# Agenda BRKSEC-2068
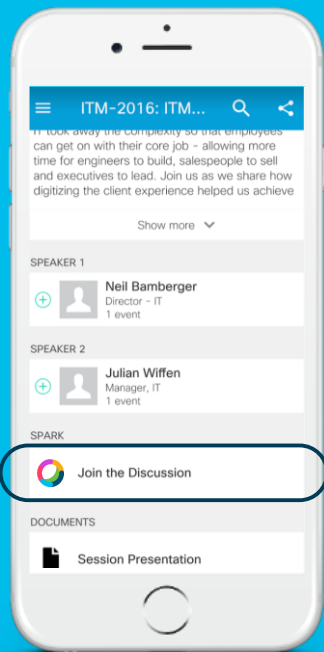
- Introduction

- Security Analytics Fundamentals

- Telemetry, Synthesis/Analytics, and Outcomes

- The Age of Artificial Intelligence & Machine Learning

- Trends and Changes That Shape the Future

- The Future of Security Analytics

- Conclusion & Takeaways

# Cisco Webex Teams
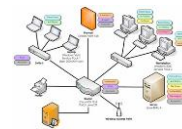
**cs.co/ciscolivebot#BRKSEC-2068**

## Questions?
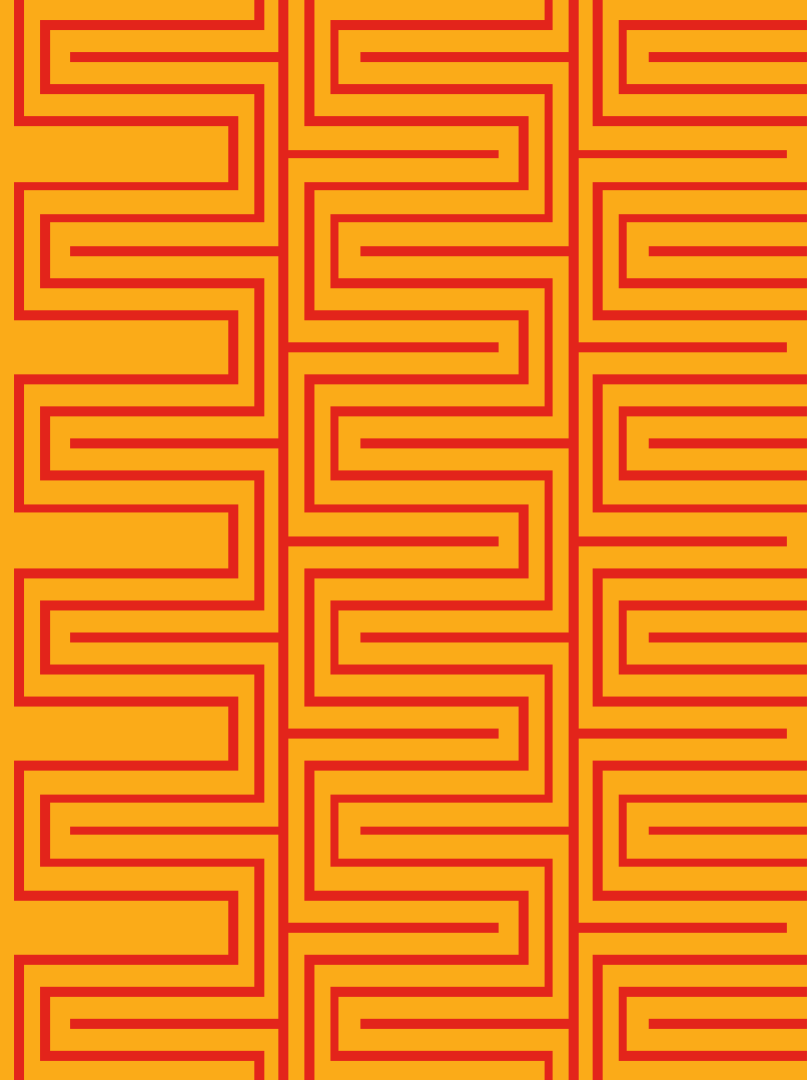Use Cisco Webex Teams (formerly Cisco Spark) to chat with the speaker after the session

## How
1. Find this session in the Cisco Events Mobile App
2. Click "Join the Discussion"
3. Install Webex Teams or go directly to the team space
4. Enter messages/questions in the team space

# Hello My Name is TK Keanini
*(Pronounced Kay-Ah-Nee-Nee)*

# Fundamentals

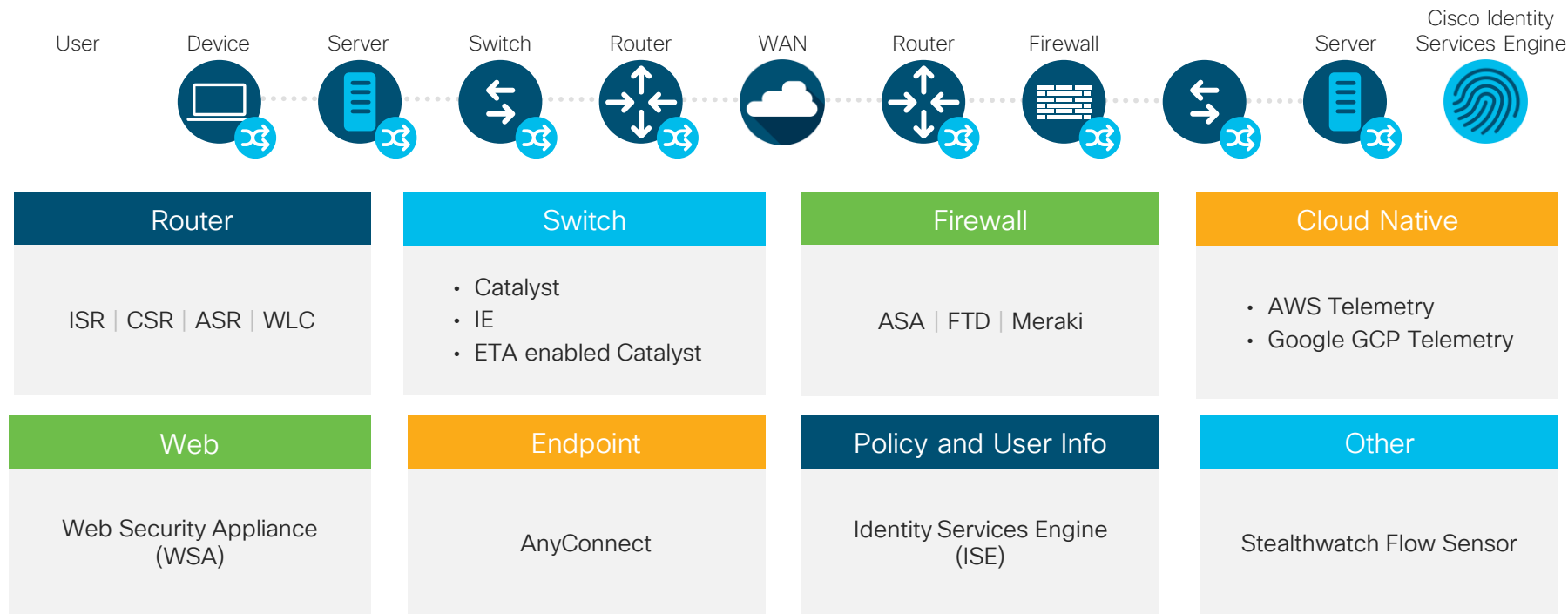# Security Analytics versus Other Analytics

