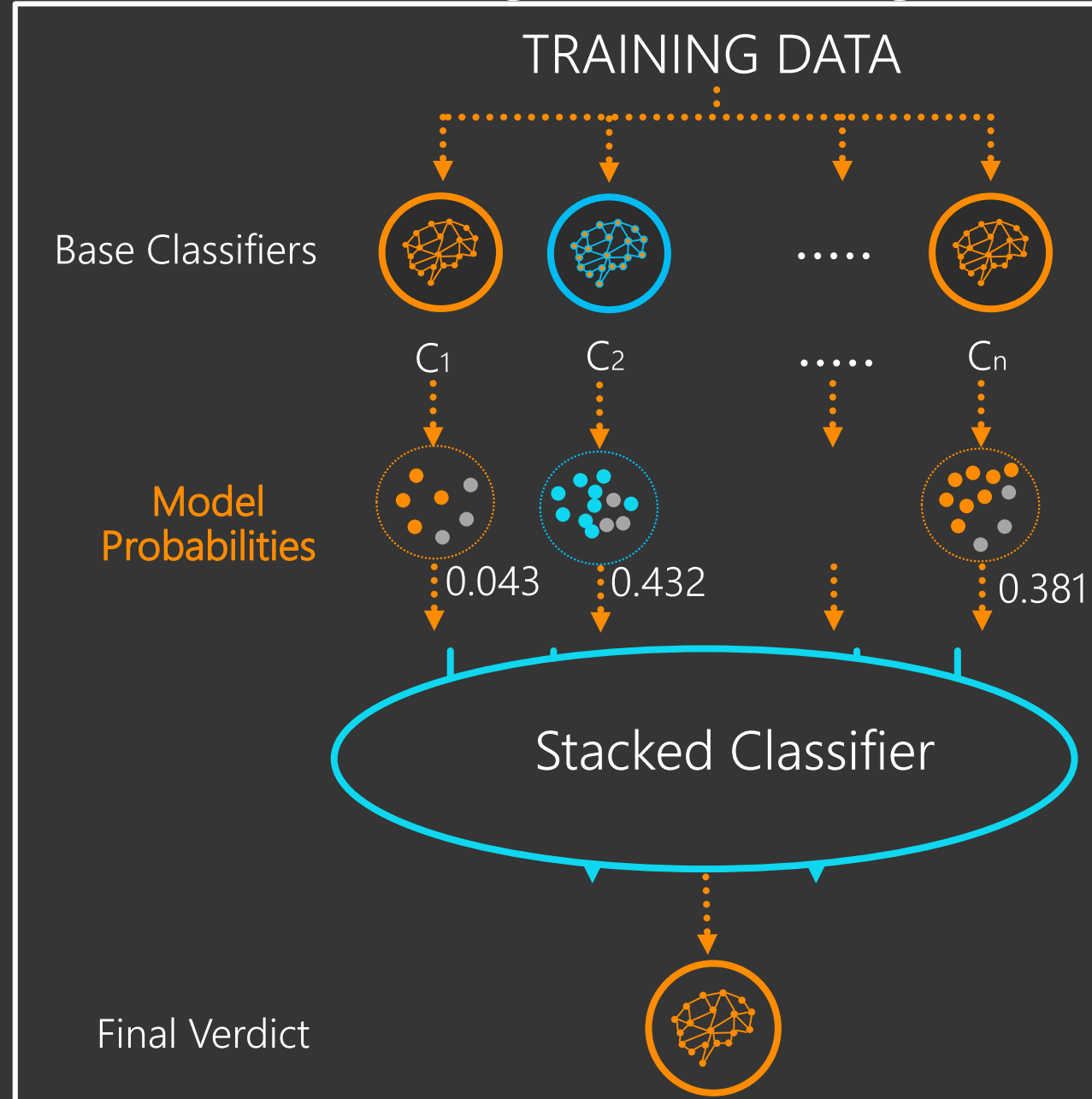
Stacked Ensemble

1. Boolean Stacking
2. Linear/ Logistic Stacking

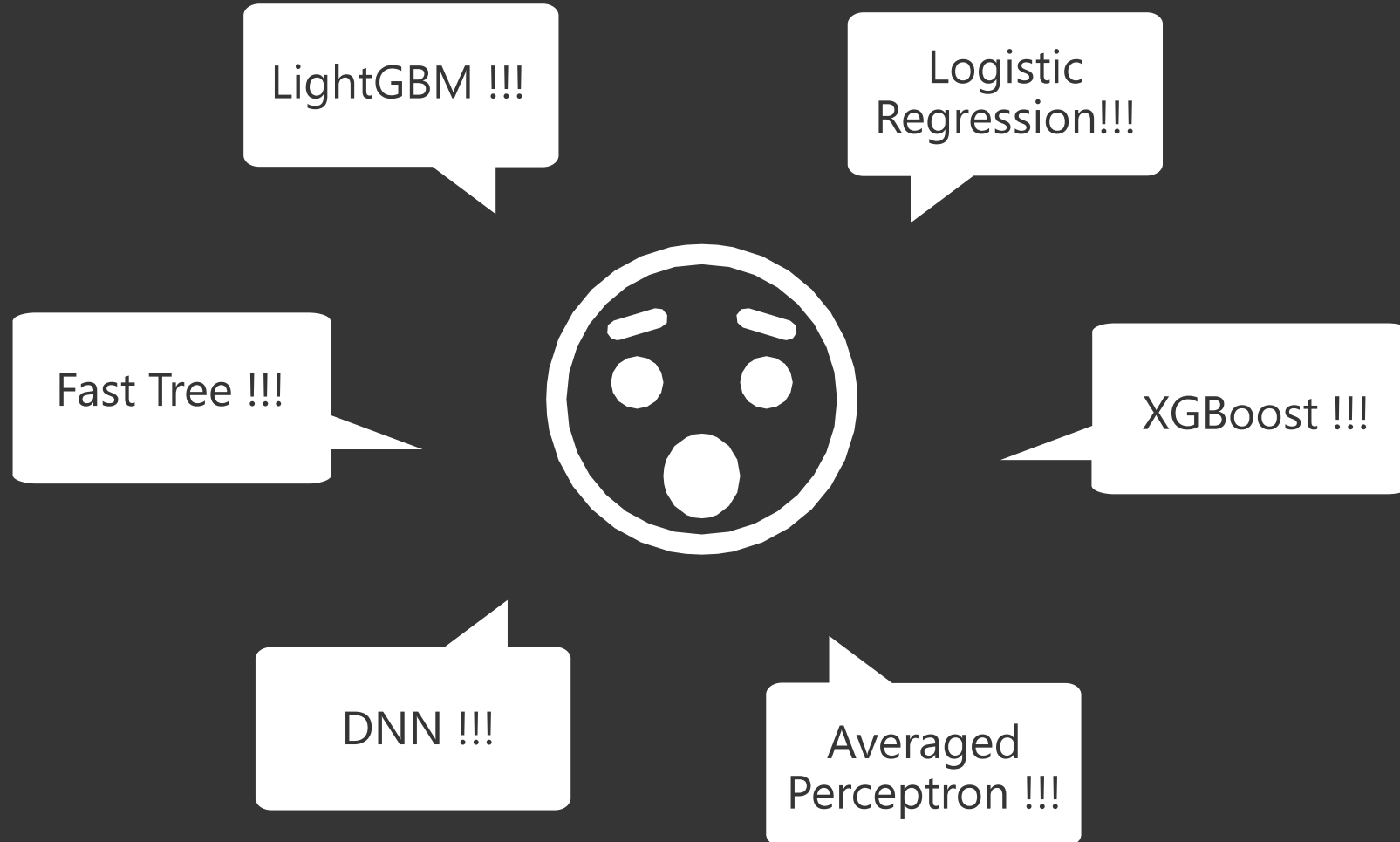
Boolean Stacking



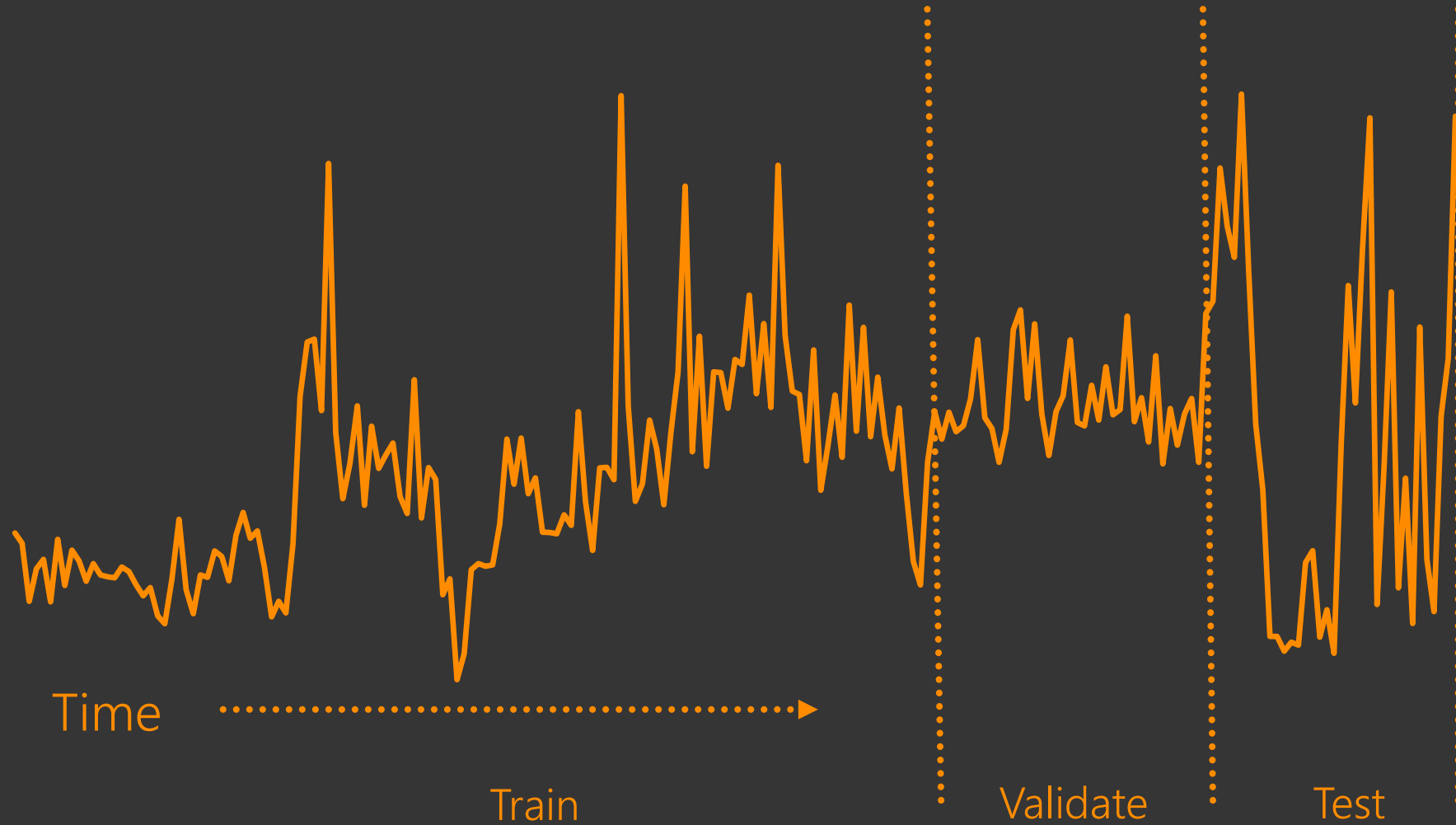
Logistic Stacking



Model Selection

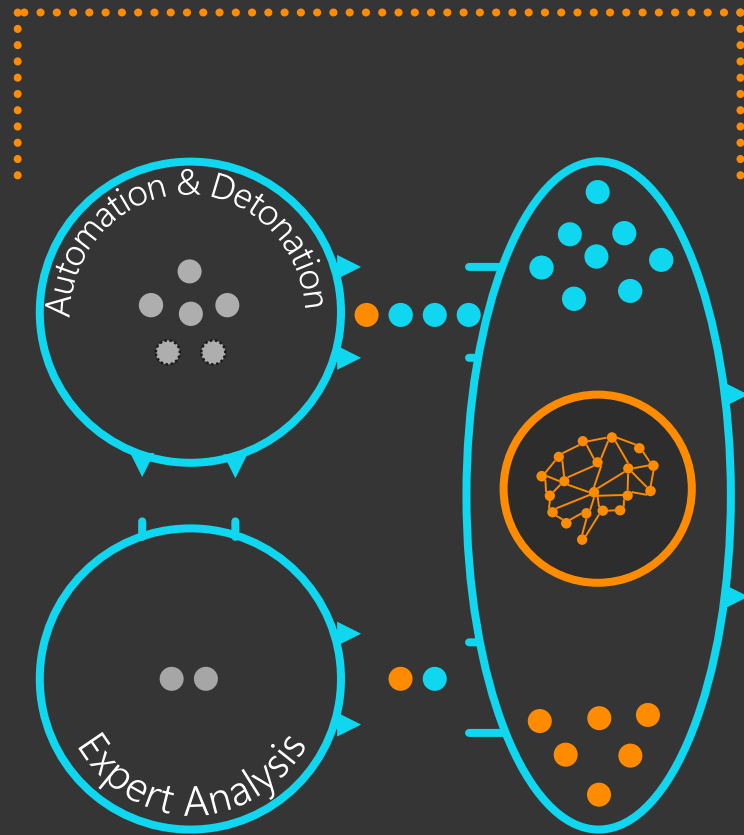


Experiment Design



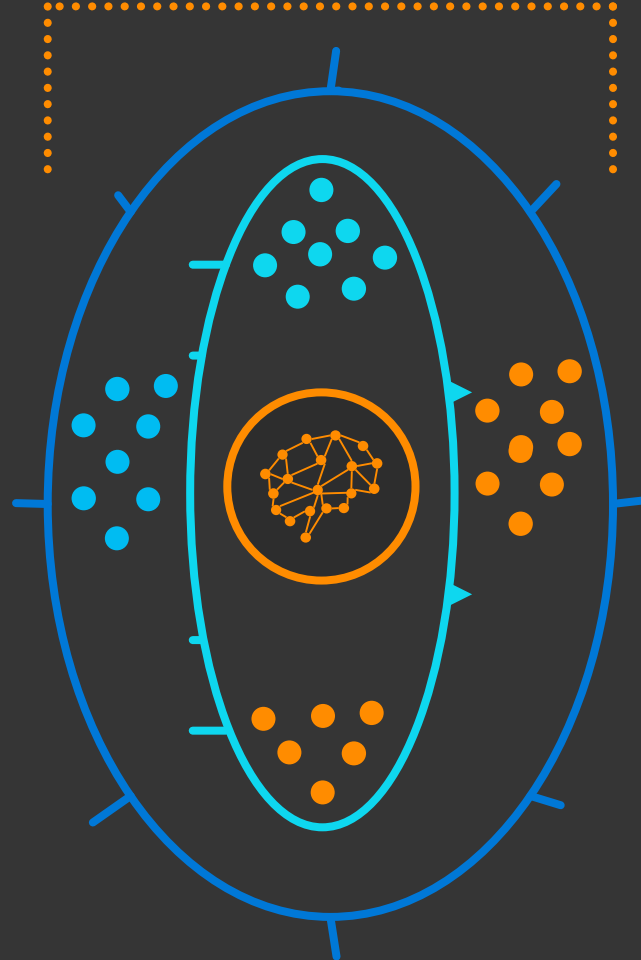
Experiment Design

Supervised Training



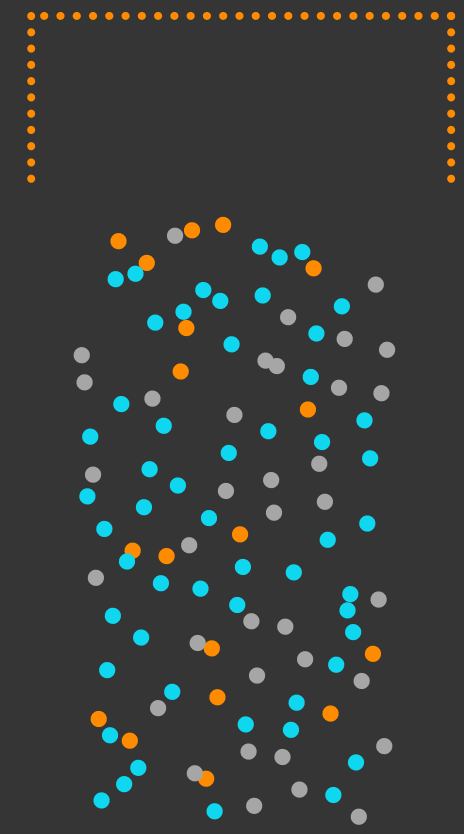
Train

Optimization



