

Uncovering Priority Anomalies Using Pattern Discovery as a Roadmap for Contextual Analysis

Thomas Henretty
henretty@reservoir.com

FloCon 2020
Savannah, GA
9 January 2018

Reservoir Labs
New York, NY
www.reservoir.com

Presentation Outline

Part 1: Background

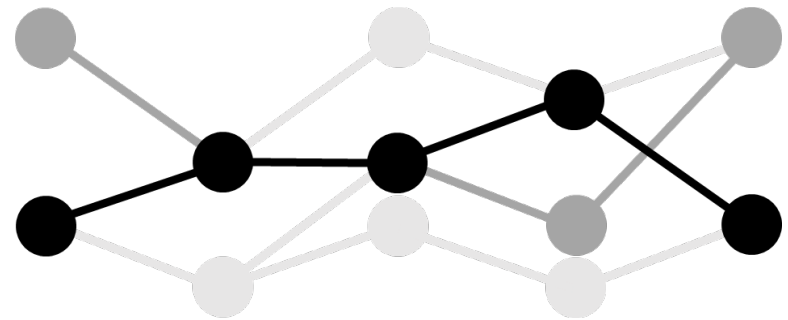
- Tensor Decomposition Basics
- Pattern Discovery in Network Flows
- MITRE ATT&CK Framework

Part 2: Anomaly Ranking

- Decompositions as Documents
- Topic Modeling for Anomaly Ranking
- Other Techniques

Part 3: Graphs and Databases

- Constructing a Targeted Query



Pattern Discovery

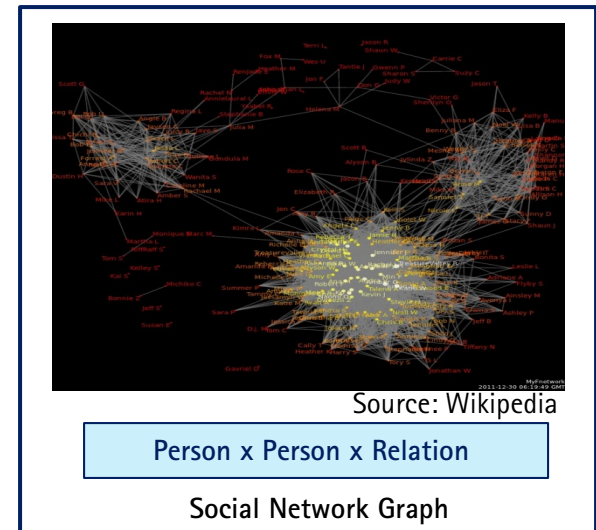
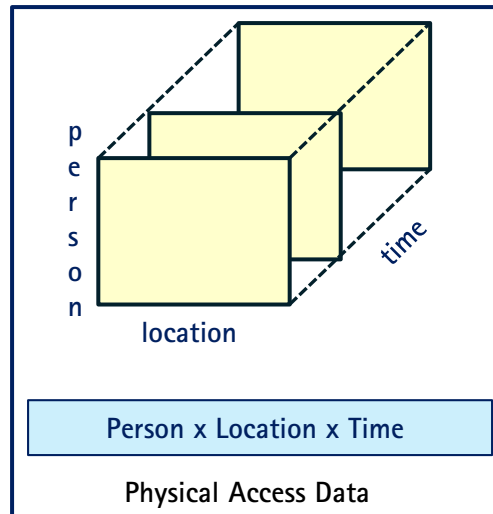
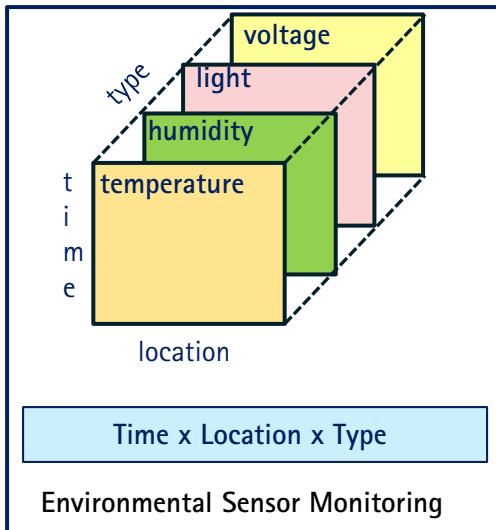
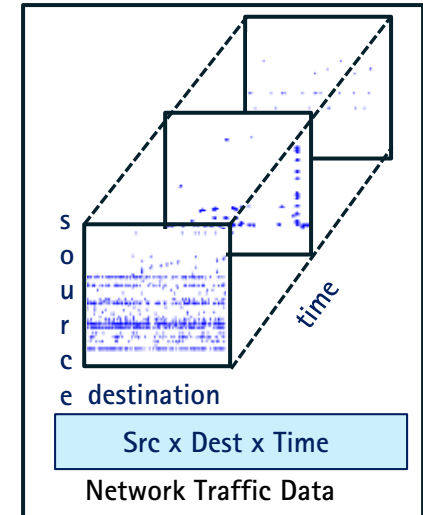
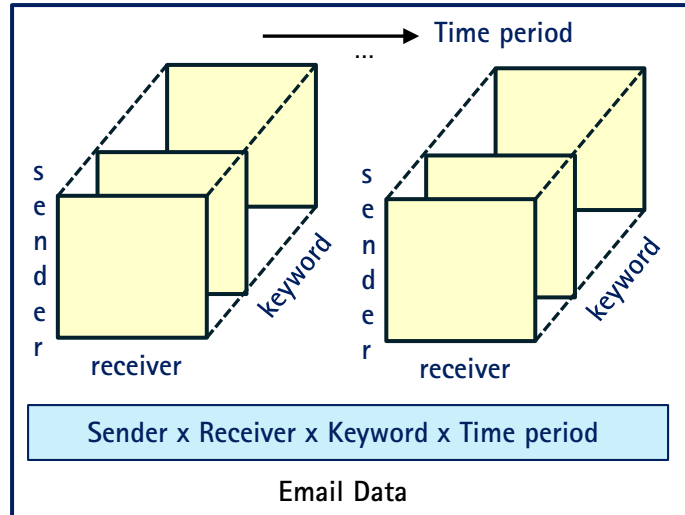
Tensor decomposition provides a model for Zeek log data that allows behaviors to be separated as coherent patterns