**Outcomes**

**Synthesis/Analytics**

**Telemetry**

Security Analytics focus on augmenting or automating these functions

- Incident Responder
- Security Analyst
- Security Operations
- Threat Hunter
- Compliance and Policy
- Business Continuity
- Cybercrime fighting

# Telemetry (changes within an observational domain)

| User | Device | Server | Switch | Router | WAN | Router | Firewall | | Server | Cisco Identity Services Engine |

| Router | Switch | Firewall | Cloud Native |
|---|---|---|---|
| ISR \| CSR \| ASR \| WLC | • Catalyst<br>• IE<br>• ETA enabled Catalyst | ASA \| FTD \| Meraki | • AWS Telemetry<br>• Google GCP Telemetry |

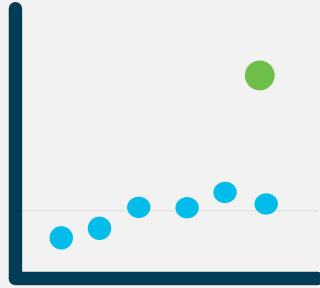| Web | Endpoint | Policy and User Info | Other |
|---|---|---|---|
| Web Security Appliance (WSA) | AnyConnect | Identity Services Engine (ISE) | Stealthwatch Flow Sensor |

## All Telemetry is Data but not all Data is Telemetry

# What Did We Do Before Machine Learning?

Simple Pattern
Matching

Statistical Methods

Rules and First
Order Logic (FoL)

Use in Combination with Machine Learning

# When to Use Machine Learning?

❌ If the domain is **static**, has **limited variability**, and is **well-understood**, then machine learning would not be needed.

✅ If the domain is **evolving**, has **a large amount of variability**, or is **not well-understood**, then we can use machine learning to either **help understand the domain** or **efficiently make predictions of unseen instances**.

# Why Use Machine Learning for Security Analytics

- Advanced Threat inherently is not static and evolving

- The data sets are often very large at scale (the 1% that matters)

- The most advanced threats are not well-understood and novel

- Machine Learning is not magic and still has problems!

**The key is to use its strengths along side other techniques in a analytics pipeline.  This makes it difficult to evade and delivers the highest fidelity!**

# Insider Threats & Behavioral Security Analytics

### Attackers
They're not breaking in, they are logging in

### Detecting
Through novelty and outliers

### Events
Turn weak signals into a strong ones

# Using the Analytical Stack to explain Encrypted Traffic Analytics
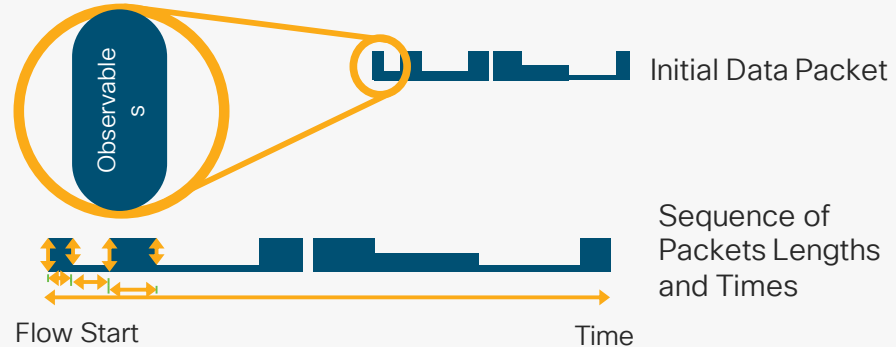
# Example: Encrypted Traffic Analytics