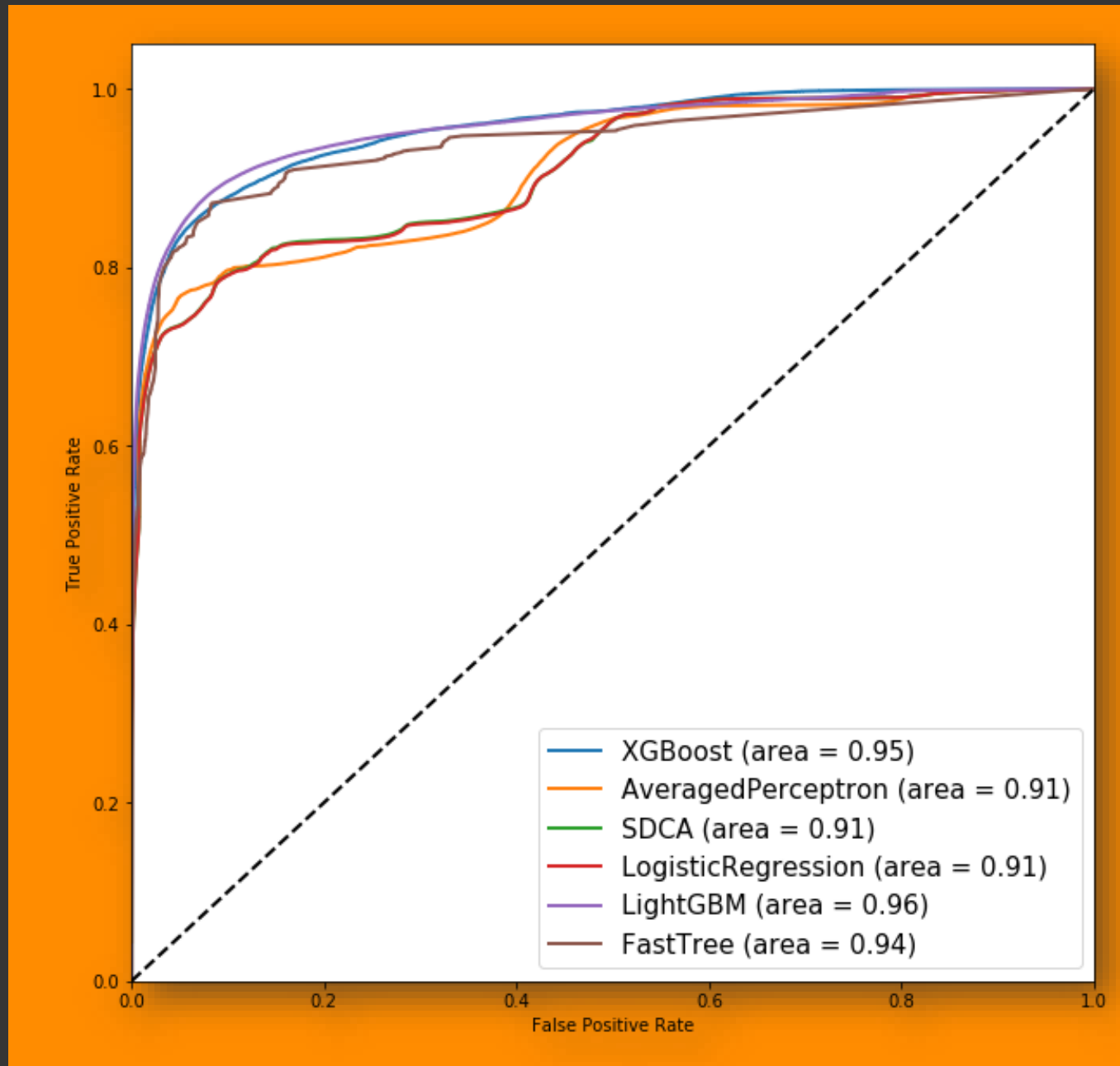
Validate

Evaluation



Test

Model Selection

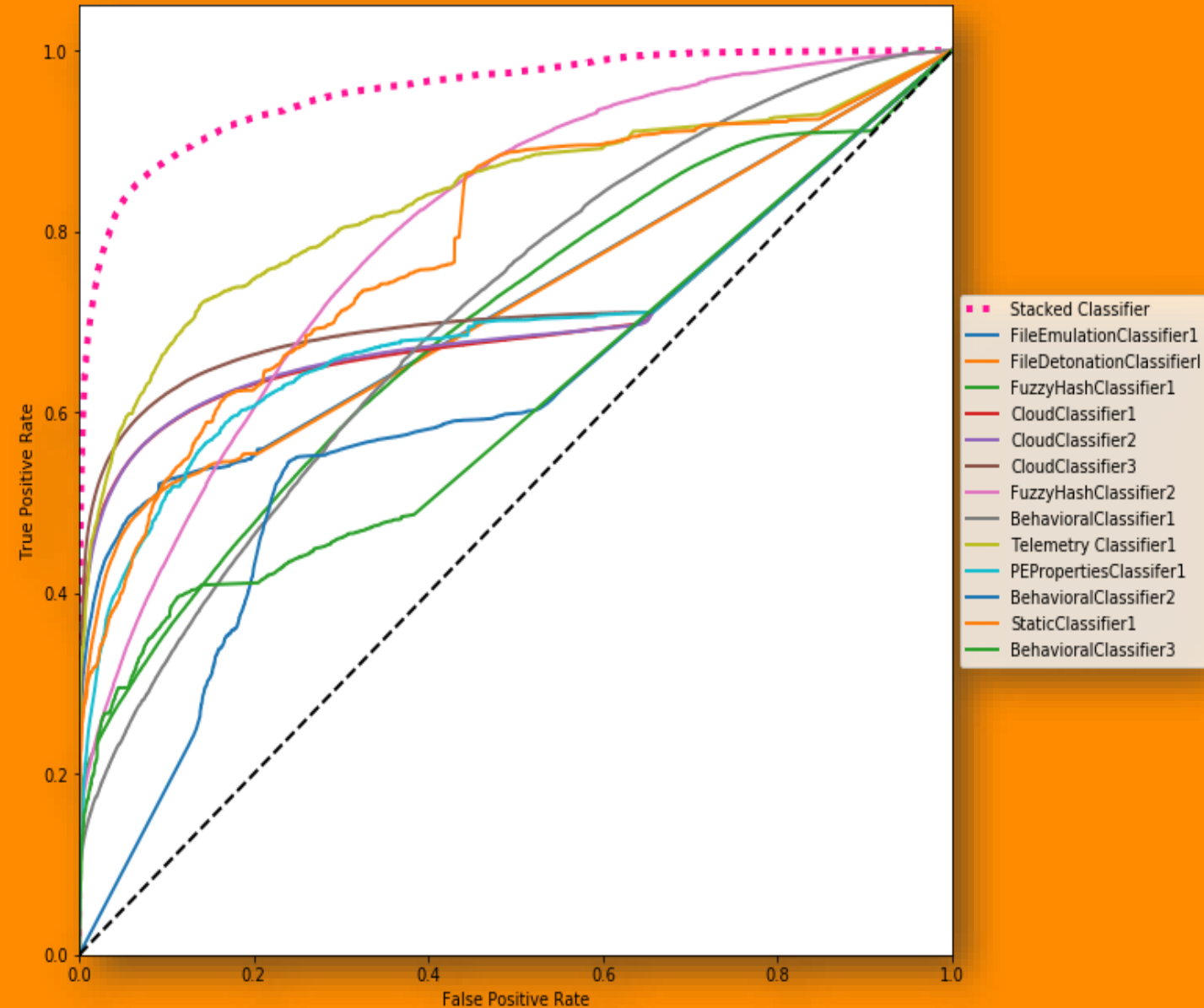


Final Training

- Generate diverse set of base classifiers
- Use model probabilities as input features to train the Stacked Classifier
- Use LightGBM to train the Stacked Classifier
- Plot the ROC curve for Stacked Classifier vs. Top Base Classifiers

Results

Stacked Ensemble Performance against top base classifiers



Evaluating on Live Data?