PART 1: BACKGROUND

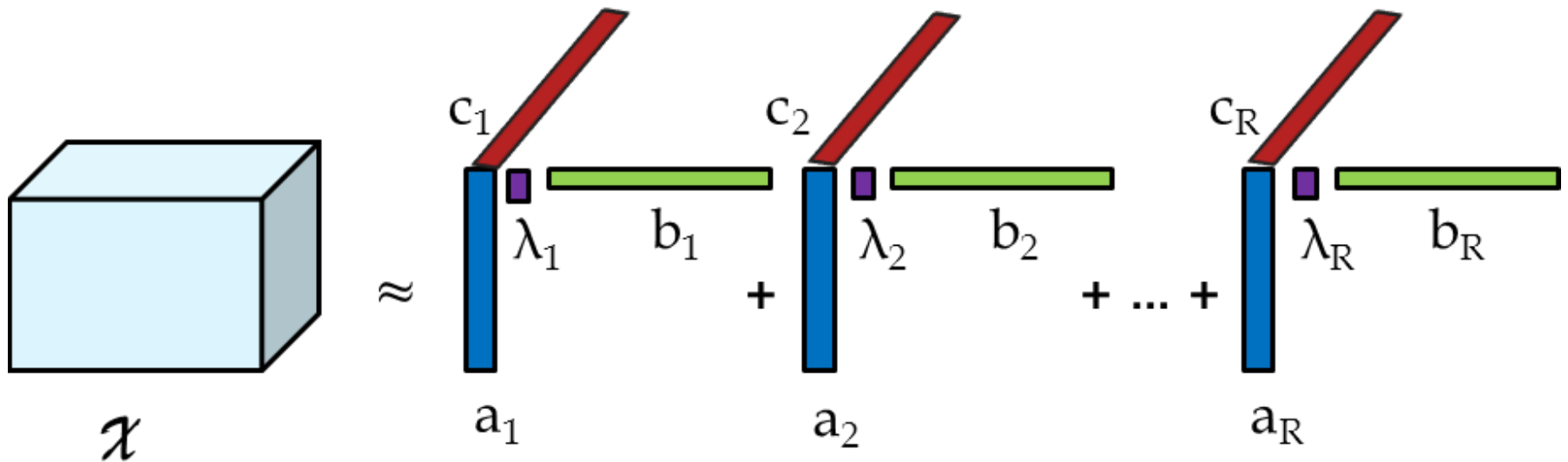
Tensors: Representing Multidimensional Data

Real World Data

- Multidimensional
- Heterogeneous
- Large
- Sparse



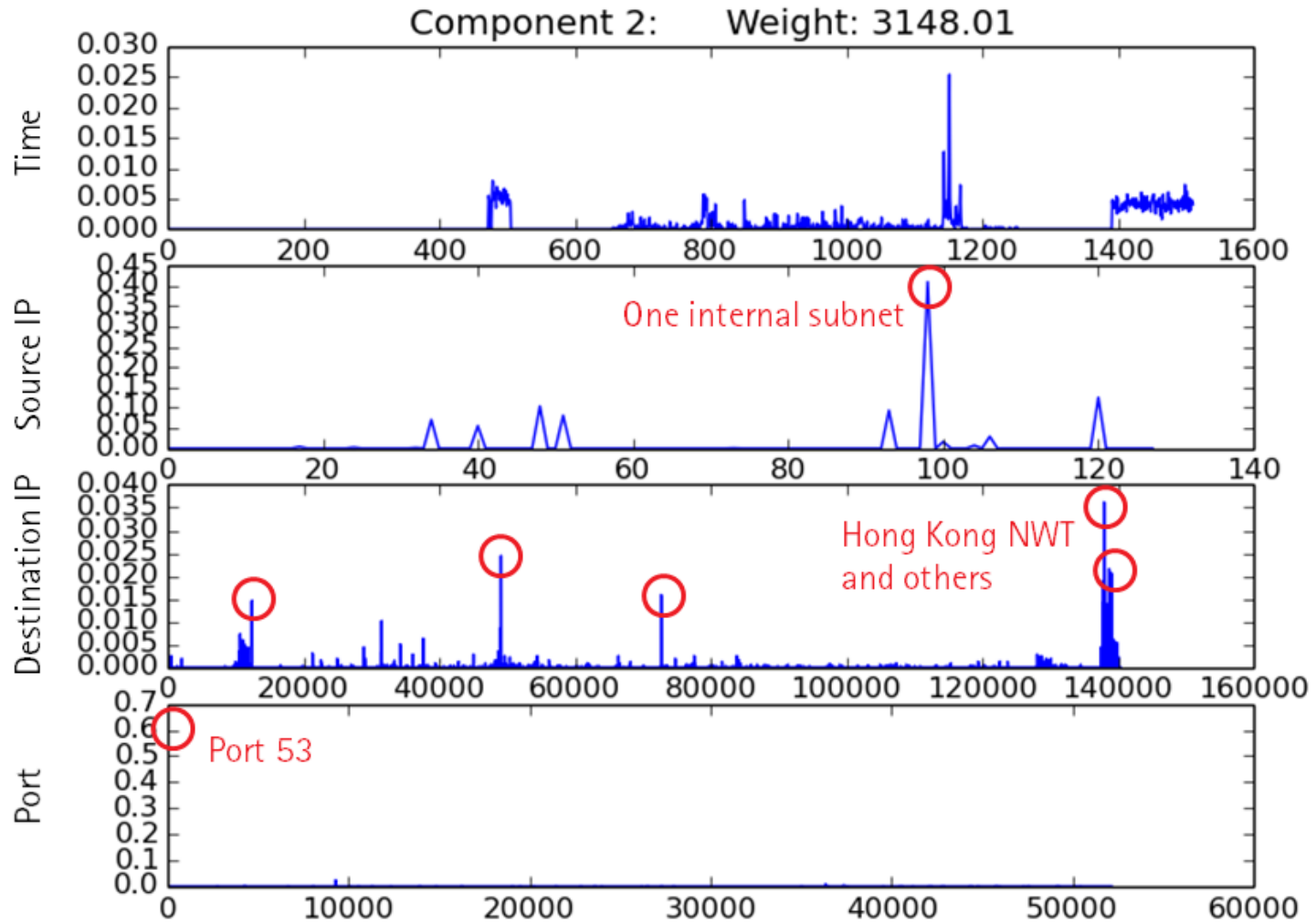
Basic CP Tensor Decomposition



CP tensor decomposition

- Multidimensional analog to matrix factorization
- Break tensor into R components
- Components represent correlated data (quantitatively)
- Can reconstruct tensor from subset of components

