# Leveraging other data sources with flow to identify anomalous network behavior

Peter Mullarkey, Peter.Mullarkey@ca.com
Mike Johns, Mike.Johns@ca.com
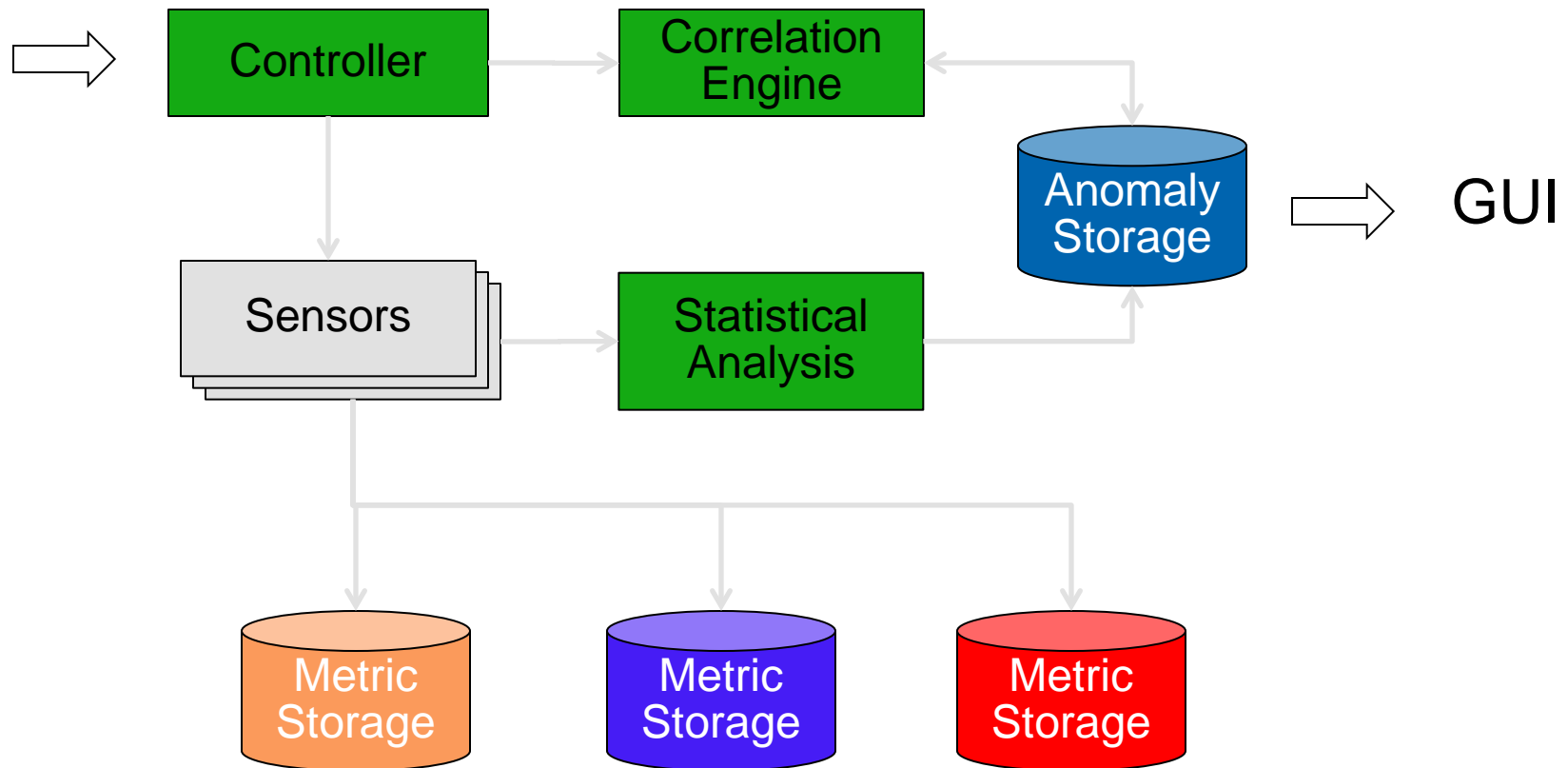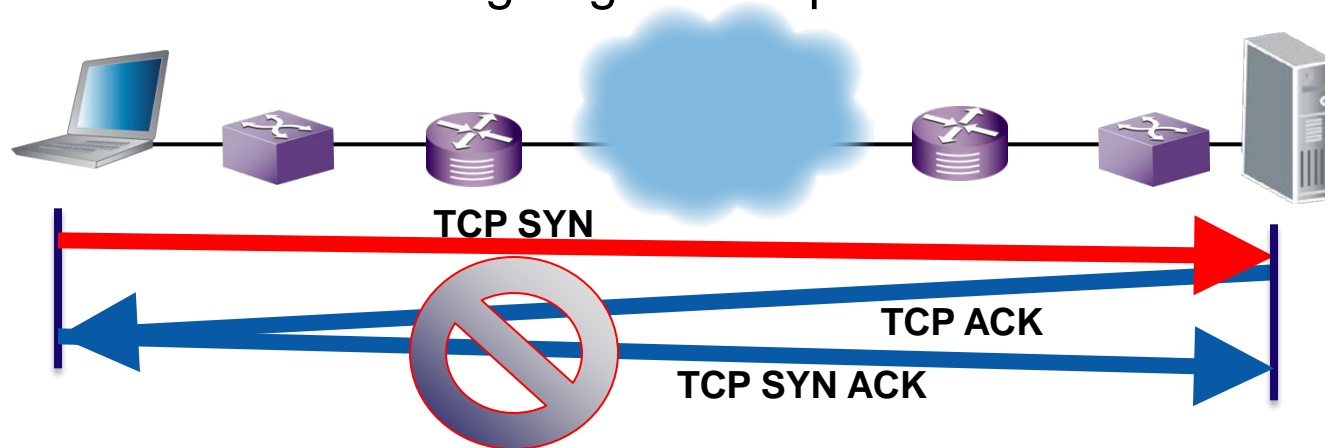Ben Haley, Ben.Haley@ca.com