**Outcomes**
- Detection of Malware without Decryption
- Cryptographic Compliance

**Synthesis/Analytics**
- Analytics Pipeline of Diverse Methods

**Telemetry**
- Observables
- Initial Data Packet
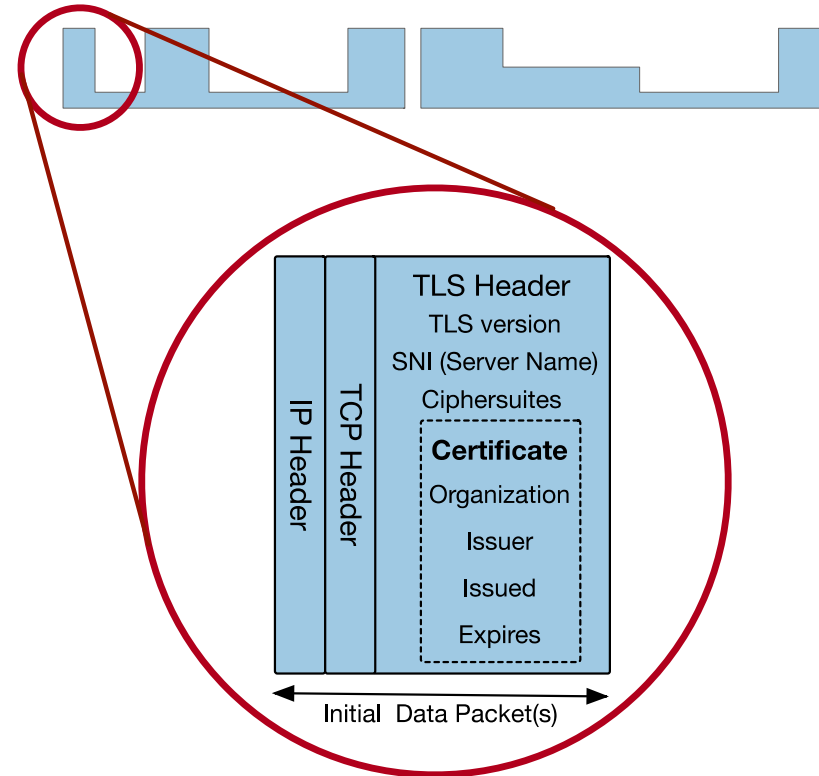- Sequence of Packets Lengths and Times
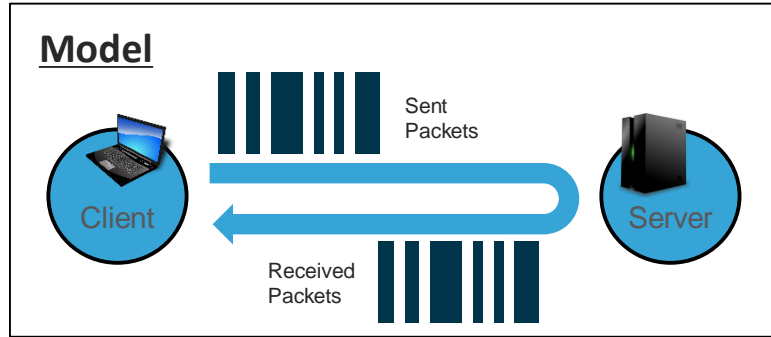- Flow Start
- Time

# Initial Data Packet (IDP)

- HTTPS header contains several information-rich fields

- Server name provides domain information

- Crypto information educates us on client and server behavior and application identity

- Certificate information is similar to *whois* information for a domain

- And much more can be understood when we combine the information with global data

Initial Data Packet



TLS Header
TLS version
SNI (Server Name)
Ciphersuites

**Certificate**
Organization
Issuer
Issued
Expires

IP Header

TCP Header

Initial Data Packet(s)

# Sequence of Packet Lengths and Times (SPLT)
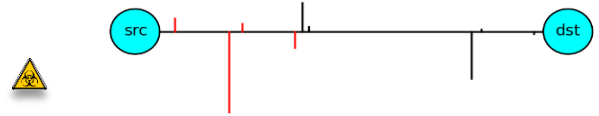
**Model**



Sent Packets

Received Packets

Client

Server

Packet lengths, arrival times and durations tend to be inherently different for malware than benign traffic.

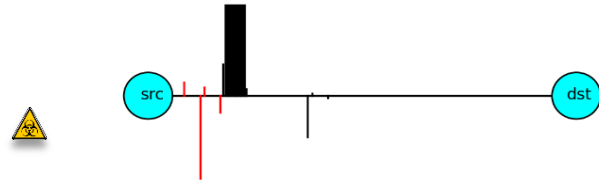**Google search Page Download**



src — dst

**Initiate Command & Control**



src — dst

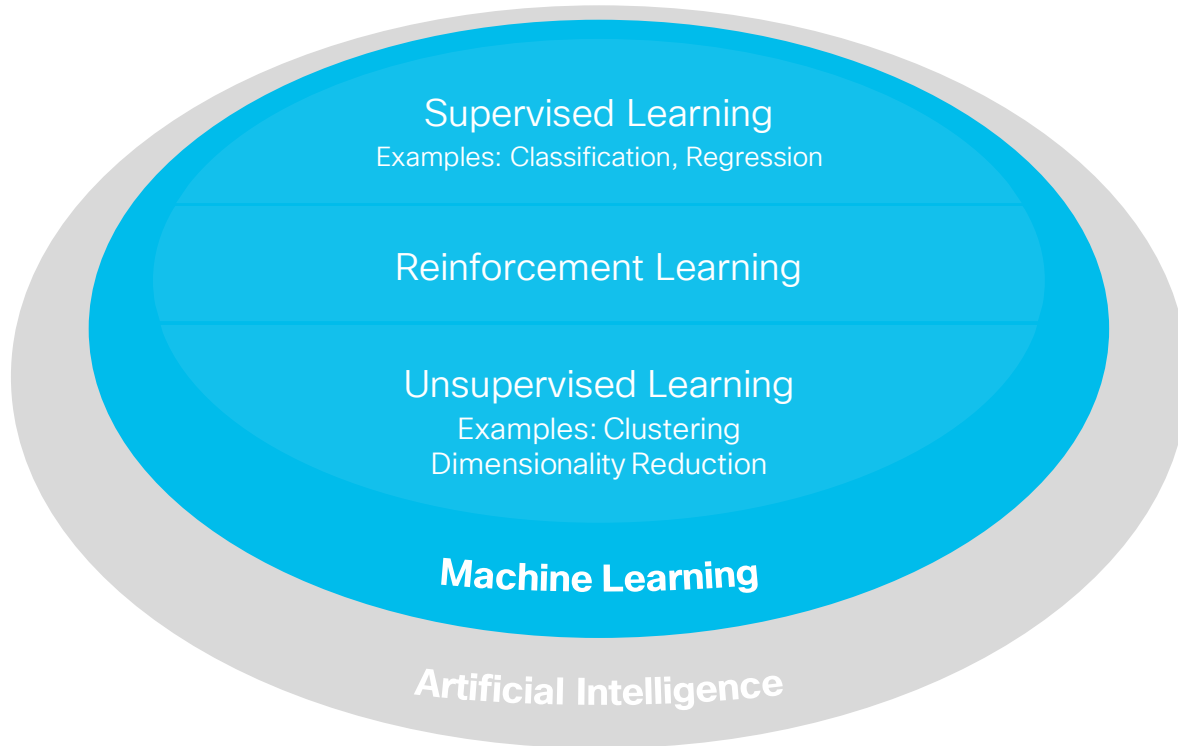**Exfiltration & Keylogging**



src — dst

# Artificial Intelligence & Machine Learning

*"Field of study that gives computers the ability to learn without being explicitly programmed."*