Example Component: Suspicious DNS Traffic



Time x Source IP x Destination IP x Port

Tensor Library for Cybersecurity

Tensor Name	Bro Log	Tensor Dimensions
Connections	The connections log (conn.log)	Time x Sender IP x Receiver IP x Port
Outgoing	Connections log entries with local sender and external receiver	Time x Sender IP x Receiver IP x Port
Incoming	Connections log entries with local receiver and external sender	Time x Sender IP x Receiver IP x Port
Time Independent	The connections log	Sender IP x Receiver IP x Port x Connection State
File Transfer	The file transfer log (files.log)	Time x Sender IP x Receiver IP x MIME-Type
HTTP	The HTTP traffic log (http.log)	Time x Sender IP x Receiver IP x URI x User Agent
DNS Query	All queries from the DNS log (dns.log)	Time x Sender IP x Receiver IP x Query x Query Type

Tensor Decompositions in MITRE ATT&CK

Relevant techniques in the MITRE ATT&CK framework

- Depends on data decomposed
- Focus on network flows
 - **Netflow** – Techniques detected via Netflow/Enclave Netflow
 - **Zeek logs** – Netflow + Network Protocol Analysis + Network Intrusion Detection

Relevant tactics


- When decomposing Zeek logs ...
 - **Initial Access** (3 of 11 techniques)
 - **Execution** (3 of 34)
 - **Persistence** (5 of 62)
 - **Privilege Escalation** (1 of 32)
 - **Defense Evasion** (5 of 69)
 - **Credential Access** (3 of 21)
 - **Discovery** (4 of 23)
 - **Lateral Movement** (4 of 18)
 - **Collection** (0 of 13)
 - **Command and Control** (20 of 22)
 - **Exfiltration** (3 of 9)
 - **Impact** (4 of 16)


Substantially increase coverage by adding host data (e.g., Sysflow, Event Log, ...)