Model evaluated on time-split test set

Confusion table

PREDICTED	positive	negative	Recall
TRUTH			
positive	703,140	59,131	0.9224
negative	2,1030	8,013,623	0.9974
Precision	0.9710	0.9927	
OVERALL ACCURACY: 0.9811			

Model evaluated on Live Data for 60 mins without any calibrations

Confusion table

PREDICTED	positive	negative	Recall
TRUTH			
positive	2,177	21,182	0.0932
negative	14,004	2,097,228	0.9934
Precision	0.1345	0.9900	
OVERALL ACCURACY: 0.9835			

Testing the Model

Data Leaks

- Information from the target inadvertently works its way into the model-checking mechanism
- Causes an overly optimistic assessment of generalization performance
- Filtering features that directly correlate to the training labels

With some data leaks

Confusion table			
PREDICTED	positive	negative	Recall
TRUTH			
positive	703,140	59,131	0.9224
negative	2,1030	8,013,623	0.9974
Precision	0.9710	0.9927	
OVERALL ACCURACY: 0.9811			

Filtering known data leaks

Confusion table			
PREDICTED	positive	negative	Recall
TRUTH			
positive	625,324	136,947	0.8203
negative	29,540	8,005,113	0.9963
Precision	0.9549	0.9832	
OVERALL ACCURACY: 0.9668			

10% drop in Recall

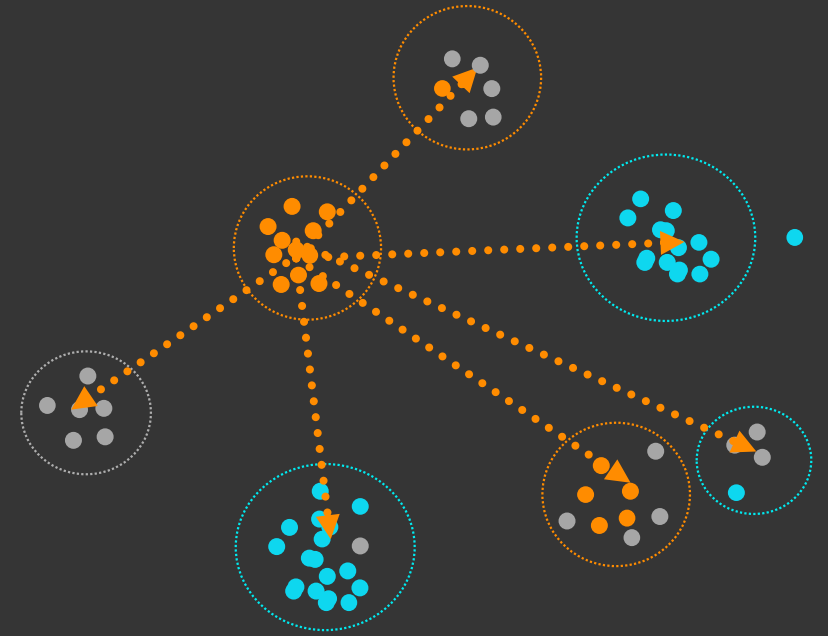
Handling Missing Values

- Not all Base Classifiers classify every threat scenario
- What you can do:
 - Retaining the instance but..
 - Adding Boolean features indicating what features were missing
 - Cross Join between features
 - Interpretable models

Model	Probability	Verdict
FileEmulation1	N/A	Unknown
FileDetonation	N/A	Unknown
FuzzyHash1	N/A	Unknown
FuzzyHash2	0.014020299	Clean
CloudClassifier1	N/A	Unknown
CloudClassifier2	N/A	Unknown
CloudClassifier3	N/A	Unknown
CloudClassifier4	N/A	Unknown
CloudClassifier5	N/A	Unknown
CloudClassifier6	N/A	Unknown
.	.	.
.	.	.
.	.	.
ResearcherExpertise	0.07285905	Clean
PEPropertiesClassifier	N/A	Unknown
FileMetaDataClassifier1	N/A	Unknown
FileMetaDataClassifier2	N/A	Unknown
FileMetaDataClassifier3	N/A	Unknown
BehavioralClassifier1	N/A	Unknown
BehavioralClassifier2	N/A	Unknown
Stacked Ensemble	0.92	Malware

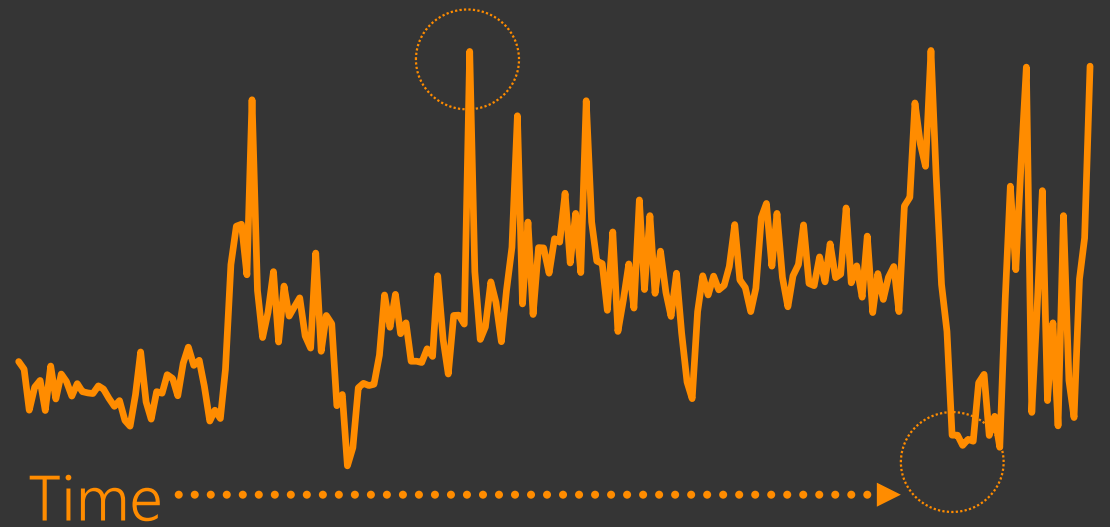
Using Unsupervised Features

- Adding K-means distance for each instance from the centroid of each cluster as an input feature



Other Improvements

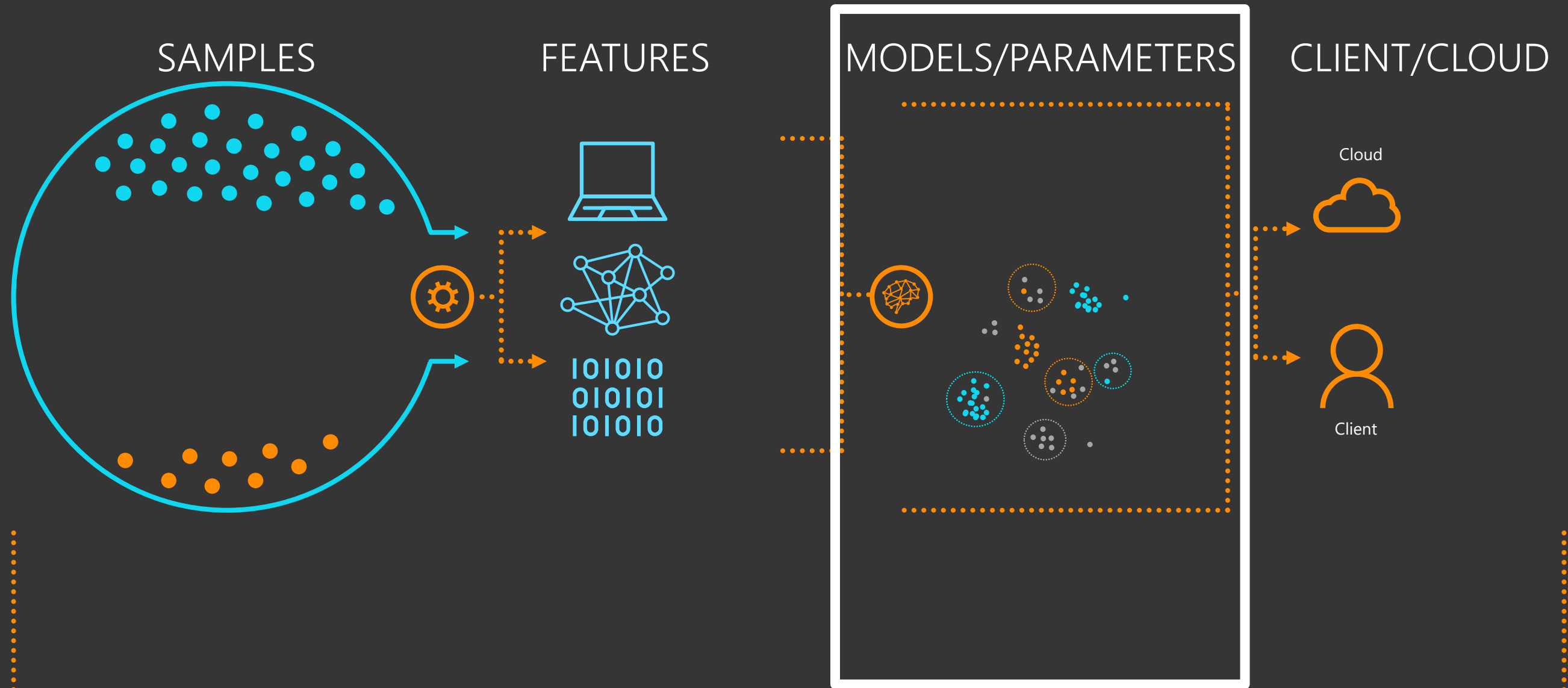
- Maintaining a fixed label distribution for training
- Continuous monitoring of incoming telemetry to catch anomalies/ outliers before training



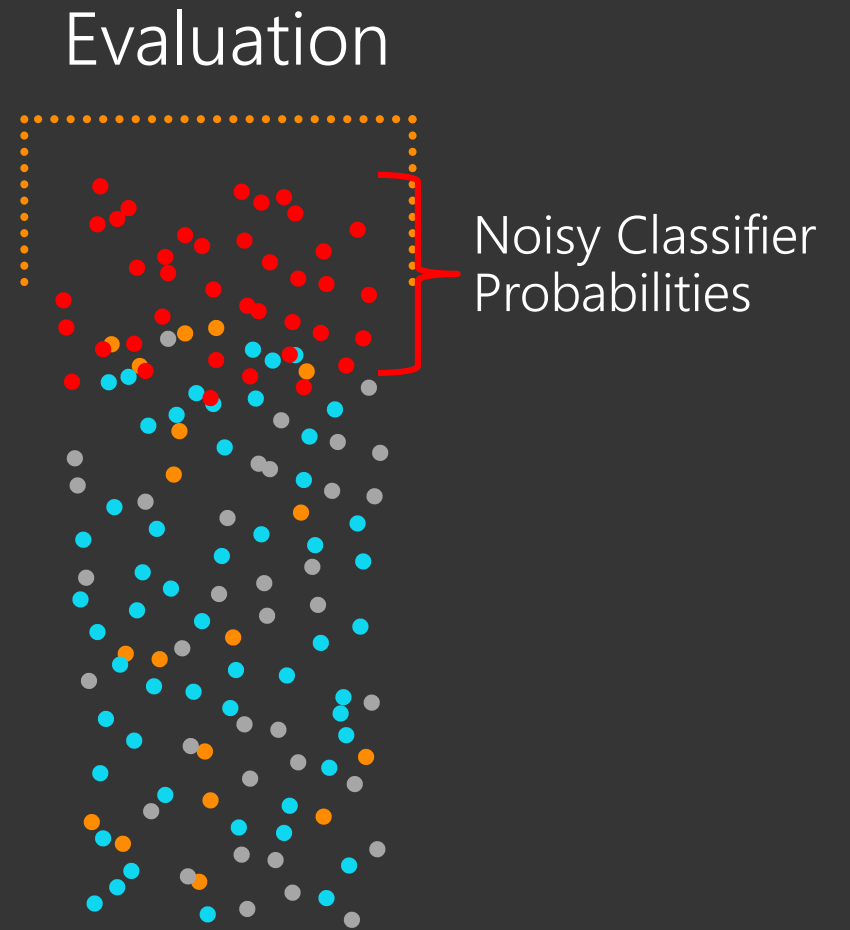
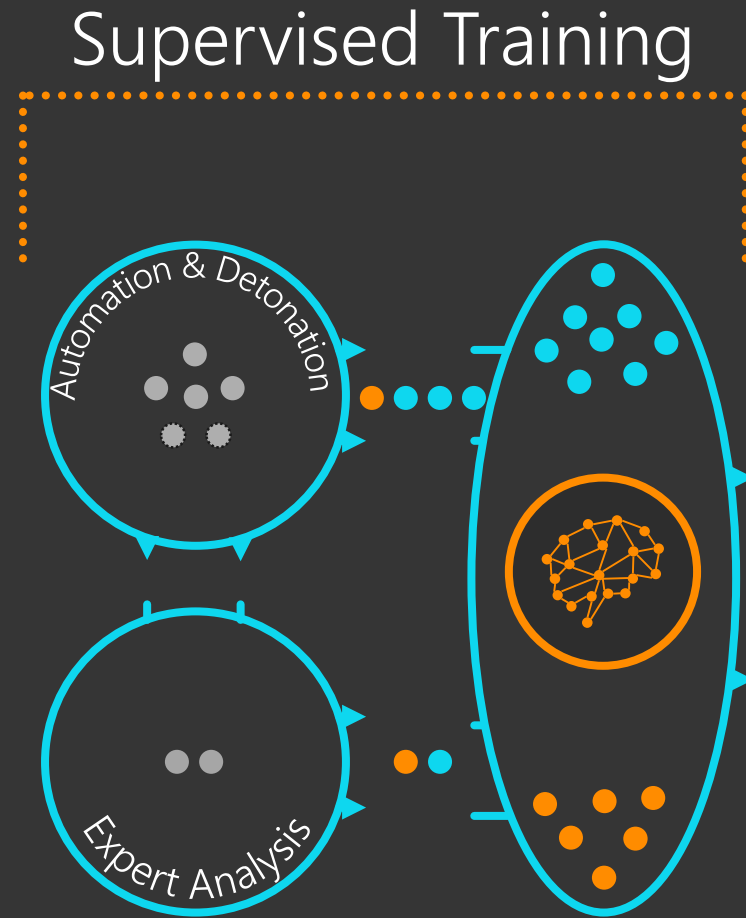
Model Deployed !!!