Arthur Samuel's definition of machine learning in 1959

# Machine Learning Big Picture

**Supervised Learning**
Examples: Classification, Regression

**Reinforcement Learning**

**Unsupervised Learning**
Examples: Clustering
Dimensionality Reduction

**Machine Learning**

**Artificial Intelligence**

Machine Learning is **one** of the fields in Artificial Intelligence, where machines learn to act autonomously, and react to new situations **without being pre-programmed**. It is about designing algorithms that allow computers to learn aimed at some outcome.

- Learn to identify faces, learn to drive a car, etc

- Learning to detect malware, learning to identify a threat actors, etc.

# Supervised

- used **when you know the question you are trying to ask**

- and **have examples of it being asked and answered correctly**

- If you can phrase a problem as **'we know this is right, learn a way to answer more questions of this type'**

# Unsupervised

- Less structured & know little about the structure

- You don't have answers and may not fully know the questions

- Unsupervised techniques act as a tool for gaining an understanding of how elements of the set relate to each other

# Reinforced Learning

- sometimes called RL and is really the 'other' category

- learns the optimal solution by repeated trial and error

- If you can formalize your problem even at a level above even what supervised learning calls for then RL has some powerful tools for solving it.

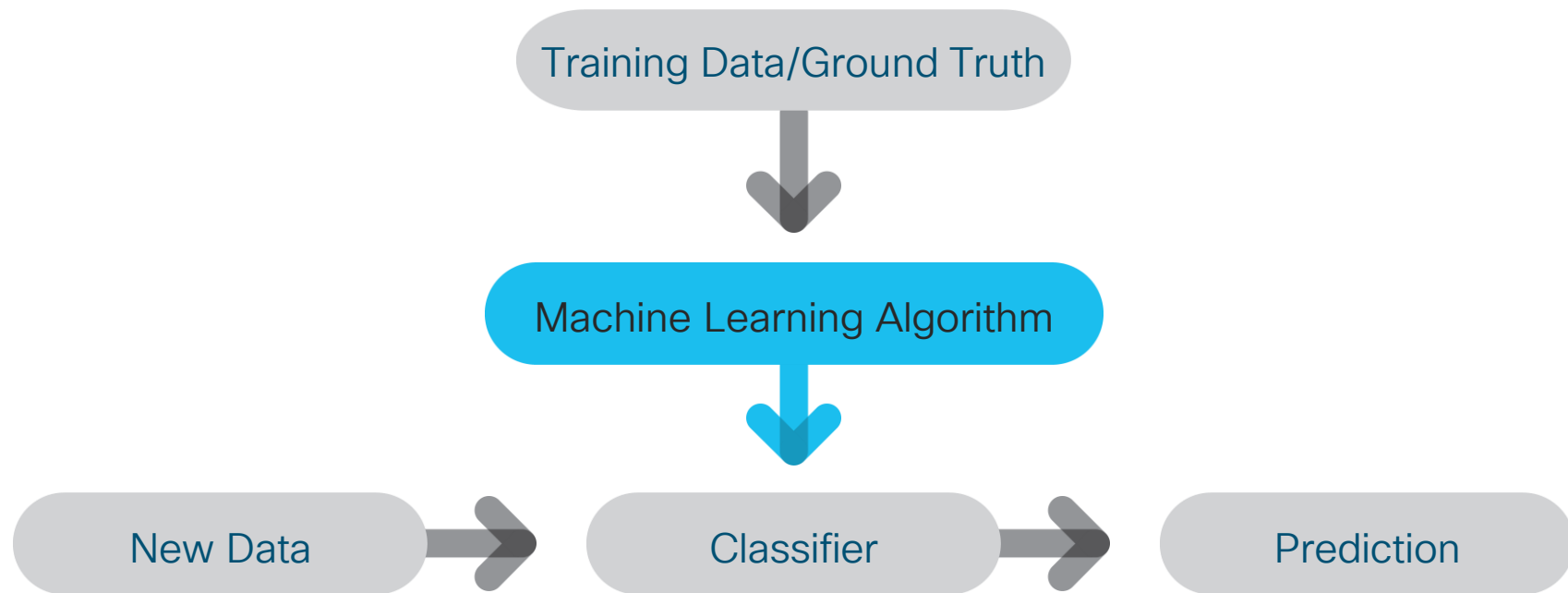# Ground Truth Used in Supervised Learning

- The **'Ground Truth'** is the pairing of example questions and answers

- If you can phrase a problem as **'we know this is right, learn a way to answer more questions of this type'**

- Success depends greatly on the dataset expressing the Question -> **Answer mapping**

*"Field of study that gives computers the ability to learn without being explicitly programmed."*

*"Field of study that gives computers the ability to be implicitly programmed."*

# Training Classifiers

# Pitfalls to avoid with Machine Learning

# How the data will explain itself?

## Regression

- Answer is a real number

- Example: given the weather conditions and temperatures from the previous 10 days, attempting to predict the exact temperature of the following day

## Classification

- Answer is a binary/n-ary set of labels

- Example: given the weather conditions and temperatures from the previous 10 days, attempting to predict rainy vs sunny vs windy vs snowing, etc. (labels)

# The Problem with **'Just'** Machine Learning...

**Hal-9000**

"I have found a threat actor operating at your branch office in Los Angeles, would you like me to remove that device from the network?"

"Yes, I would like to quarantine this device but please tell me how you arrived at this conclusion?"

**Dave**

**Hal-9000**

"I'm sorry Dave I only have the computation paths of my ML algorithms, would you still like me to remove this device from the network?"