FloCon 2011

ca
technologies

# Goal and Approach

— Goal:  Create high quality events without sacrificing scalability

— Approach:  Create a system that
- Is more abstract than a signature-based approach
- Leverages domain knowledge more than a pure statistical approach
- Makes use of all available data to increase event quality
- Relies only on readily available data – no new collection

# Sensors

— *Sensors* are a level of abstraction above signatures

  – leveraging knowledge of network behavior

— Sensors describe behavior to watch for

  – Is this host contacting more other hosts than usual?

  – Is this host transmitting large ICMP packets?



**TCP SYN**

**TCP ACK**

**TCP SYN ACK**

— Sensors can be created and modified in the field

# Example Sensors

— SYN-only Packet Sources
 − Looking at flows with SYN as the only flag. SYN flood, denial of service attack, worm infection

— High Packet Fan Out
 − Looking at hosts talking to many more peers tan usual. Virus or worm infection

— Large DNS and/or ICMP Packet Sources
 − Looking at  volume/packet, compared to typical levels for these protocols. Data ex-filtration – discretely attempting to offload data from internal network to an external location

— TTL Expired Sources
 −      Network configuration issue – routing loops, heavy trace route activity

— Previously Null Routed Sources
 − Traffic discovered from hosts that have had previous traffic null routed

# Example Sensor (non-Flow data sources)

— **Incoming Discard Rate**

The Incoming Discard Rate sensor look for patterns where incoming packets were dropped even though they contained no errors. Can be caused by: Overutilization, Denial of service, or VLAN misconfiguration

— **Voice Call DoS**

This sensor looks for patterns where a single phone is called repeatedly over a short period of time. This type of attack differs from other Denial of Service (DoS) attacks and traditional IDS may not catch it because it is so low volume. It only takes about 10 calls per minute or less to keep a phone ringing all the time.

— **Packet Load**

This sensor looks for a pattern in bytes per packet to server. Applications running on servers generally have a fairly constant ratio between the number of packets they receive in requests for their service and the volume of those packets. This sensor looks for anomalous changes in that ratio.

# SQL Interface to Metric Data (including flow)

— Very helpful for exploring the data – to look for interesting patterns, and develop sensors

— *Example: top talkers (by flows)*

```
SELECT srcaddr as source,
        count(*) as flowsPerSrc,
        count(*)/ ((max(timestamp) - min(timestamp)) / 60 ) as avgPerMin
FROM AHTFlows
        group by source order by flowsPerSrc desc limit 10
```

# SQL Interface to Metric Data (including flow)

— *More in-depth example: looking at profiling SSL traffic (as a basis for identifying exfiltration)*

*Select    inet_ntoa(srcaddr) as srcHostAddr,  count(if(dstport = 443, inbytes, 0)) as samples,*