But is it Resilient to Adversarial
Attacks...

What if ... We evaluate on rogue/ noisy classifiers as features



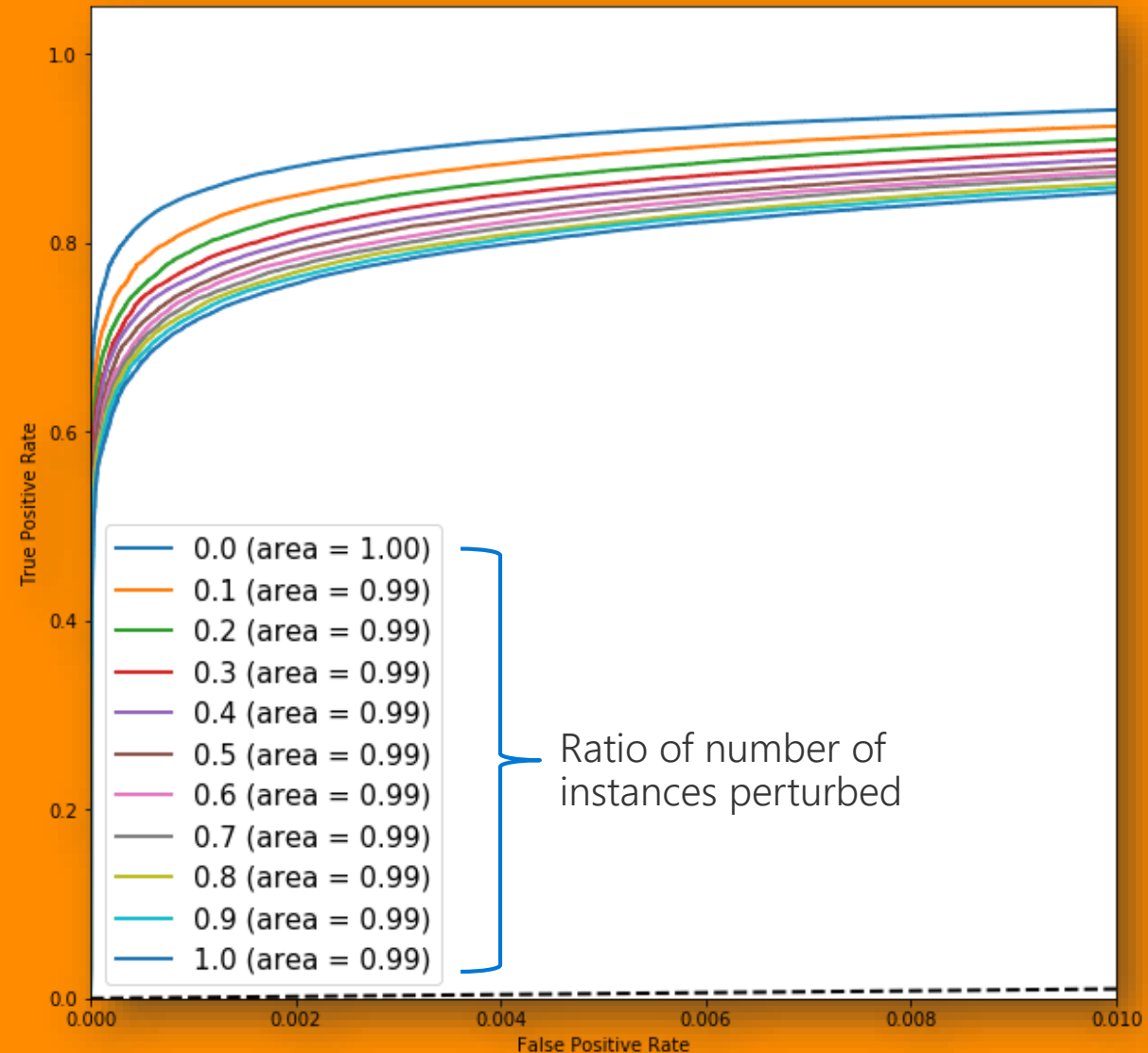
Experiment Design



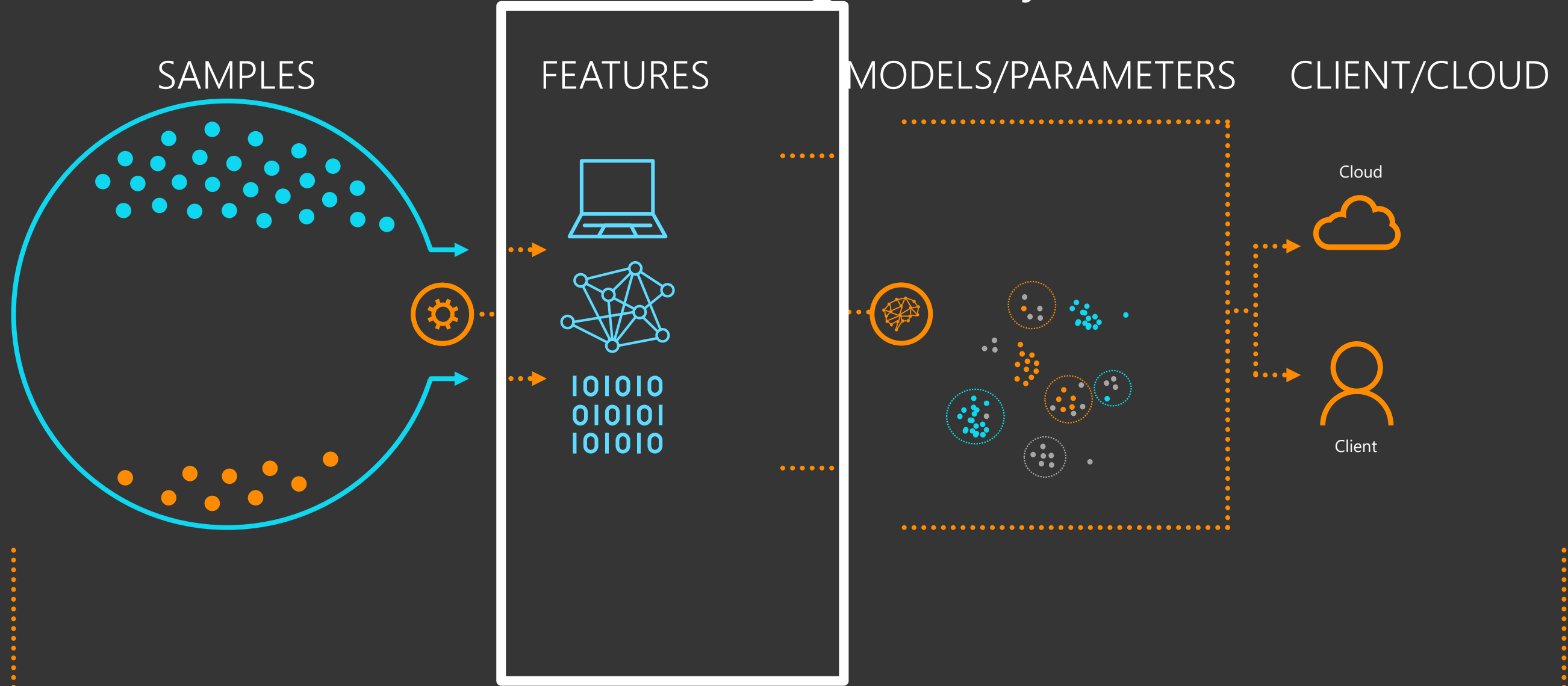
Experimental Verification

Trust, but Verify!!!

The classifier is robust to one of the classifiers being compromised at 1% FPR.