Tensor Decomposition Coverage in ATT&CK

Covered: Data can be converted to tensors, decomposed, and anomalies identified

Initial Access 11 Items	Execution 34 Items	Persistence 69 Items	Privilege Escalation 32 Items	Defense Evasion 69 Items	Credential Access 21 Items	Discovery 23 Items	Lateral Movement 18 Items	Collection 13 Items	Command And Control 22 Items	Exfiltration 9 Items	Impact 16 Items
Security Compromise	AppScript	Web Jobs and Jobs	Access Token Manipulation	Account Manipulation	Account Manipulation	Account Discovery	AppScript	Screening Local Disk	Command and Control	Automated Exfiltration	Account Access Removal
Exploit Public-Facing Application	CMS3P	Accessibility Features	Accessibility Features	Binary Hijacking	Binary Hijacking	Application Window Discovery	Application Deployment Software	Automated Collection	Communication Through Removable Media	Data Compressed	Data Destruction
External Remote Services	Command-Line Interface	Account Manipulation	AppCert DLLs	BITS Jobs	Brute Force	Browser Bookmarks Discovery	Component Object Model and Related COM	Clipboard Data	Connection Proxy	Data Encrypted	Data Encrypted for Impact
Hardware Additions	Compiled HTML File	AppCert DLLs	Apprnt DLLs	Bypass User Account Control	Credential Dumping	Domain Trust Discovery	Exploitation of Remote Services	Data from Information Repositories	Custom Command and Control Protocol	Data Transfer Size Limits	Defacement
Replication Through Removable Media	Apprnt DLLs	Apprnt DLLs	Application Shimming	Clear Command History	Credentials from Web Browsers	File and Directory Discovery	Internal Spearphishing	Data from Local System	Custom Cryptographic Protocol	Exfiltration Over Alternative Protocol	Disk Content Wipe
Searchsploit	Control Panel Items	Application Shimming	Bypass User Account Control	CMS3P	Credentials in Files	Network Service Scoping	Login Scripts	Data from Network Attached Drive	Data Encoding	Exfiltration Over Command and Control	Disk Structure Wipe
Searchsploit via Service	Dynamic Data Exchange	Authentication Package	DLL Search Order Hijacking	Code Signing	Credentials in Registry	Network Share Discovery	Pass the Hash	Data from Removable Media	Data Obfuscation	Exfiltration Over Other Network Medium	Endpoint Denial of Service
Searchsploit via Service	Execution through API	BITS Jobs	Dylib Hijacking	Compte After Delivery	Exploitation for Credential Access	Network Sniffing	Pass the Ticket	Data Staged	Domain Fronting	Exfiltration Over Physical Medium	Firmware Corruption
Supply Chain Compromise	Execution through Module Load	Bootkit	Elevated Execution with Prompt	Compiled HTML File	Component Firmware	Password Policy Discovery	Remote Desktop Protocol	Email Collection	Domain Generation Algorithms	Schedules Transfer	Inhibit System Recovery
Trusted Relationship	Exploitation to Client Execution	Browser Extensions	Enroll	Component Firmware	Component Firmware	Peripheral Device Discovery	Remote File Copy	Input Capture	Farback Channels		Network Denial of Service
Valid Accounts	Graphical User Interface	Change Default File Association	Exploitation for Privilege Escalation	Component Object Model Hijacking	Input Capture	Process Discovery	Remote Services	Man in the Browser	Multi-hop Proxy		Resource Hijacking
	InstallINI	Component Firmware	Extra Window Memory Injection	Connection Proxy	Input Prompt	Query Registry	Replication Through Removable Media	Screen Capture	Multi-Stage Channels		Runtime Data Manipulation
	Launchctl	Component Object Model Hijacking	File System Permissions Weakness	Control Panel Items	Kerberoasting	Quota Registry	Shared Webroot	Video Capture	Multistep Communication		Service Stop
	Local Job Scheduling	Crista Account	Hooking	DCShadow	Keychain	Remote System Discovery	SSH Hijacking		Multistep Encryption		Steered Data Manipulation
	LSASS Driver	DLL Search Order Hijacking	Image File Execution Options Injection	Desktops/Codecs Files or Information	LSASSNTLM Poisoning and Relay	Security Software Discovery	Tare Shared Content		Port Knocking		System Shutdown/Reboot
	Mhta	Dylib Hijacking	Disabling Security Tools	Network Sniffing	Network Sniffing	Software Discovery	Third-party Software		Remote Access Tools		Transmitted Data Manipulation
	PowerShell	Enroll	New Service	DLL Search Order Hijacking	Password Filter DLL	System Information Discovery	Windows Admin Shares		Remote File Copy		
	Regsvr32	External Remote Services	Parent PID Spoofing	D.L. Side-Loading	Private Keys	System Network Configuration Discovery	Windows Remote Management		Standard Application Layer Protocol		
	Regsvr32	File System Permissions Weakness	Path Interception	Evolution Guardians	Secured Memory	System Network Connections Discovery			Standard Cryptographic Protocol		
	Runas	Hidden Files and Directories	Plist Modification	Exploitation for Defense Evasion	Steal Web Session Cookie	System Owner/User Discovery			Standard Non-Application Layer Protocol		
	Scheduled Task	Hooking	Port Monitors	Extra Window Memory Injection	Two-Factor Authentication Interception	System Time Discovery			Uncommonly Used Port		
	Scripting	Hypervisor	PowerShell Profile	File Deletion	File Deletion	System Time Discovery			Web Service		
	Service Execution	Image File Execution Options Injection	Process Injection	Scheduled Task	File System Logical Offsets	Virtualization/Sandbox Evasion					
	Signed Binary Proxy Execution	Kernel Modules and Extensions	Service Registry Permissions Weakness	Gatekeeper Bypass	Group Policy Modification						
	Signed Script Proxy Execution	Launch Agent	Launch Daemon	Behold and Setgid	Hidden Files and Directories						
	Source	Launchctl	LSASS Driver	LSASSNTLM Poisoning and Relay	Hidden Users						
	Space after Filename	LSASS Driver	LSASSNTLM Poisoning and Relay	Hidden Users	Hidden Window						
	Third-party Software	Local Job Scheduling	Local Job Scheduling	Local Job Scheduling	HISTCONTROL						
	Trap	Local Job Scheduling	Local Job Scheduling	Local Job Scheduling	Image File Execution Options Injection						
	Trusted Developer Utilities	Login Rem	Login Rem	Login Rem	Indicator Blocking						
	User Execution	Login Scripts	Login Scripts	Login Scripts	Indicator Removal from Tools						
	Windows Management Instrumentation	LSASS Driver	LSASS Driver	LSASS Driver	Indicator Removal on Host						
	Windows Remote Management	Modify Existing Service	Modify Existing Service	Modify Existing Service	Indirect Command Execution						
	XSL Script Processing	Netsh Helper DLL	Netsh Helper DLL	Netsh Helper DLL	Install Root Certificate						
		New Service	New Service	New Service	Install Root Certificate						
		Other Application Startup	Other Application Startup	Other Application Startup	Launchctl						
		Path Interception	Path Interception	Path Interception	LSASSNTLM Poisoning and Relay						
		Plist Modification	Plist Modification	Plist Modification	LSASSNTLM Poisoning and Relay						
		Port Knocking	Port Knocking	Port Knocking	LSASSNTLM Poisoning and Relay						
		Port Monitors	Port Monitors	Port Monitors	LSASSNTLM Poisoning and Relay						
		PowerShell Profile	PowerShell Profile	PowerShell Profile	LSASSNTLM Poisoning and Relay						
		Rc common	Rc common	Rc common	LSASSNTLM Poisoning and Relay						
		Re-opened Applications	Re-opened Applications	Re-opened Applications	LSASSNTLM Poisoning and Relay						
		Relevant Access	Relevant Access	Relevant Access	LSASSNTLM Poisoning and Relay						
		Registry Run Keys / Startup Folder	Registry Run Keys / Startup Folder	Registry Run Keys / Startup Folder	LSASSNTLM Poisoning and Relay						
		Scheduled Task	Scheduled Task	Scheduled Task	LSASSNTLM Poisoning and Relay						
		Screen saver	Screen saver	Screen saver	LSASSNTLM Poisoning and Relay						
		Security Support Provider	Security Support Provider	Security Support Provider	LSASSNTLM Poisoning and Relay						
		Server Software Component	Server Software Component	Server Software Component	LSASSNTLM Poisoning and Relay						
		Service Registry Permissions Weakness	Service Registry Permissions Weakness	Service Registry Permissions Weakness	LSASSNTLM Poisoning and Relay						
		Send and Setgid	Send and Setgid	Send and Setgid	LSASSNTLM Poisoning and Relay						
		Shortcut Modification	Shortcut Modification	Shortcut Modification	LSASSNTLM Poisoning and Relay						
		SIP and Trust Provider Hijacking	SIP and Trust Provider Hijacking	SIP and Trust Provider Hijacking	LSASSNTLM Poisoning and Relay						
		Startup Items	Startup Items	Startup Items	LSASSNTLM Poisoning and Relay						
		System Firmware	System Firmware	System Firmware	LSASSNTLM Poisoning and Relay						
		Systemd	Systemd	Systemd	LSASSNTLM Poisoning and Relay						
		Time Providers	Time Providers	Time Providers	LSASSNTLM Poisoning and Relay						
		Trap	Trap	Trap	LSASSNTLM Poisoning and Relay						
		Valid Accounts	Valid Accounts	Valid Accounts	LSASSNTLM Poisoning and Relay						
		Web Shell	Web Shell	Web Shell	LSASSNTLM Poisoning and Relay						
		Windows Management Instrumentation Event Subscription	Windows Management Instrumentation Event Subscription	Windows Management Instrumentation Event Subscription	LSASSNTLM Poisoning and Relay						
		Winlogon Helper DLL	Winlogon Helper DLL	Winlogon Helper DLL	LSASSNTLM Poisoning and Relay						