"Don't do anything until you can share with me the logical path of your investigation!"

**Dave**

# Success is Domain Specific

Other ML Application  $\neq$  Security

# Transformational Trends

# The 5 Transformative Trends

1.  Overlay-based Networking and the Associated Control Planes

2.  Cloud Native Architectures Like: Kubernetes Service Mesh, Lambda, etc.

3.  Zero Trust Architecture

4.  Transit Inspection Opacity (TLS 1.3, HTTPS by default, etc.)

5.  TOR

While the security objectives have not changed, what we defend, how we defend it, & where we defend are all shifting

# SD-Access/SDWAN is Networking-as-a-Service



End-point flexibility:
- Physical or virtual
- Rich services or lite
- Branch, Agg, Cloud

**4**

Cloud Delivered    Analytics

**1** Cloud delivered WAN with operational simplicity & analytics

**3** Application QOE

USERS

DEVICES

THINGS

SD-WAN    **5** Cloud OnRamp    ···    Use-Cases

**WAN**

DNA Center
Policy  Automation  Analytics

Intent-based Network Infrastructure

DC

IaaS

SaaS

vDC

Apps

**0** Transport Independent WAN Fabric

**2** Superior security architecture – cloud based & on-prem

# Serverless (uber for code)

# Zero Trust Architectures

## Zero Trust Fundamentals

- Firewall enforced perimeters/zones are gone

- All hosts are treated as Internet facing (no firewalls or VPNs)

- Every device, user, and network flow is authenticated and authorized



BeyondCorp components and access flow

# ETA Data Features, <= TLS 1.2



Application Information

Server Information

Behavioral Information

client_key_exchange

change_cipher_spec

encrypted_handshake_message

app_data

app_data

client_hello

server_hello

cont.

certificate

server_key_exchange

server_hello_done

change_cipher_spec

encrypted_handshake_message

app_data

encrypted_alert

# ETA Data Features, TLS 1.3

Application Information

Server Information

Behavioral Information

**client_hello**

**server_hello**

**app_data**

**app_data**

**app_data**

**app_data**

**app_data**

**app_data**

**app_data**

**app_data**

# The Onion Router

Open source SW / public design specs

Data is constantly encrypted at multiple layers

Sent through multiple routers. Each router decrypts the outer layer and finds routing instructions

Sends the data to the next router

Result is a completely encrypted path using random routers

# How is the Tor Network built?

- The Tor network consists of relays

- Relays are just nodes where the Tor software is installed

- They build encrypted connections to other relays, forming an overlay network

- Everyone can run a Tor relay and contribute to the network...

# Tor Relay

# List of all Tor Relays

https://torstatus.blutmagie.de/



The public listed TOR relays (of which there are about 7000) lists about 70% of actual relays, the rest are intentionally withheld and not publicly listed

The remaining 30% must be computed using security analytics

Just one more thing...The Problem with Numbers

Cisco live!

It is not your fault that you don't understand this.

# Numbers Help Us Group Things



- Credit-worthiness Class

- Legal to drink / Legally drunk

- Weight Class

- Socioeconomic Class

- Age Class

Given a number, within a social context, we are able to infer membership to a set

* The terms 'Set' and 'Class' are synonymous in this presentation

# Syntax and Semantics



- Numbers digitize certain aspects of an observable domain
  - They also help ignore what is not being counted!

- Unlike the physical domain, before we can count things in the information domain, we must all agree on what is being counted.
  - The challenge is that we don't share the same domain expertise and understanding across an enterprise

- Number systems are dependent on social processes that institutionalize semantics
  - They often fall short when asked to support multiple perspectives and points of view

# Summary of the Transformational Trends

- Through the lens of these changes, will your solutions **remain effective**?

- What **new telemetry** becomes necessary and sufficient for your analytics?

- What **integrations** become deprecated or more valued?

- We no longer are able to view 'X' for evaluation, we must infer 'X'!

- While your analytical outcome may remain the same, what are you are defending and the multiplicity of telemetry will change!

# Challenges for Security Analytics We Must Solve

- Network Overlays and Observational Opacity

- Each observation point has less observables

- Serverless (securing a server when there is no server)

- Numbers fail us when we don't have stable semantics on what is being counted

# Future Security Analytics

# Direct Versus Indirect Observations

**Old Method**

Observations

Client → Server

**New Method**

App — L7 App/User
Net — L3/L4 Overlay/Underlay
App — L7 App/User

Observations

Client ↔ Server

# Late-binding Modeling to Detect Security Events

Dynamic Entity Modeling

| Collect Input | Perform Analysis | Draw Conclusions |
| --- | --- | --- |

**Collect Input**
- IP Meta Data
- System Logs
- Security Events
- Passive DNS
- External Intel
- Vulnerability Scans
- Config Changes

**Dynamic Entity Modeling**

**Perform Analysis**
- Role
- Group
- Consistency
- Rules
- Forecast

**Draw Conclusions**
- What is the role of the device?
- What ports/protocols does the device continually access?
- What connections does it continually make?
- Does it communicate internally only? What countries does it talk to?
- How much data does the device normally send/receive?

# Classify the Observable World and Infer the Rest

Threat Actor
Activity

Weird Stuff
(but not threat related)

Normal Activity

# Multi-layer Analytical Pipeline

Cascade of Specialized Layers of **Machine Learning** Algorithms

Billions of connections

Anomaly Detection and Trust Modeling

Event Classification and Entity Modeling

Relationship Modeling

- Statistical Methods
- Information-Theoretical Methods
- 70+ Unsupervised Anomaly Detectors
- Dynamic Adaptive Ensemble Creation

- Multiple-Instance Learning
- Neural Networks
- Rule Mining
- Random Forests
- Boosting
- ML: Supervised Learning

- Probabilistic Threat Propagation
- Graph-Statistical Methods
- Random Graphs
- Graph Methods
- Supervised Classifier Training