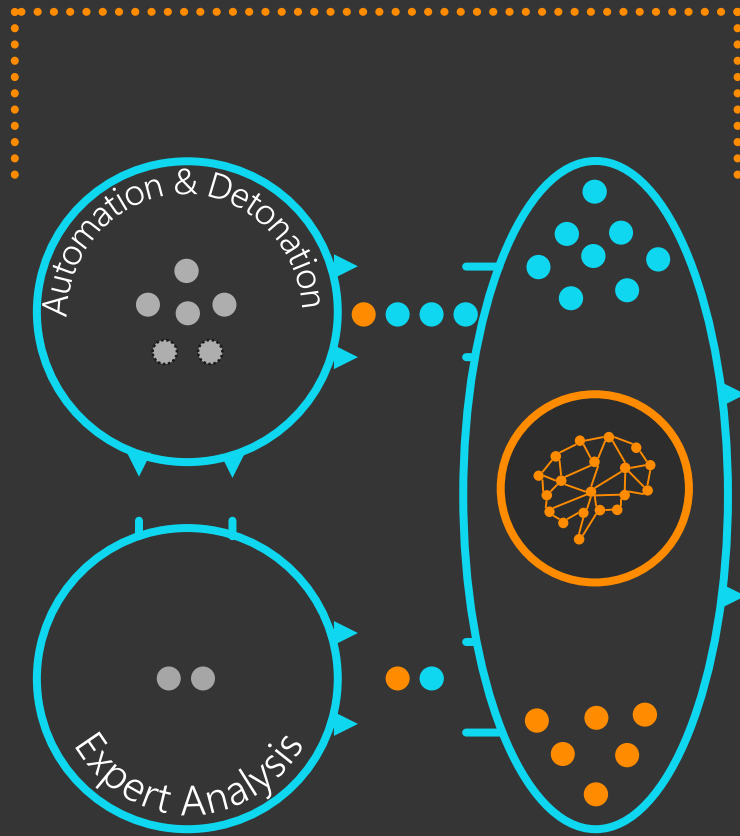


What if ... We train on new rogue/ noisy classifiers as features

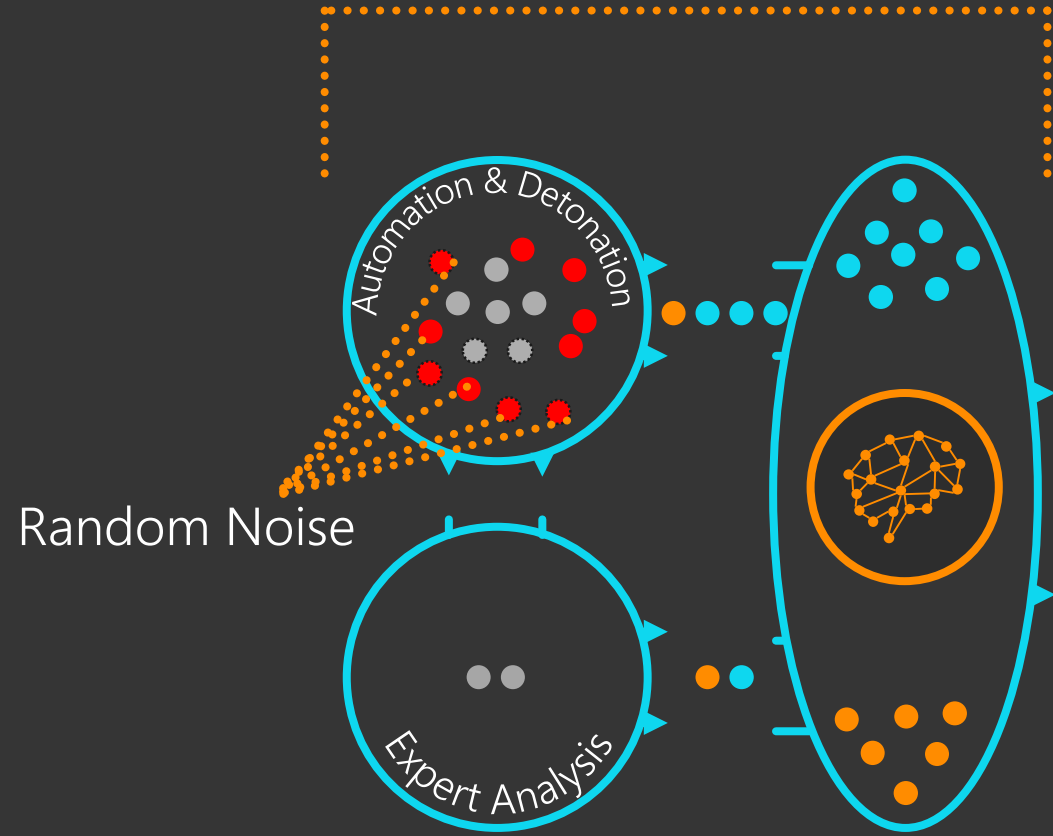


Experimental Verification

Supervised Training



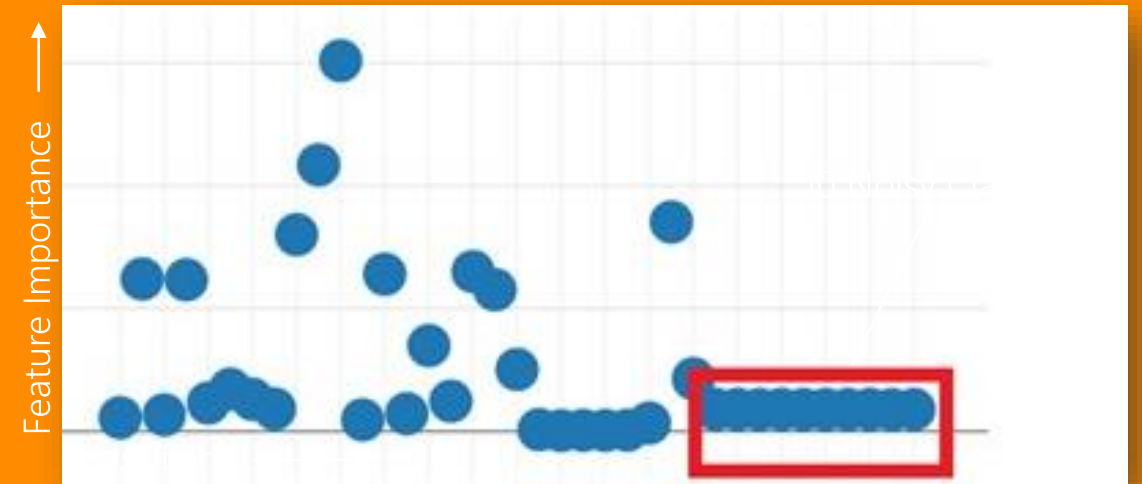
Adding Rogue Features



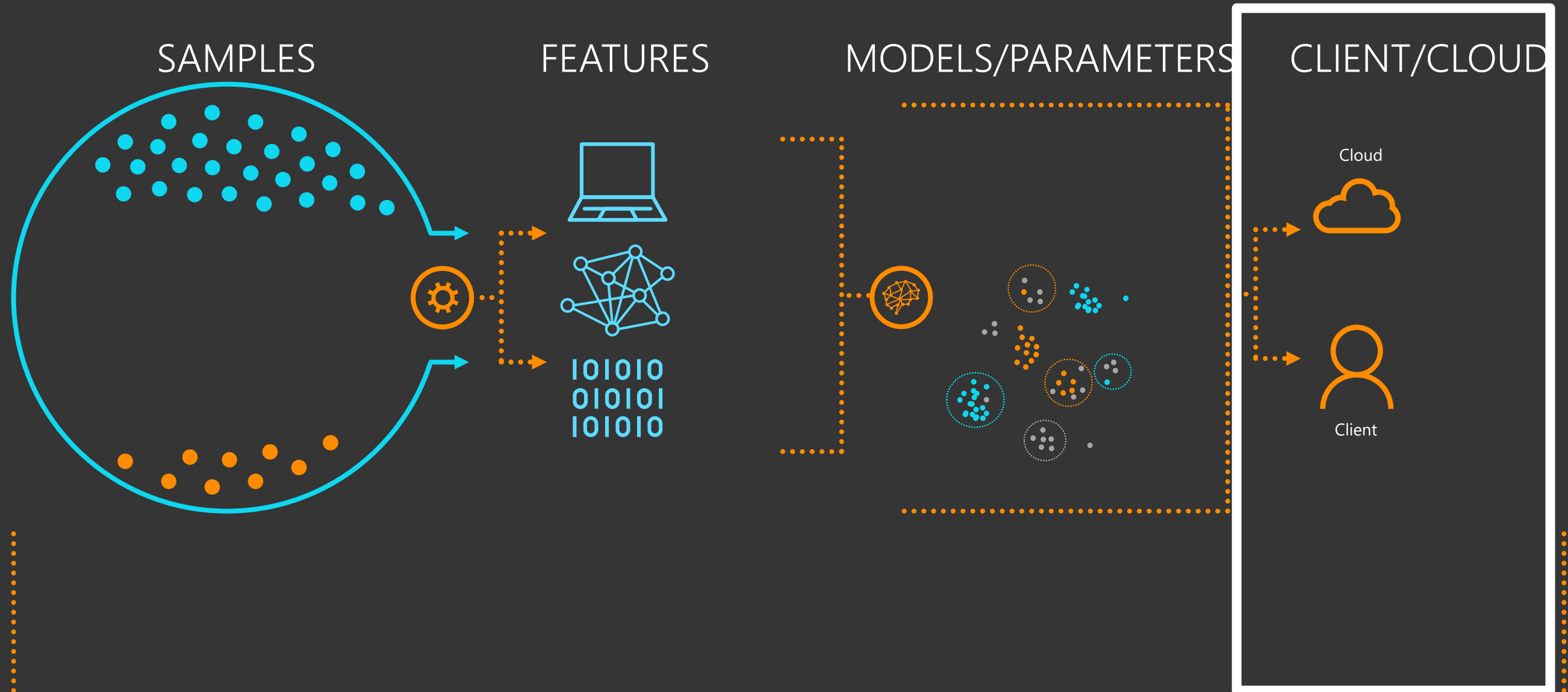
Experimental Verification

# of Random Classifiers	False Positive Rate	True Positive Rate
0	0.8746%	96.1824%
2	0.8834%	96.1222%
4	0.8912%	96.0385%
6	0.8939%	95.8932%
8	0.8974%	95.8462%
10	0.9131%	95.8462%

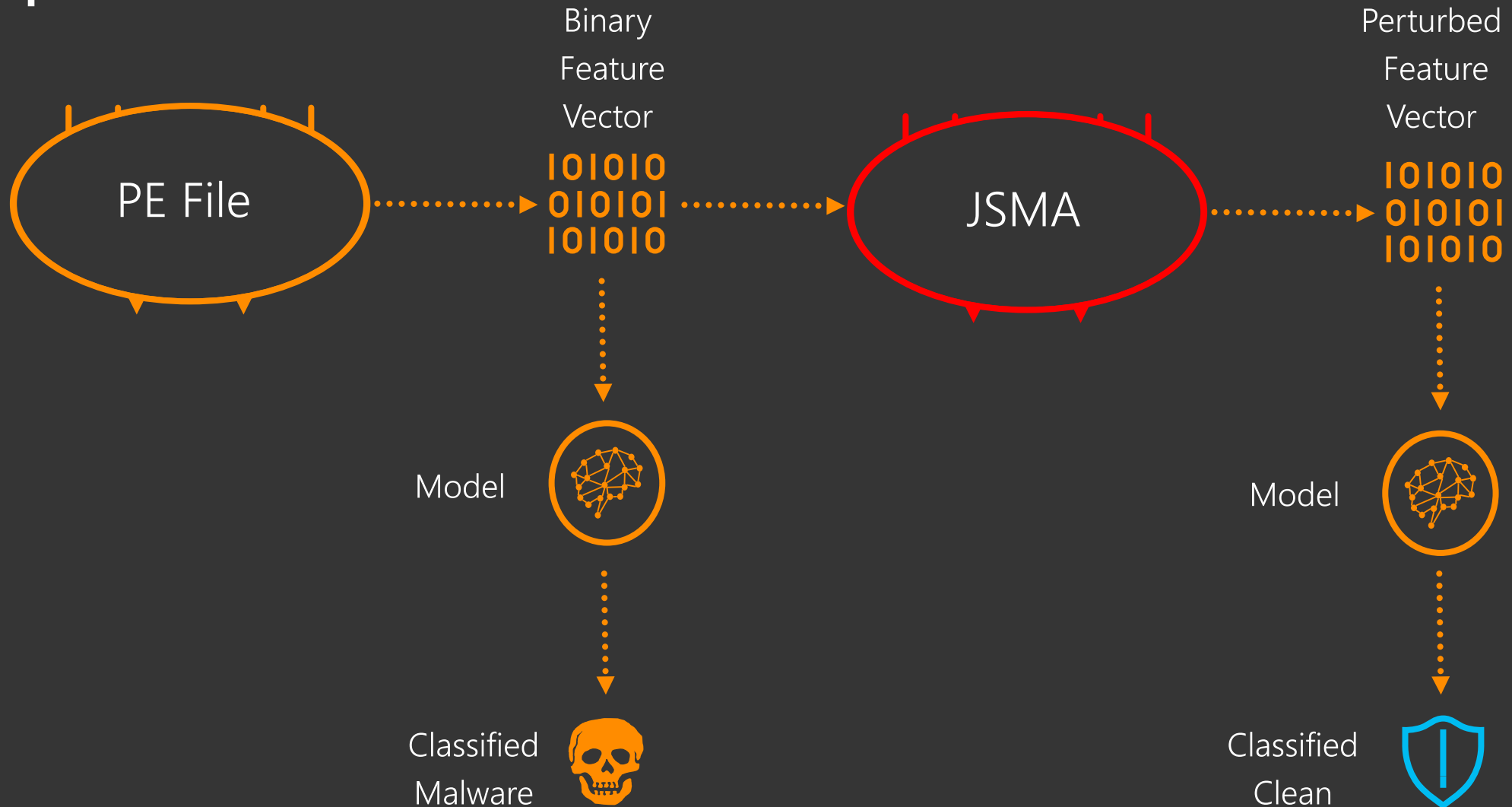
The classifier can detect these random noises and the performance drop is negligible.



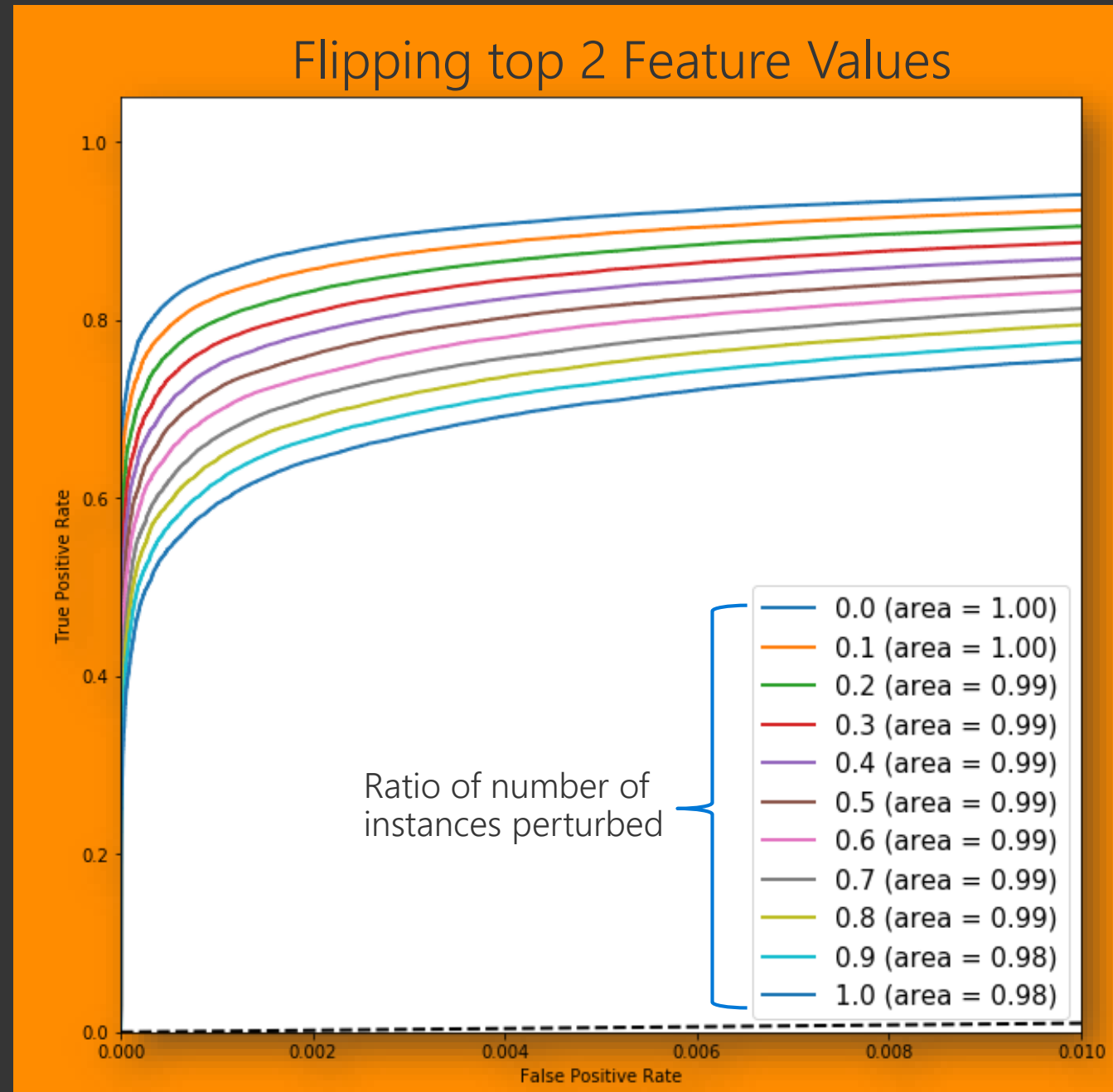
What if ... Attacker crafts adversarial samples to flip verdicts