 Covered by Zeek log tensor decompositions

 Covered by host data tensor decompositions

Example Detection of ATT&CK Technique

Tactic and Technique

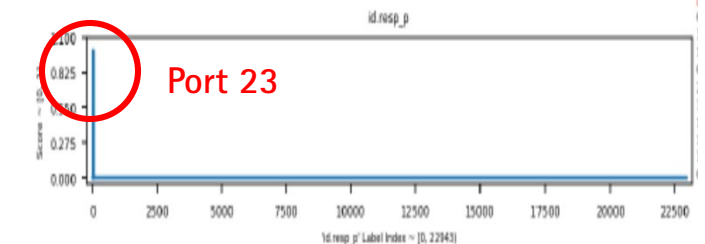
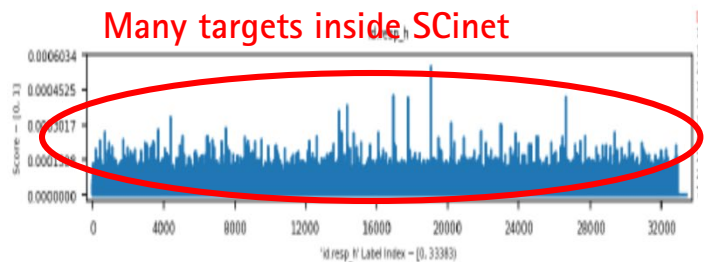
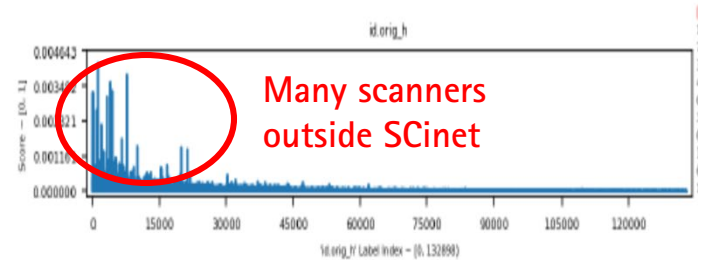
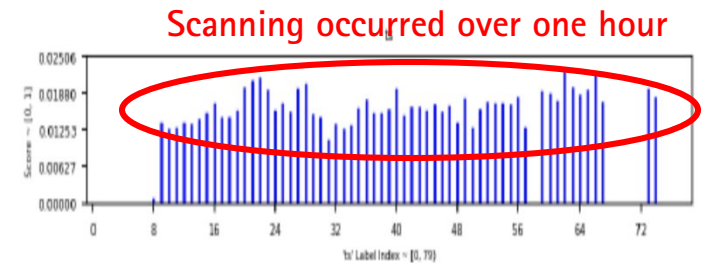
- Discovery – Network Service Scanning

Context

- SCinet 2019
- Network for Supercomputing conference
- All IP addresses public (no firewalls)
- No authentication / authorization
- ~8 Million flows per hour

Details

- Large number of external hosts scanning SCinet
- ~176K flows on port 23
- Potential coordination
- Scan evaded other scan detection tools



PART 2: ANOMALY DETECTION

Need to Automate Anomaly Detection



Often **100+** components needed to characterize network traffic

Most components are **benign**

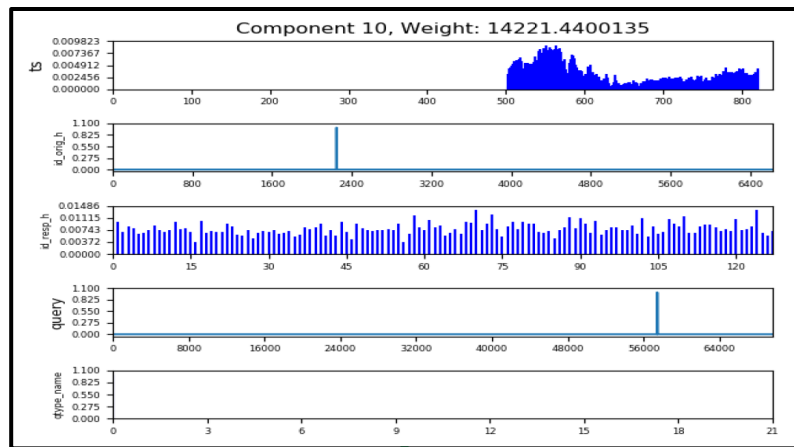
Challenge is to **identify and rank** components representing anomalous behavior

Components are **trailheads** for further investigation

Each component can take minutes or hours to manually investigate

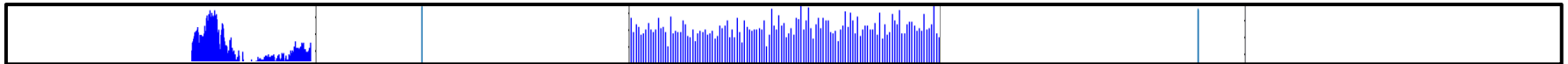
Which components are interesting?

Topic Modeling for Component Classification



Latent Dirichlet Allocation (LDA)