*count(distinct(dstAddr)) as numOfDestsPerSrcHost,*

*min(if(dstport = 443, inbytes/inpkts, 0)) as minBytesPerPacketPerSrcHost,*

*avg(if(dstport = 443, inbytes/inpkts, 0)) as avgBytesPerPacketPerSrcHost,*

*std(if(dstport = 443, inbytes/inpkts, 0)) as stdBytesPerPacketPerSrcHost,*

*max(if(dstport = 443, inbytes/inpkts, 0)) as maxBytesPerPacketPerSrcHost,*

*sum(if(dstport = 443, inbytes, 0)) / sum(inbytes)as sslRatioPerSrcHost,*

*group_concat(inet_ntoa(dstAddr)) as destAddrsPerSrcHost*

*from AHTFlows  where  protocol = 6  and timestamp > (unix_timestamp(now()) - 30*60)*

*group by hostAddr  having sslBytes > 0  and numOfDestsPerSrcHost < 10*

*order by sslBytes desc*

# Correlation Engine

— Multiple anomaly types for the same monitored item within the same time frame combine into a *correlated anomaly*

— These can span data from disparate sources
  – NetFlow, Response Time, SNMP, etc

— An index is calculated that aids in ranking the correlated anomalies

# Types of Problems Found

The developed system has found issues that are beyond single issue description

— Spreading Malware

— Router overload causing server performance degradation (Example #1)

— Data exfiltration

— Interface drops causing downstream TCP retransmissions

— Unexpected applications on the network (Example #2)

ca
technologies

# Customer Example 1: Unexpected Performance Degradation

## Ny1-x.x.100.52

# Customer Example 1: Unexpected Performance Degradation

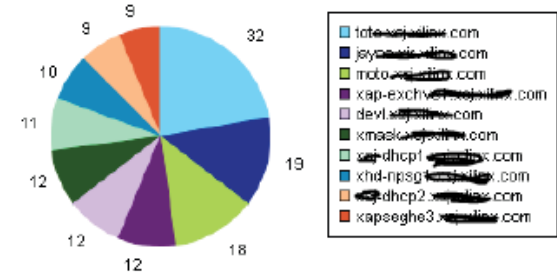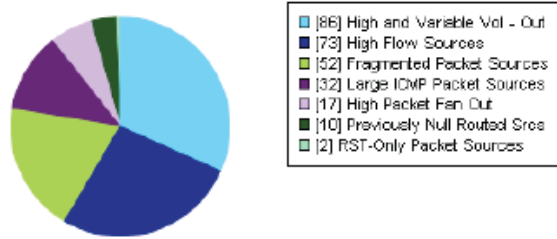# Customer Example 2: What is really happening on your network?

High quality anomalies can be found without sacrificing scalability

—Key aspects

- − Embodying domain knowledge in sensors
- − Leveraging statistical analysis approach, separating domain knowledge from data analysis
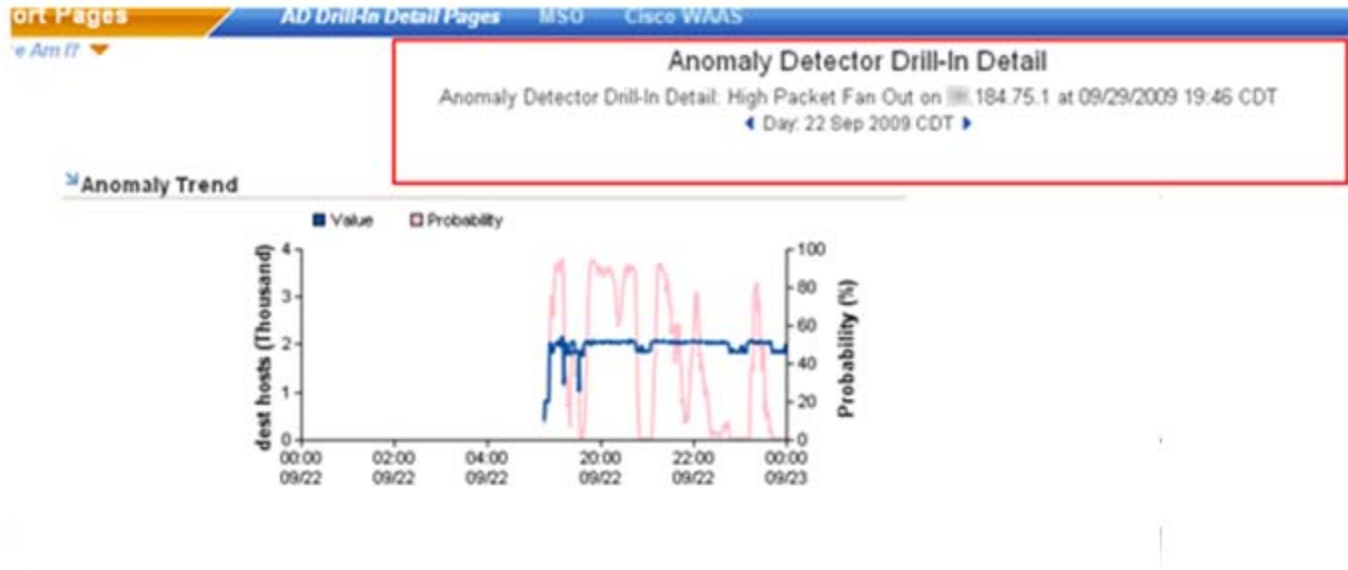- − Using simple, fast event correlation