Experimental Verification



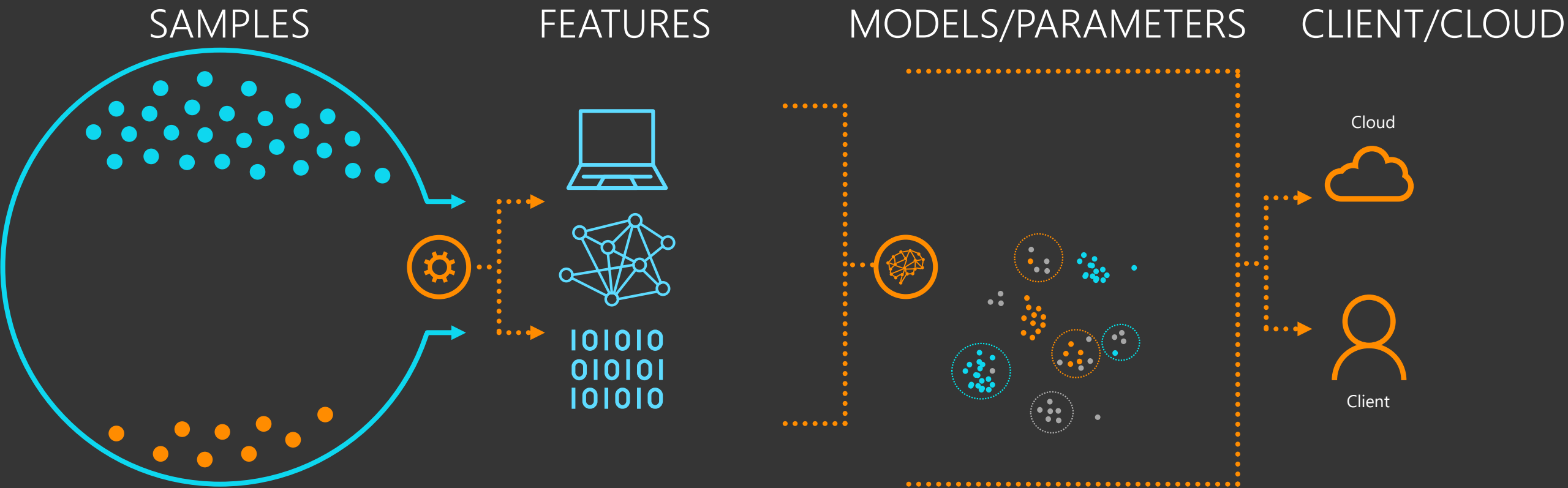
Experimental Verification



What if...

Attacker is highly motivated
to somehow just break our
Stacked Ensemble

Realtime Monitoring



of Instances
Data Distribution
Bias
Anomalies

Relevant features
Threat Landscape

Metrics Requirements
Diversity Requirements
Developing the Model

Continuous Monitoring
Overall Blocks
Telemetry

Ensemble Model Development and Testing

Ensemble ML Primer

Diversity Requirements

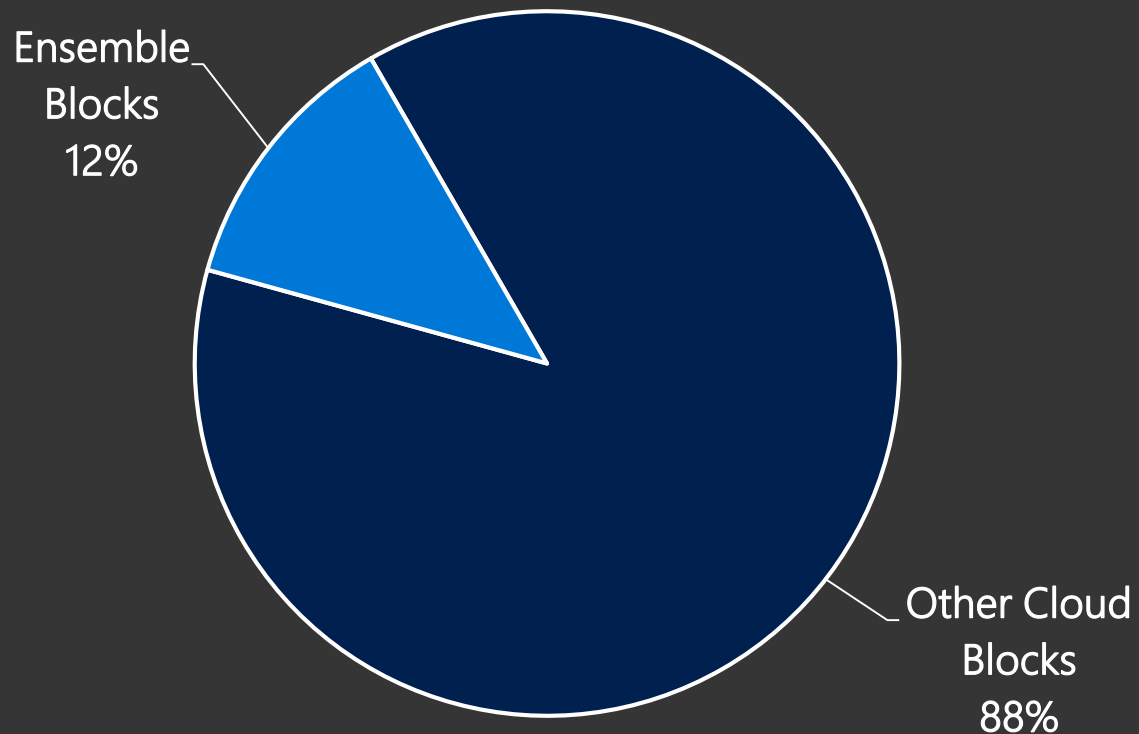
Developing the Model

Testing the Model

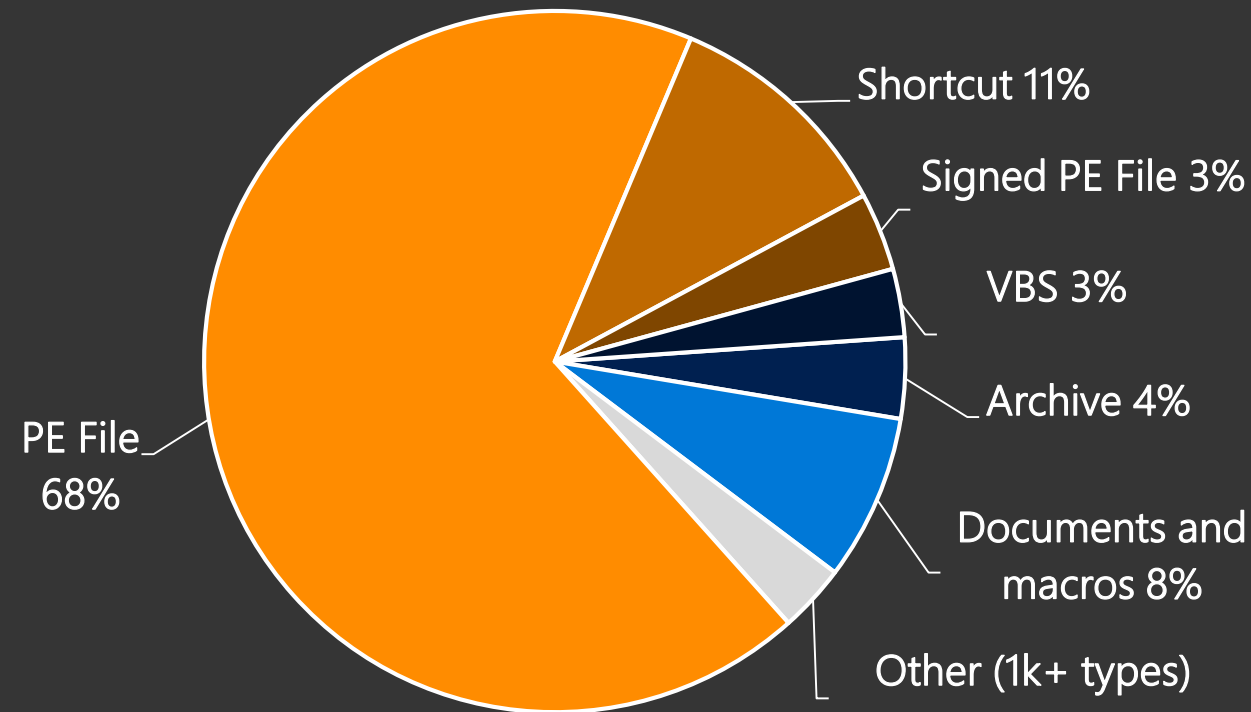
Results!