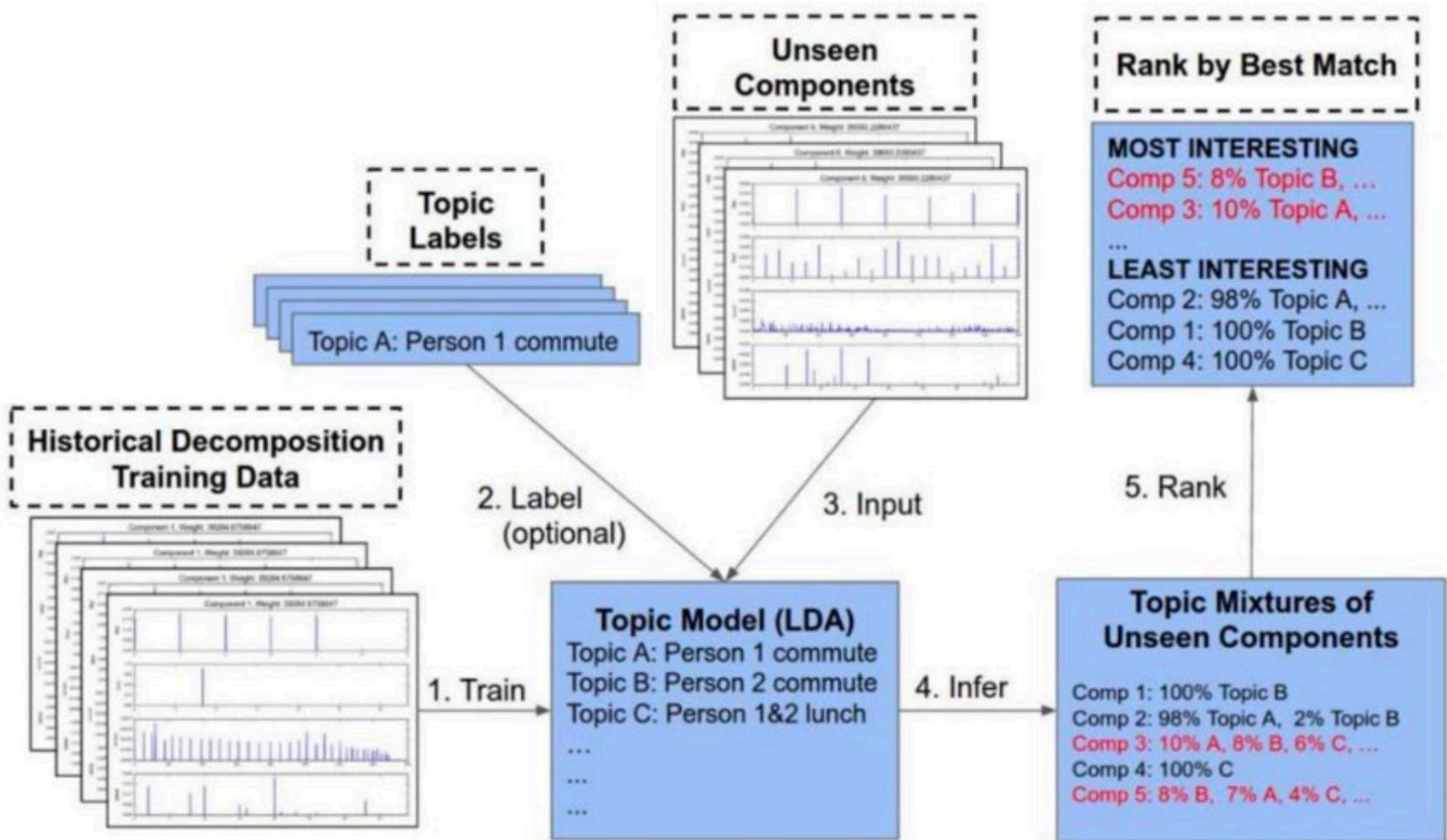
- Well-known Bayesian topic modeling algorithm
- Learns topic model from a corpus of documents
- Infers topic mixture of new documents
- Online updates of topic model
- Commonly used in other applications
 - Bioinformatics
 - Image, video, and sound processing
 - Collaborative filtering



Mapping tensor decompositions to LDA concepts

- Component (as vector) = "document"
- Label = "word"
- Score = "word count"
- Topic = recognizable pattern of network behavior