# Security that Shows its Work



New

Oct. 3

5

3

3 5

Spam tracking
#CSPM02

Oct. 4
C&C url

Oct. 15
Anomalous http

8

Information Stealer
#CDCH01

7 3

Oct. 16
Heavy
uploader
**Dropbox.com**

8 7

Oct. 25

Oct. 28
Malicious http

Recurring

8

Malware: sality
Dec. 9 | 28 days

# Serverless Security

**How Can You Secure a Server When There is No Server?**

**Serverless Computing** is a cloud computing execution model in which the cloud provider dynamically manages the allocation of machine resources (ie the servers)

# Serverless Anomaly Detection

Amazon Lambda function that normally connects to two internal resources connecting to an unexpected third

## Historical Outlier Observation ➲

One of the source's metrics deviated significantly from its historical baseline.

| Time ▾ | Source ⬍ | Time Window ⬍ | Type ⬍ | Metric ⬍ | Expected Value ⬍ | Outlier ⬍ | Probability ⬍ | Sample Size ⬍ | |
|---|---|---|---|---|---|---|---|---|---|
| 3/13/17 12:00 AM | ❗ lambda:RDSQueryLogger ▾ | 1d | device | Bytes Out | 12,097,460.313 | 142,117,719 | 0.37% | 42 | ✖ |
| 3/13/17 12:00 AM | ❗ lambda:RDSQueryLogger ▾ | 1d | device | Internal Bytes Out | 12,097,460.313 | 142,117,719 | 0.37% | 42 | ✖ |

## Static Connection Set Deviation Observation ➲

Device normally talks to a static set of (internal/external) devices, but has recently started/stopped talking to new/normal devices.

| | | | Normal Connections | | New Connections | | Lost Connections | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| Time ▾ | Source ⬍ | Type ⬍ | Set ⬍ | Count ⬍ | Set ⬍ | Count ⬍ | Set ⬍ | Count ⬍ | History Length (Days) ⬍ | |
| 3/13/17 12:00 AM | ❗ lambda:RDSQueryLogger ▾ | internal | ⓘ 10.0.10.193 ▾ , ⓘ 10.0.12.134 ▾ | 2 | 🚩 10.0.255.29 ▾ | 1 | - | 0 | 35 | ✖ |

Cisco live!

# Serverless Detection of an Unusual API Call



## AWS CloudTrail Event Observation ➔

AWS CloudTrail event reported for the device.

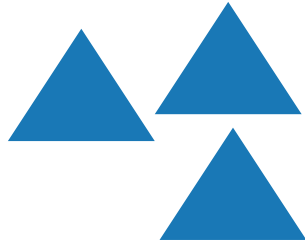| Time ▾ | Source ⬍ | Account ID ⬍ | User ⬍ | Source IP ⬍ | Event ⬍ |
|--------|----------|--------------|--------|-------------|---------|
| 3/28/17 8:23 AM | ❗ Network ▾ | 757972810156 | 👤 awslambda_963_20170328112232282 ▾ | 🇺🇸 54.91.191.63 ▾ | DeleteNetworkInterface |
| 3/26/17 12:44 PM | ❗ Network ▾ | 757972810156 | 👤 awslambda_346_20170326162935979 ▾ | 🇺🇸 54.91.191.63 ▾ | DeleteNetworkInterface |

# Serverless Behavioral Analytics
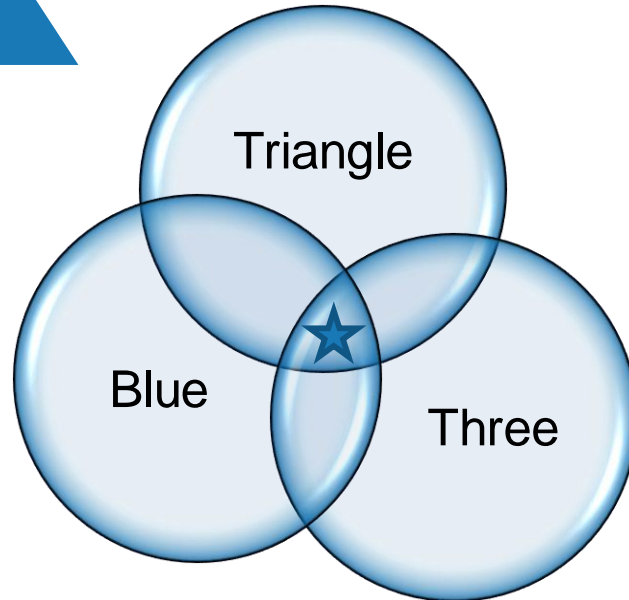
## AWS Lambda Metric Outlier Observation

An AWS Lambda function had unusual activity on one of its metrics.

| Time ▾ | Source ⇕ | Account ID ⇕ | Function name ⇕ | Metric ⇕ | Old value ⇕ | New value ⇕ |
|---|---|---|---|---|---|---|
| 3/30/17 9:00 PM | ❶ 192.168.43.147 ▾ | 23456789012 | lambda:rds-poller | Invocations | 21 | 182 |

# Thinking in Sets/Class and Membership

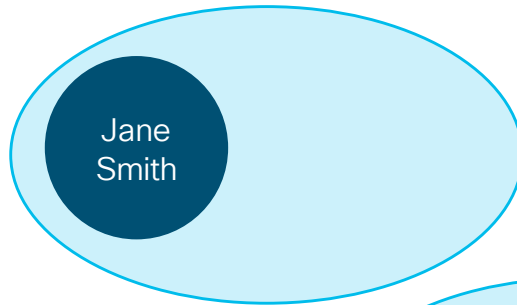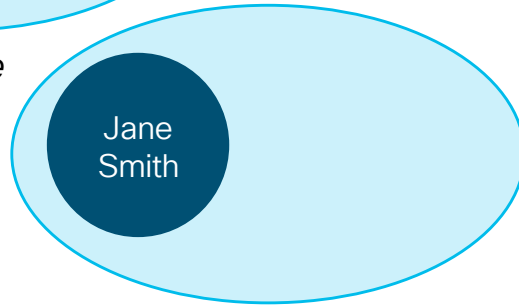There are 3 blue triangles