Impact of Ensemble Models

Percent of Threats Blocked by Cloud Protection
June 2018

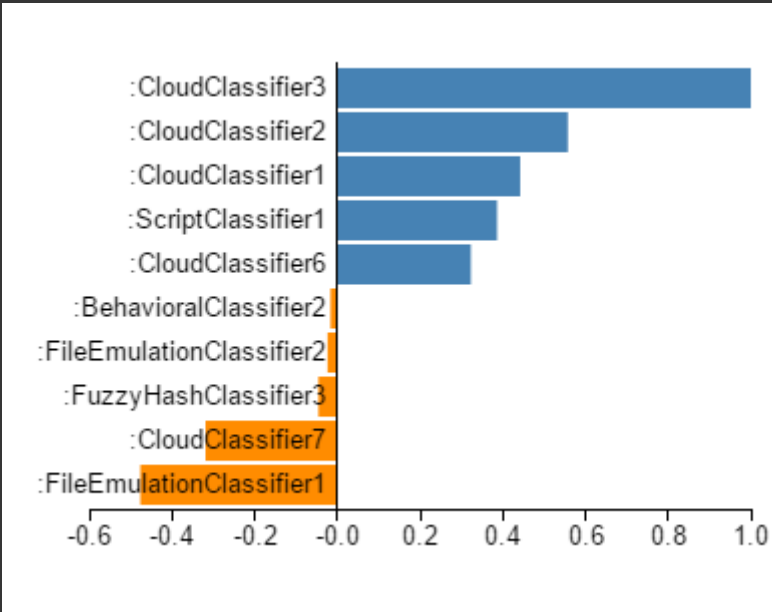


Ensemble Model Blocks by File Type
June 2018

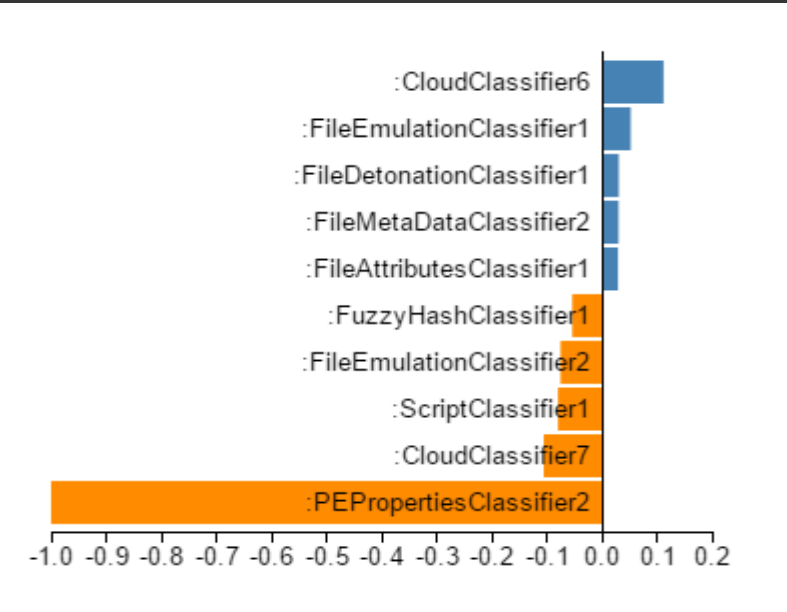


Bonus: Interpretability

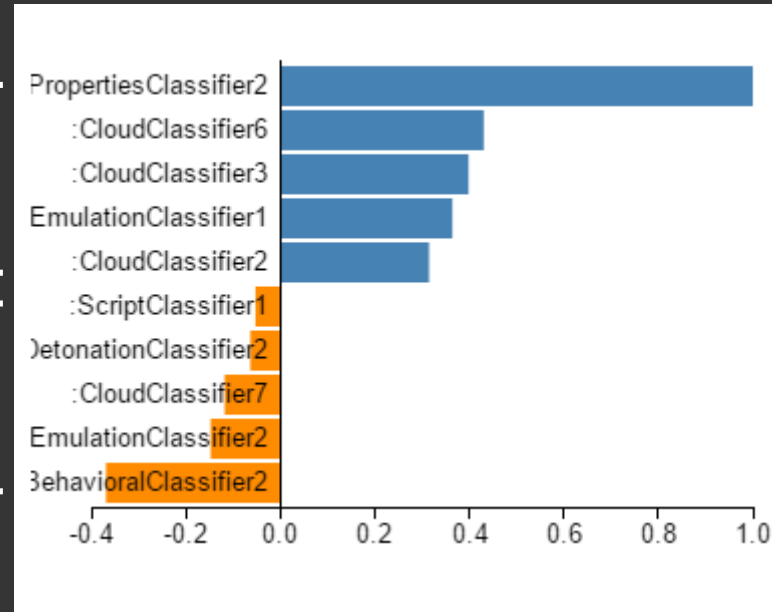
True Positive



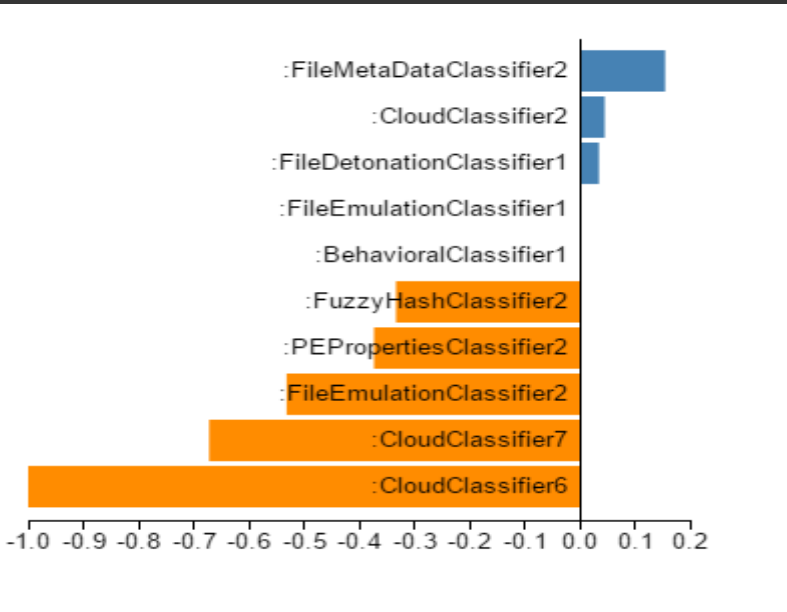
False Negative



False Positive



True Negative



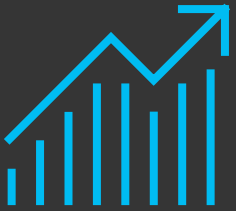
Top classifiers contributing to **malware** verdict

Top classifiers contributing to **clean** verdict

Benefits of an Ensemble Model



Filters out noisy signals from an occasionally underperforming model

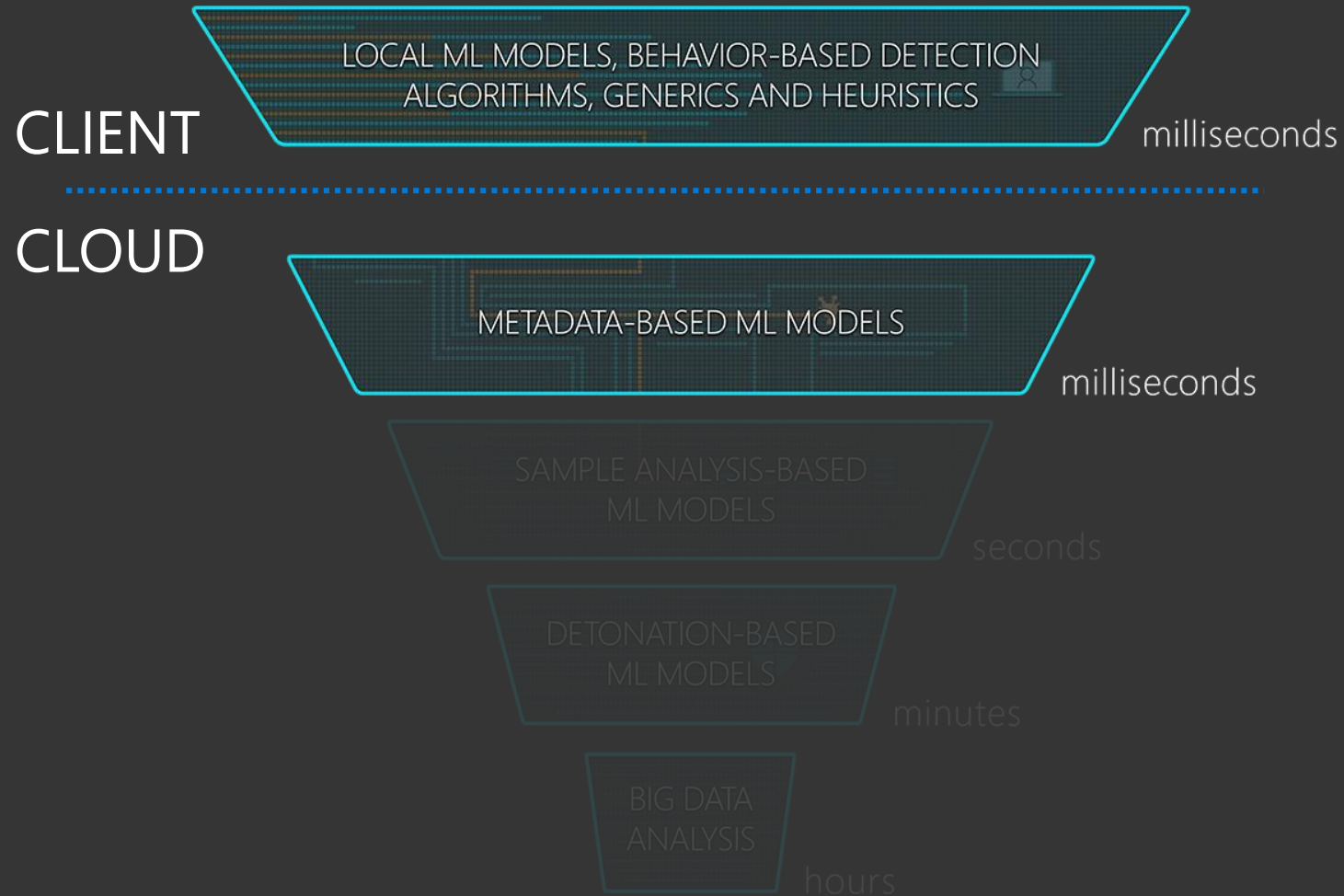


Increases predictive power with easy interpretability



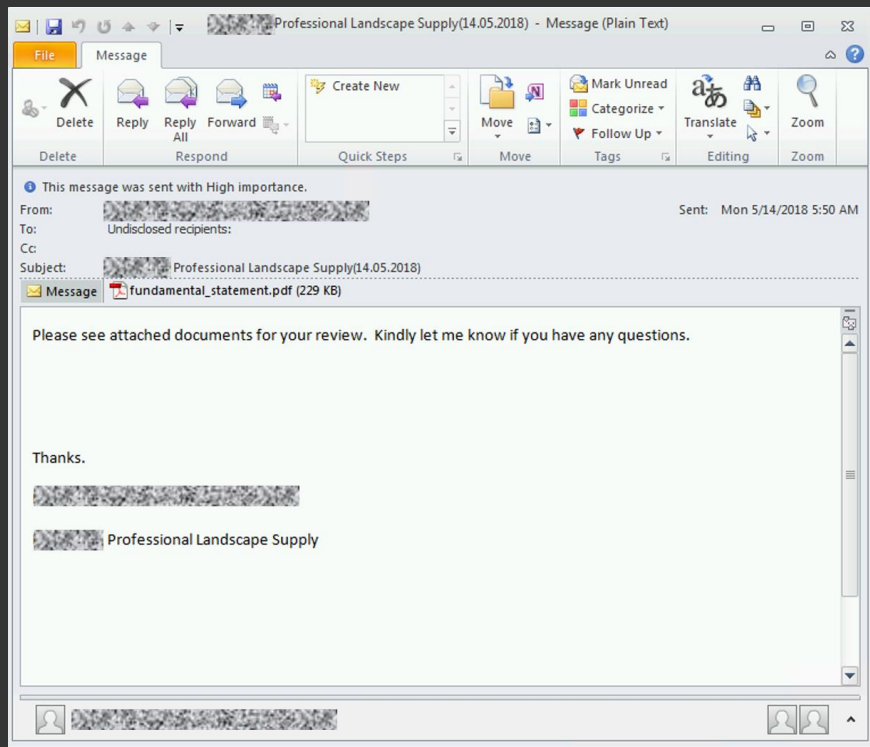
Adds resilience against attacks on individual models

Recent Realworld Case Studies (2)