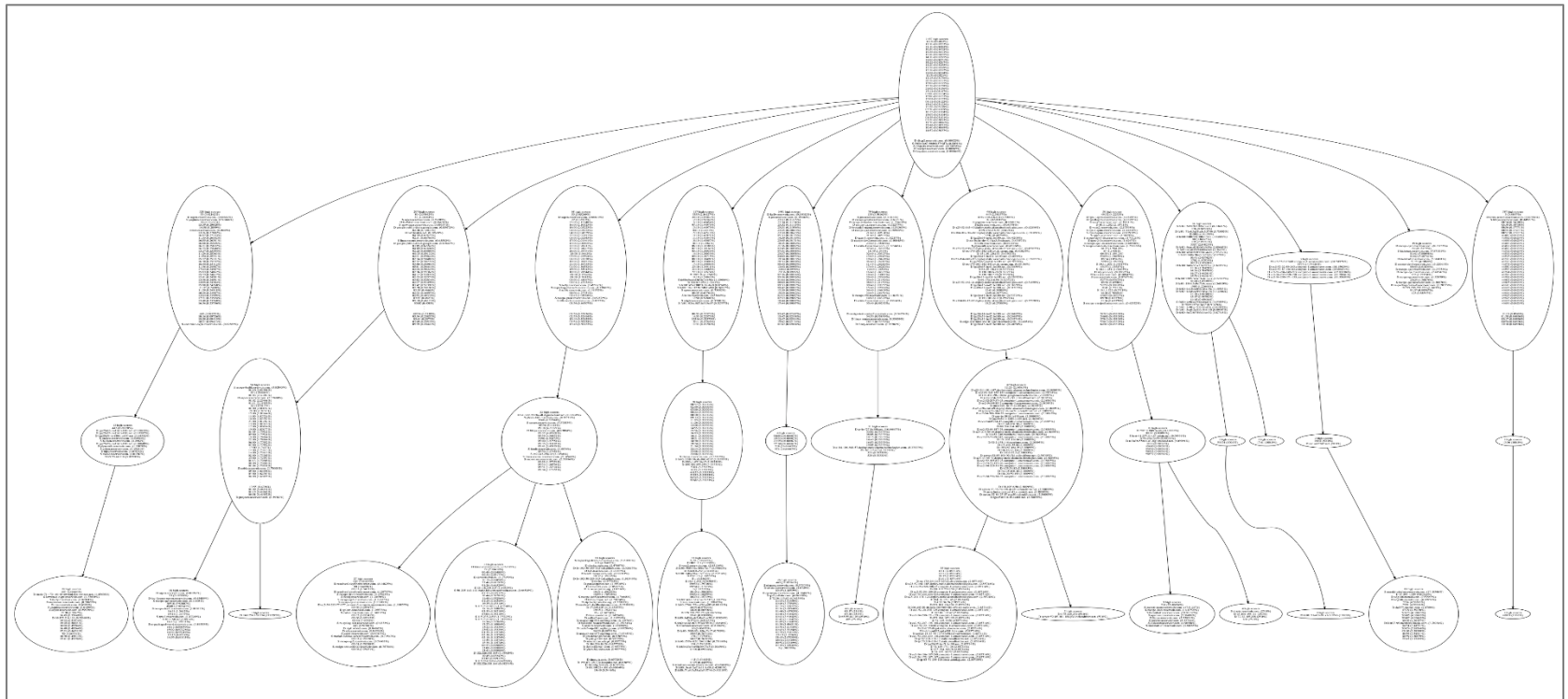
LDA Dominant Topic Approach



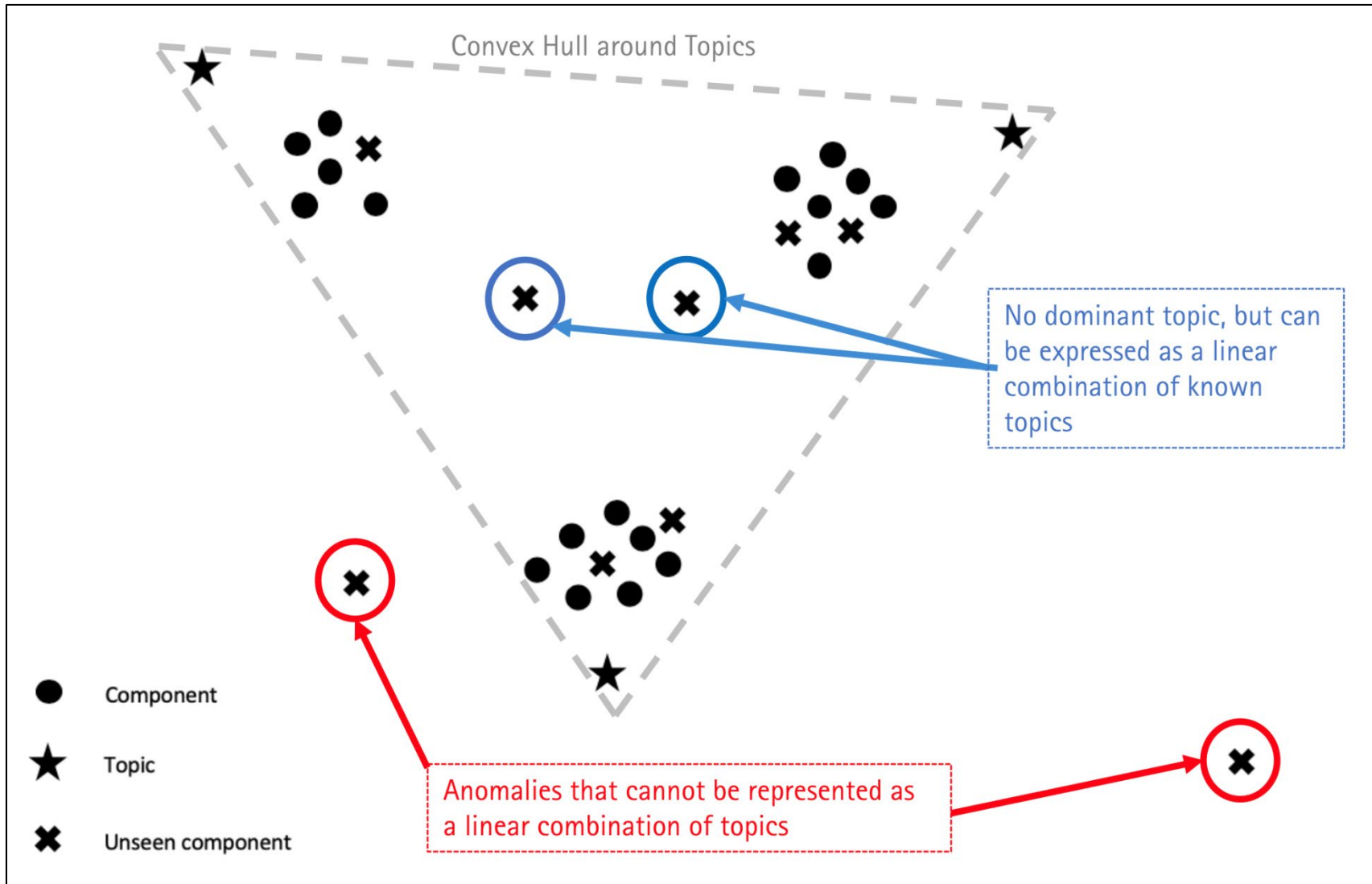
Hierarchical LDA Approach

Learn topics in tree

- Coarse grain behavior at root, fine grain at leaves
- Topic is weighted mixture of root-to-leaf paths in tree
- Same approach as dominant topic otherwise