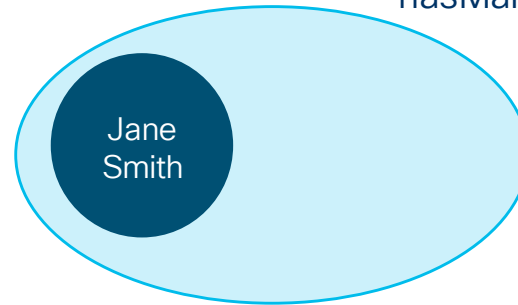...is a member of the intersection of the set Blue, the set Triangle, and the set Three

Triangle

Blue

Three

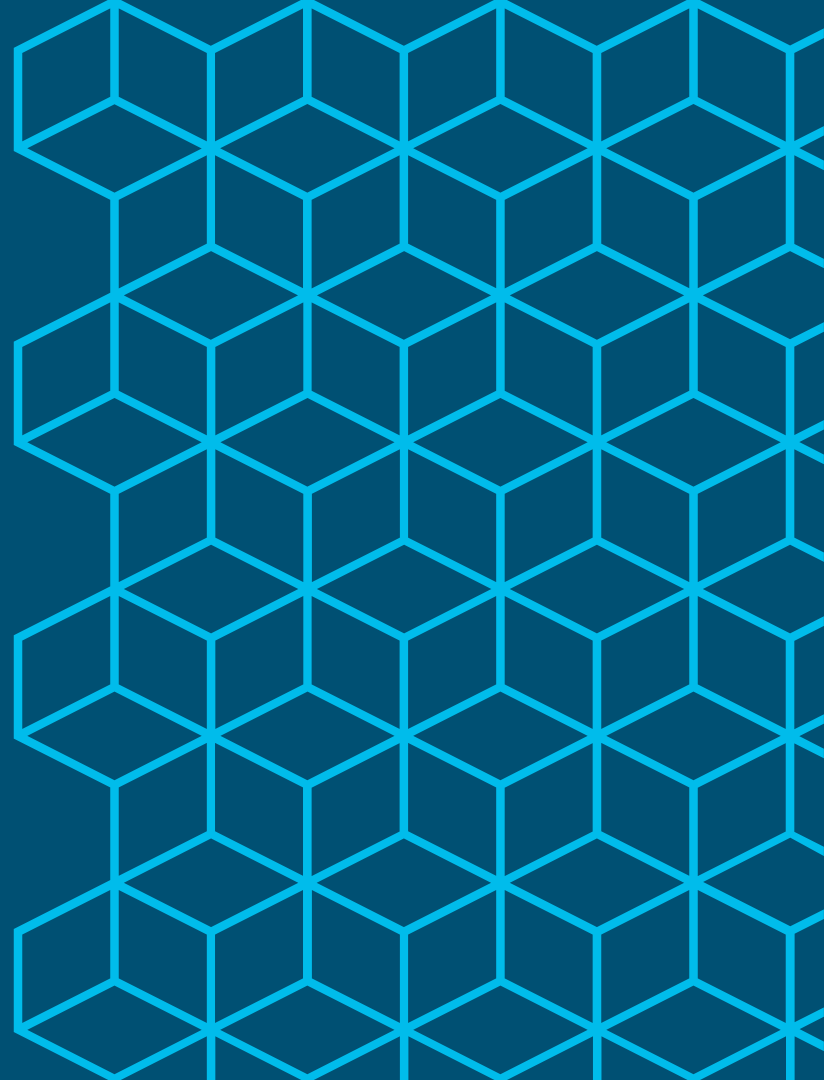# Reasoners (side step the numbers problem with first order logic)

Jane

hasMaidenName →

Smith

## Semantic Models

| *DOMAIN* | hasMaidenName | *RANGE* |
|---|---|---|
| *Female* | hasMaidenName | |
| *Married* | hasMaidenName | |
| | hasMaidenName | *SirName* |

Jane Smith

*Female*

Jane Smith

*Married*

Jane Smith

*SirName:Smith*

While syntax can be right or wrong, analytical outcomes are helpful or not helpful to you

Cisco *live!*

# How Helpful Was This Alert?



| 2018 | Stealthwatch Cloud Alerts Marked Helpful by Customers (%) |
|---|---|
| Jan | 95.91% |
| Feb | 94.52% |
| Mar | 94.75% |
| **Q1 (Jan-Mar)** | **94.45%** |
| Apr | 97.23% |
| May | 94.97% |
| Jun | 91.70 |
| **Q2 (Apr-Jun)** | **94.63%** |

# What to Ask Your Vendor

How are you applying Machine Learning in your product and why?

How do you measure its effectiveness?

Regarding supervised learning, what are you using for 'ground truth'?

What non-machine learning are you using and why?

What papers or open-source have you published regarding your analytics?

For the ML based assertions, what entailments are provided?

# Closing Thoughts

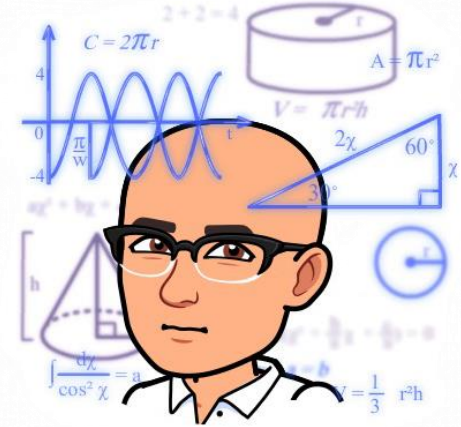✅ Be Pragmatic

✅ Provide Entailments

✅ Analytical pipeline, over single technique

✅ Measure helpfulness, not mathematical accuracy

✅ Be Transparent with your science, publish papers and open source

# Recommended Sessions
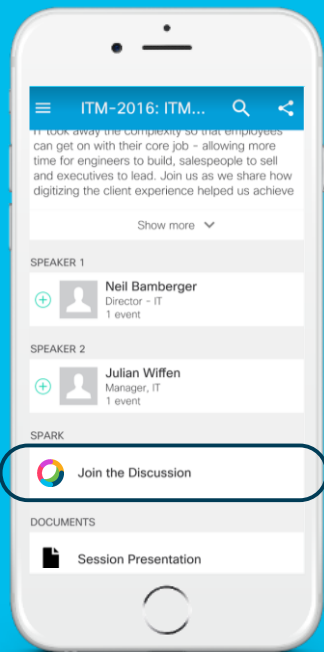
World of Solutions – Cloud Protect booth (Stealthwatch with Kubernetes & Serverless Security Demo)