Effectiveness of approach has been shown by solving customer problems on real networks

# Questions?

# Backup Slides

— Extra info slides

# Customer Example 3: Malware Outbreak

# Customer Example 3: Malware Outbreak



Report Results

| Router Addr | Interface In | IP Protocol | Src Addr | Src Port | Dest Addr | Dest Port | ToS | Bytes | Rate (Bits) | % Total (Bytes) | Flows | Flow Duration | Pkts | Rate (Pkts) | % Total |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 34.184.253.73 | Gi0/0 | tcp (6) | 2D0C94D5B (184.75.43) | 1025 | 154.119.202.101 | 445 | Default Traffic (0) | 96 Bytes | 0 bps | <1.00 % | 1 | 2 secs 980 ms | 2 | 0.00 pkts/s | < |
| 34.184.253.73 | Gi0/0 | tcp (6) | 2D0C94D5B (184.75.43) | 1025 | 159.115.225.122 | 445 | Default Traffic (0) | 96 Bytes | 0 bps | <1.00 % | 2 | 0 ms | 2 | 0.00 pkts/s | < |
| 34.184.253.73 | Gi0/0 | tcp (6) | 2D0C94D5B (184.75.43) | 1025 | 160.62.141.9 | 445 | Default Traffic (0) | 96 Bytes | 0 bps | <1.00 % | 2 | 0 ms | 2 | 0.00 pkts/s | < |
| 34.184.253.73 | Gi0/0 | tcp (6) | 2D0C94D5B (184.75.43) | 1025 | 172.113.161.12 | 445 | Default Traffic (0) | 96 Bytes | 0 bps | <1.00 % | 1 | 3 secs 4 ms | 2 | 0.00 pkts/s | < |
| 34.184.253.73 | Gi0/0 | tcp (6) | 2D0C94D5B (184.75.43) | 1025 | 187.121.126.32 | 445 | Default Traffic (0) | 96 Bytes | 0 bps | <1.00 % | 1 | 2 secs 972 ms | 2 | 0.00 pkts/s | < |
| 34.184.253.73 | Gi0/0 | tcp (6) | 2D0C94D5B (184.75.43) | 1025 | 221.78.132.120 | 445 | Default Traffic (0) | 96 Bytes | 0 bps | <1.00 % | 2 | 0 ms | 2 | 0.00 pkts/s | < |
| 34.184.253.73 | Gi0/0 | tcp (6) | 2D0C94D5B (184.75.43) | 1026 | 15.11.111.124 | 445 | Default Traffic (0) | 96 Bytes | 0 bps | <1.00 % | 1 | 2 secs 916 ms | 2 | 0.00 pkts/s | < |
| 34.184.253.73 | Gi0/0 | tcp (6) | 2D0C94D5B (184.75.43) | 1026 | 21.124.185.115 | 445 | Default Traffic (0) | 96 Bytes | 0 bps | <1.00 % | 1 | 2 secs 952 ms | 2 | 0.00 pkts/s | < |
| 34.184.253.73 | Gi0/0 | tcp (6) | 2D0C94D5B (184.75.43) | 1026 | 34.185.153.43 | 445 | Default Traffic (0) | 96 Bytes | 0 bps | <1.00 % | 2 | 0 ms | 2 | 0.00 pkts/s | < |
| 34.184.253.73 | Gi0/0 | tcp (6) | 2D0C94D5B (184.75.43) | 1026 | 34.185.183.179 | 445 | Default Traffic (0) | 96 Bytes | 0 bps | <1.00 % | 1 | 2 secs 928 ms | 2 | 0.00 pkts/s | < |