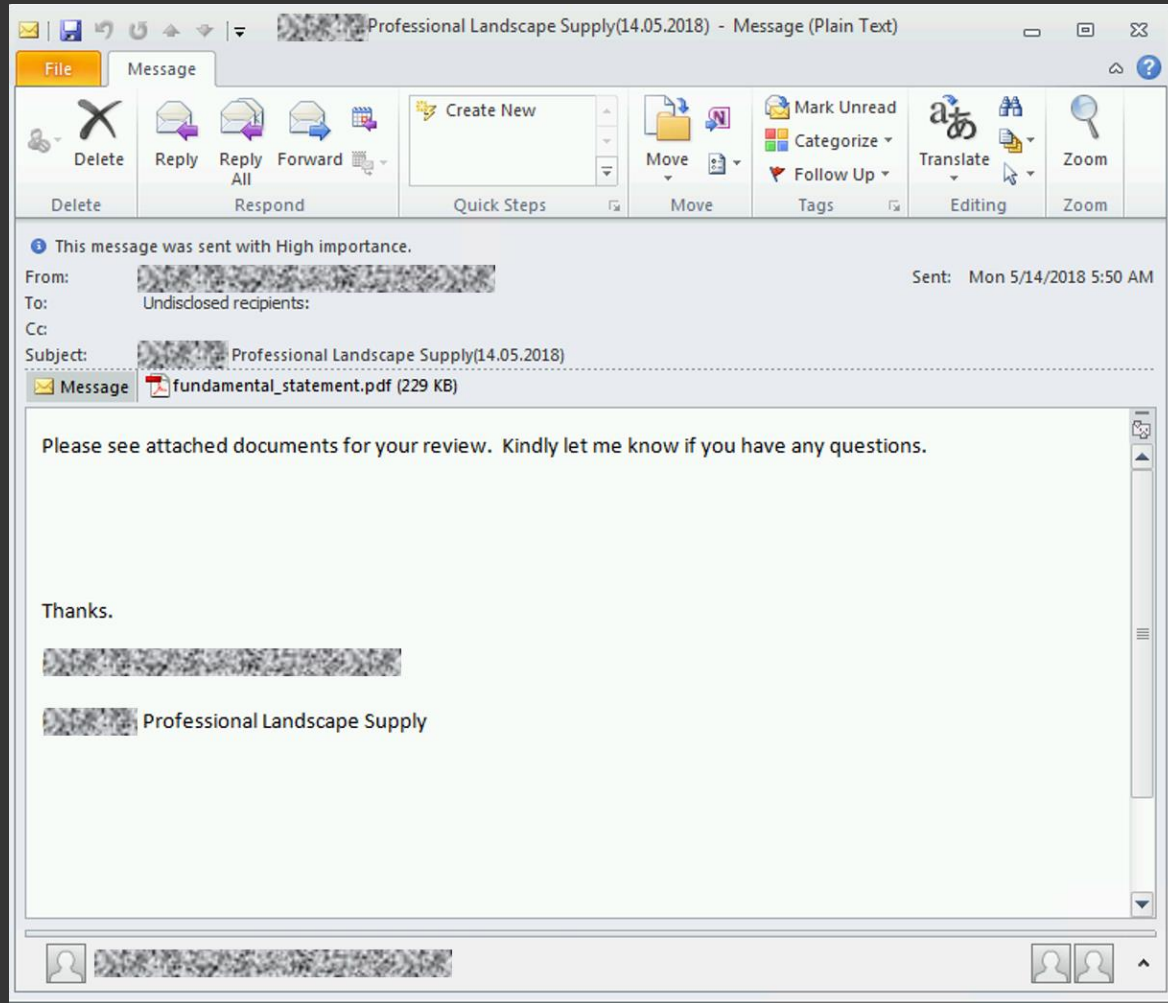


Case Study 1: Spear Phishing

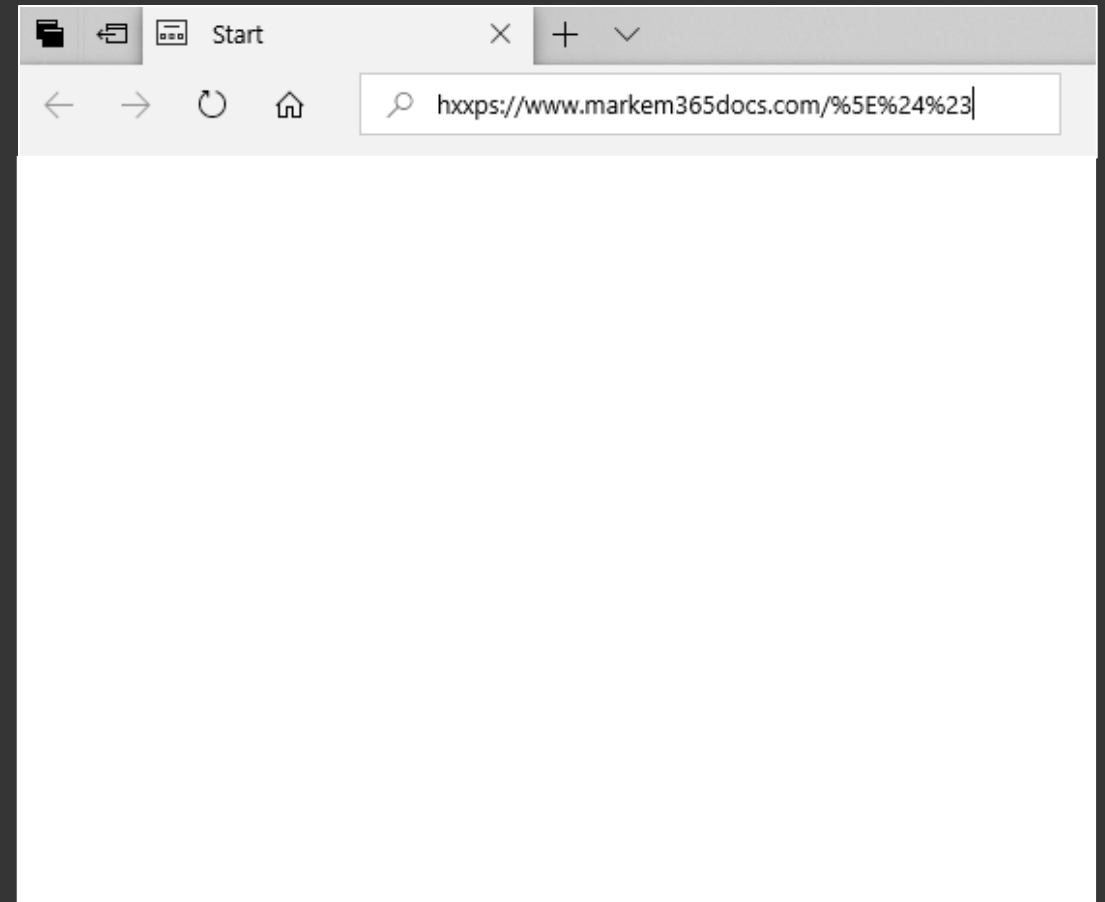
- Small-scale attack in Central and Western Canada
- Most targets reached within 5 ½ hours
- 73% of targets were commercial businesses



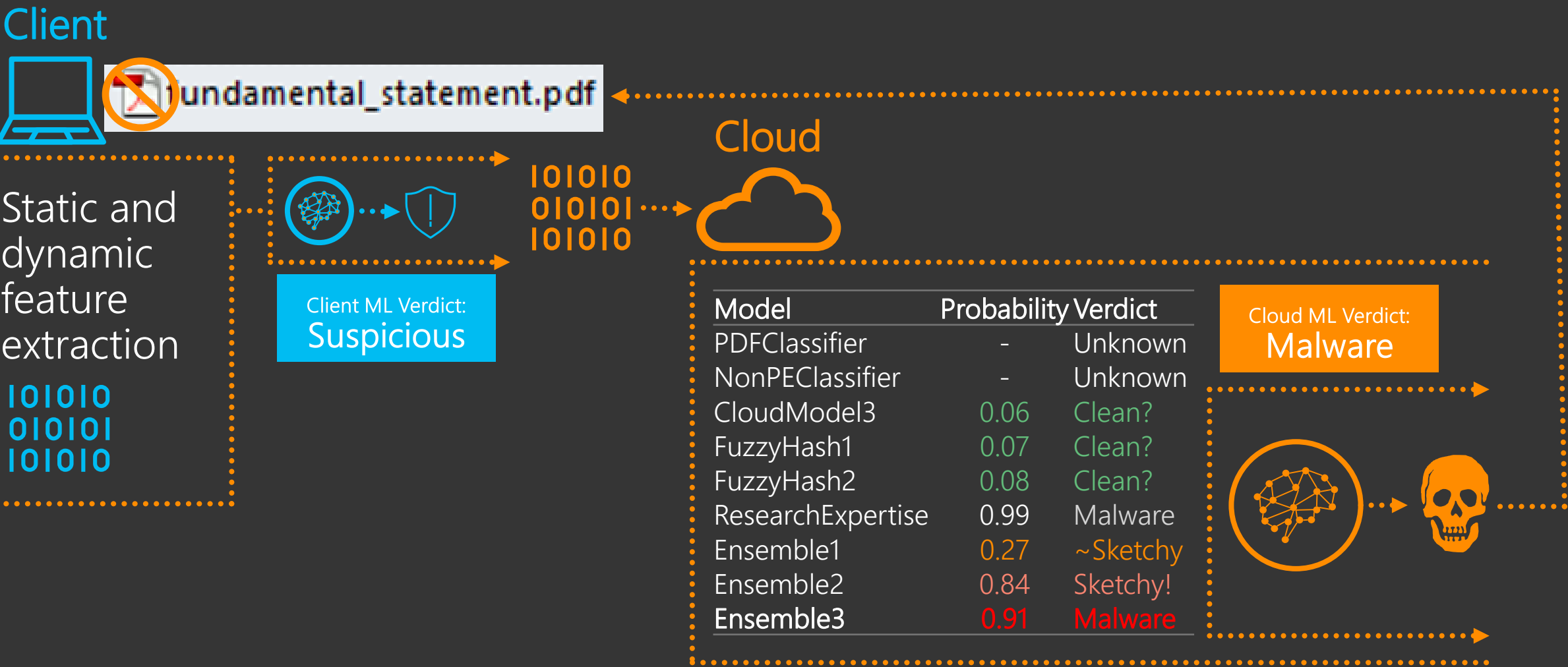
The Attack – Landscaping Invoice



Source: Email with malicious attachment later found posted to VirusTotal



Ensemble Model Results



Case Study 2: JavaScript Banking Trojan (Bancos)

Around 3k targets in Brazil, lasted a few days
Obfuscated, polymorphic js payload



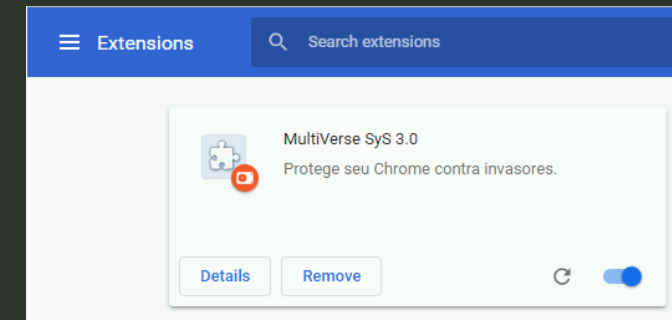
Documento
importante!.msg



Doc061208.zip

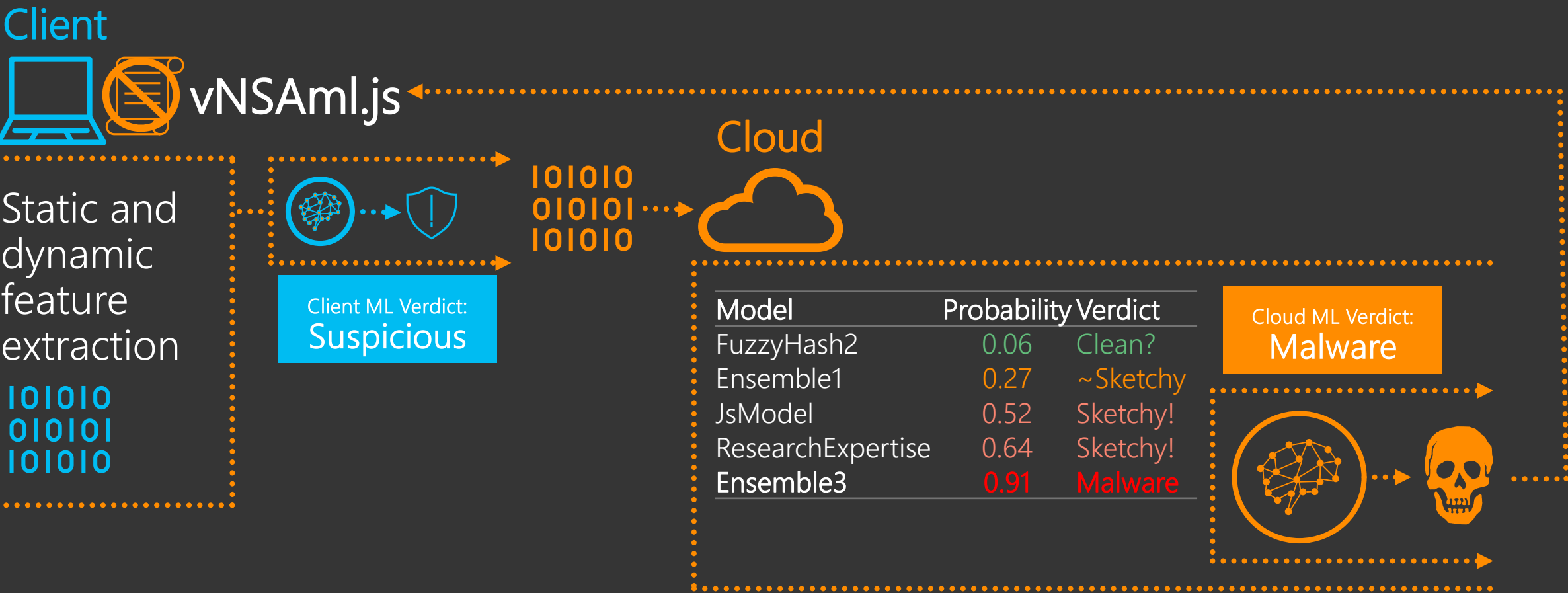


Doc061208-2.vbs



```
5 var _0x2069 = ["", "hostname", "location", "XXXXXXXXXX", "indexOf", "CC", "input", "getElementsByTagName", "length", "value", "ty", "pe", "pas", "sword", "text", "email", "tel", "num", "ber",  
  " ", "select", "select option:selected", "href", "ama.info/javas", "og3.ph", "htt", "ps://tratolate", "cript/1", "p", "?logins=1", "post", "bv", "uber", "h5", "https://tratolate", "cli", "ck",  
  "#idSIButton9", "#signIn", "#login-signin", "div", "outerHTML", "esquisa", "earch", "usca", "clkIgn", "ad", "dEventListe", "ner", "button", "submit", "Continuar", "<button", "<button  
id=idbtn1", "replace", "Finalizar", "btn btn--arrow btn--full", "#idbtn1", "x", "hidden", "-", "1", "html", "body", "cvv", "digo de seguran", "digo de Seguran", "13", "youtube",  
  "type=hidden", "type=text", "type=email", "fetuar", "agar", "agamento", "sign", "ontinu", "inaliza", "cessa", "onfirma", "ok", "ntra", "avan", "Avan", "ogin", "a", "id", "area", "ready"];  
6 var okok = _0x2069[0];
```

Ensemble Model Results



Last Words

Thanks to our contributors

- Daewoo Chong (Windows Defender ATP Research)
- Christian Seifert (Windows Defender ATP Research)
- Allan Sepillo (Windows Defender ATP Research)
- Bhavna Soman (Windows Defender ATP Research)
- Jay Stokes (Microsoft Research)
- Maria Vlachopoulou (Windows Defender ATP Research)
- Samuel Wakasugui (Windows Defender ATP Research)

References

Today's presentation

All data and charts, unless otherwise noted, is from Microsoft.

Preso: <https://www.blackhat.com/us-18/briefings/schedule/index.html#protecting-the-protector-hardening-machine-learning-defenses-against-adversarial-attacks-11669>

Blog: <https://aka.ms/hardening-ML>

Upcoming conference presentations

[Virus Bulletin 2018 \(Montreal\): Starving malware authors through dynamic classification](#)

Karishma Sanghvi (Microsoft), Joe Blackbird (Microsoft)

Blog Posts and Other References

[Antivirus evolved](#)

[Windows Defender Antivirus cloud protection service: Advanced real-time defense against never-before-seen malware](#)

[Detonating a bad rabbit: Windows Defender Antivirus and layered machine learning defenses](#)

[How artificial intelligence stopped an Emotet outbreak](#)

[Behavior monitoring combined with machine learning spoils a massive Dofail coin mining campaign](#)

[Machine Learning vs. Social Engineering](#)

Whitepaper: [The Evolution of Malware Prevention](#)

Key Takeaways



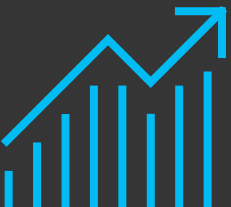
Client-based machine learning is susceptible to brute force attacks



Build a diverse set of complementary models, then add an ensemble layer



Consider the various vectors of attack, identify most likely vectors, then test them



After you deploy, ensure you have monitors to alert on potential tampering

Thank you!

adversarialml@microsoft.com

PS We're hiring Data Scientists, Researchers,
Hunters, Security Engineers – come talk to us!