World of Solutions – SOC and ThreatWall (Encrypted Traffic Analytics Live)

BRKSEC-3014 – Security Analytics with Stealthwatch: Operationalising Visibility and Machine Learning – Matt Robertson – Friday, Feb 1, 9:00 AM – 11:00 AM

BRKSEC-2323 – Claim Jumpers: Dealing with Illicit Bitcoin Miners – Matt Robertson – Thursday, Jan 31, 2:30 PM – 4:00 PM

# Cisco Webex Teams

## Questions?
Use Cisco Webex Teams (formerly Cisco Spark)
to chat with the speaker after the session

## How

1 Find this session in the Cisco Events Mobile App

2 Click "Join the Discussion"

3 Install Webex Teams or go directly to the team space

4 Enter messages/questions in the team space

cs.co/ciscolivebot#BRKSEC-2068

# Complete your online session survey

- Please complete your Online Session Survey after each session

- Complete 4 Session Surveys & the Overall Conference Survey (available from Thursday) to receive your Cisco Live T-shirt

- All surveys can be completed via the Cisco Events Mobile App or the Communication Stations

Don't forget: Cisco Live sessions will be available for viewing on demand after the event at ciscolive.cisco.com
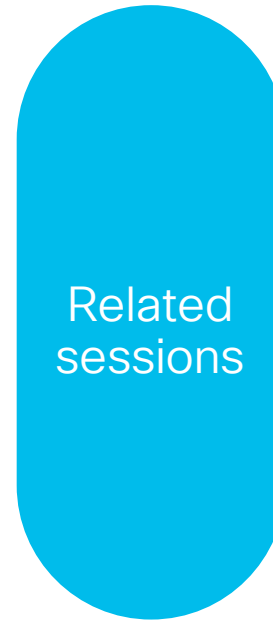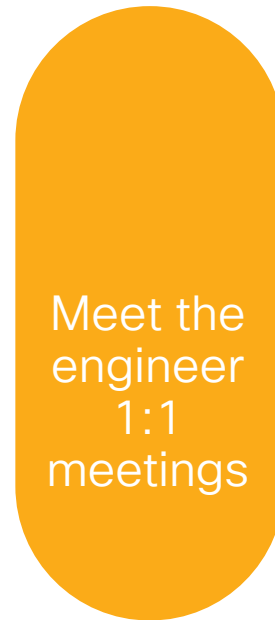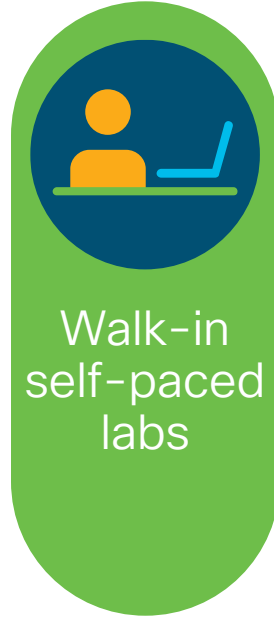
Five   * * * * *
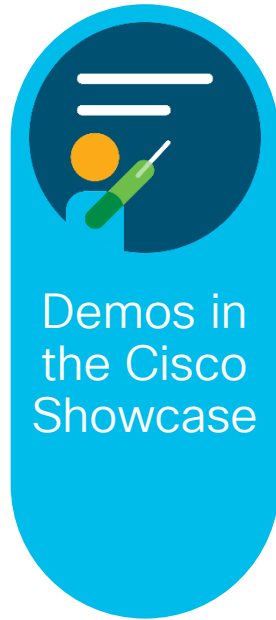
5

PLEASE!

# Continue Your Education

**Demos in the Cisco Showcase**

**Walk-in self-paced labs**

**Meet the engineer 1:1 meetings**

**Related sessions**

Thank you

# References

# Learn More….

- [Cisco Stealthwatch Enterprise](#)

- [Cisco Stealthwatch Cloud](#)

- [Encrypted Traffic Analytics](#)

# Basic References

- Blog: Detecting Encrypted Malware Traffic (Without Decryption)

- Blog: Learning Detectors of Malicious Network Traffic

- Blog: Transparency in Advanced Threat Research

- Blog: Turn Your Proxy into Security Device

- Blog: Securing Encrypted Traffic on a Global Scale

- Blog: Closing One Learning Loop: Using Decision Forests to Detect Advanced Threats

# Make Your Head Hurt Reading Material

- Identifying Encrypted Malware Traffic with Contextual Flow Data, Blake Anderson and David McGrew, AISEC '16

- Grill, M., Pevny, T., & Rehak, M. (2017). Reducing false positives of network anomaly detection by local adaptive multivariate smoothing. Journal of Computer and System Sciences, 83(1), 43-57.

- Komarek, T., & Somol, P. (2017). End-node Fingerprinting for Malware Detection on HTTPS Data. In Proceedings of the 12th International Conference on Availability, Reliability and Security (p. 77). ACM.

- Jusko, J., Rehak, M., Stiborek, J., Kohout, J., & Pevny, T. (2016). Using Behavioral Similarity for Botnet Command-and-Control Discovery. IEEE Intelligent Systems, 31(5), 16-22.

- Bartos, K., & Rehak, M. (2015). IFS: Intelligent flow sampling for network security–an adaptive approach. International Journal of Network Management, 25(5), 263-282.

- Letal, V., Pevny, T., Smidl, V. & Somol, P. (2015). Finding New Malicious Domains Using Variational Bayes on Large-Scale Computer Network Data. In NIPS 2015 Workshop: Advances in Approximate Bayesian Inference (pp. 1-10).

- Rehak, M., Pechoucek, M., Grill, M., Stiborek, J., Bartoš, K., & Celeda, P. (2009). Adaptive multiagent system for network traffic monitoring. IEEE Intelligent Systems, 24(3).