# Customer Example 4: Retransmissions traced back

# Statistical Analysis Methodology

— Define *anomaly* as a sequence of improbable events

— Derive the probability of observing a particular value from (continually updated) historical data

- Example
  - Under normal circumstances values above the 90th percentile occur 10 percent of the time

— Use Bayes' Rule to determine the probability that a sequence of events represents anomalous behavior
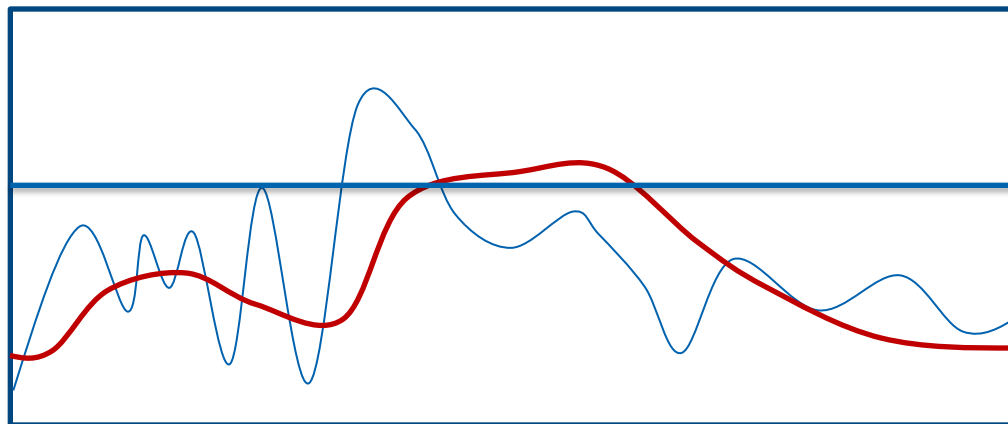
$$p(anomaly \mid point) = \frac{p(point \mid anomaly) * p(anomaly)}{p(point)}$$

Thresholding directly off of observations is difficult



We wanted an approach that could take both time and degree of violation into account, so we threshold on probability
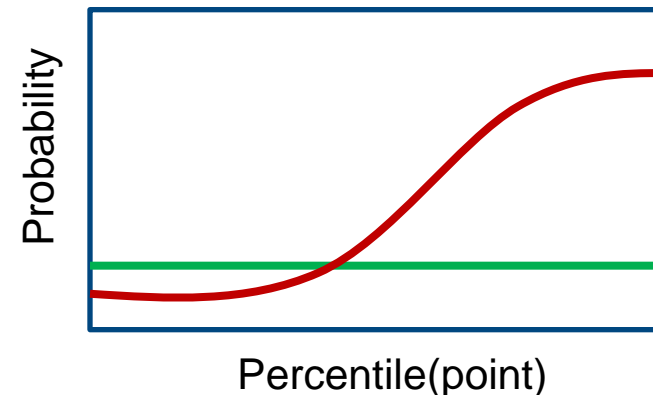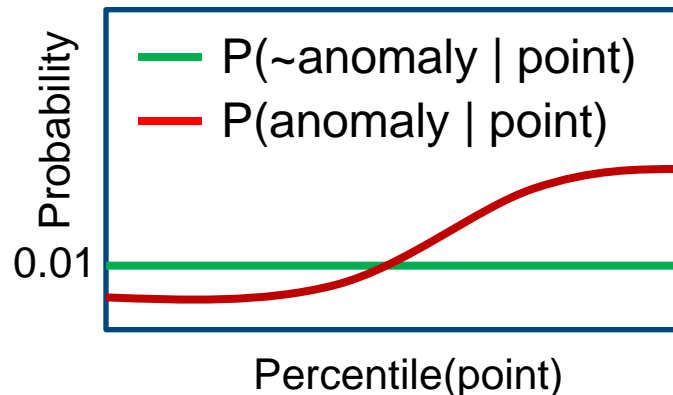
$$p(anomaly \mid point) = \frac{p(point \mid anomaly) * p(anomaly)}{(p(point \mid anomaly) * p(anomaly)) + (p(point \mid \sim anomaly) * p(\sim anomaly))}$$

**p(anomaly)** is the prior probability – either some starting value or the output from last time

**p(point|anomaly)** & **p(point|~anomaly)** are given by *probability mass functions – and are the basis for our customizable, pluggable engines*

Less Scalable
Higher Quality Events

More Scalable
Lower Quality Events

"Behavior
Analysis"

Per-metric thresholds

Baselining

Intrusion Detection Systems

Virus Scanners

Packet Inspection

Signature-Based

Statistical Methods