Limitations of Dominant Topic Approaches

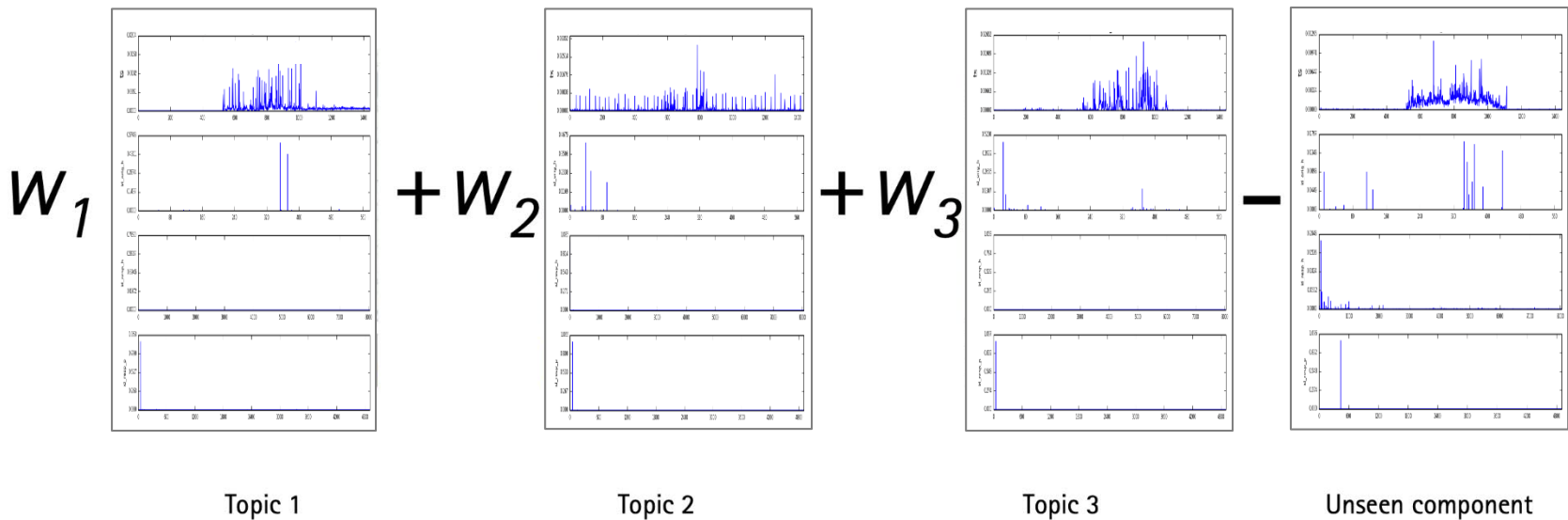


Component Reconstruction Approach

Addresses mathematical limitations of dominant topic approach

Infer topic mixtures for unseen components and reconstruct with known topics

Compare to unseen component and rank by reconstruction error



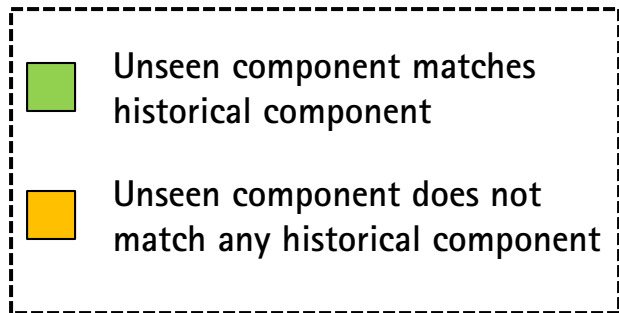
2

Decomposition Difference Approach

Compute similarity matrix between current and historical decomposition components

Component(s) dissimilar to every historical component represents anomalous behavior

Rank by max similarity



	.00	.01	.04	.01	.99
Unseen Components	.95	.02	.01	.00	.02
	.00	.01	.00	.00	.03
	.02	.98	.05	.03	.01
	.00	.02	.01	.97	.01

Historical Components

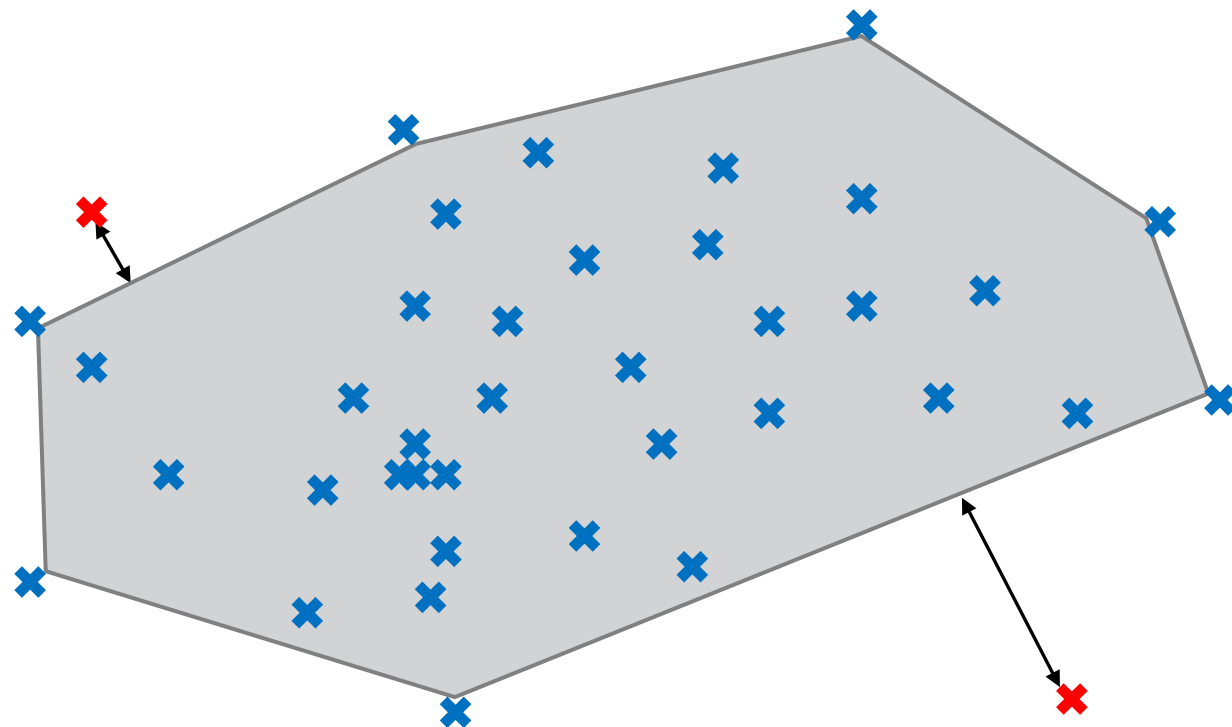
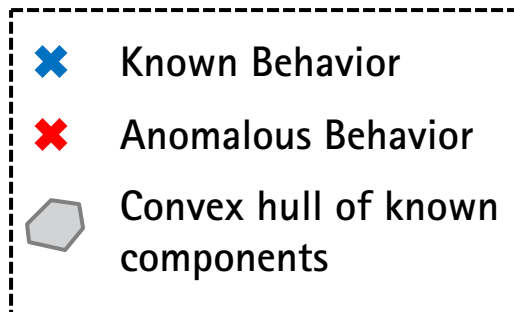
Approximate Convex Hull Approach

Compute approximate convex hull of historical decomposition components

If a component is a linear combination of historical components, it's inside the hull and we've seen all aspects of the behavior it represents

Identify anomalous components outside hull, compute distance to hull

Rank by distance to hull

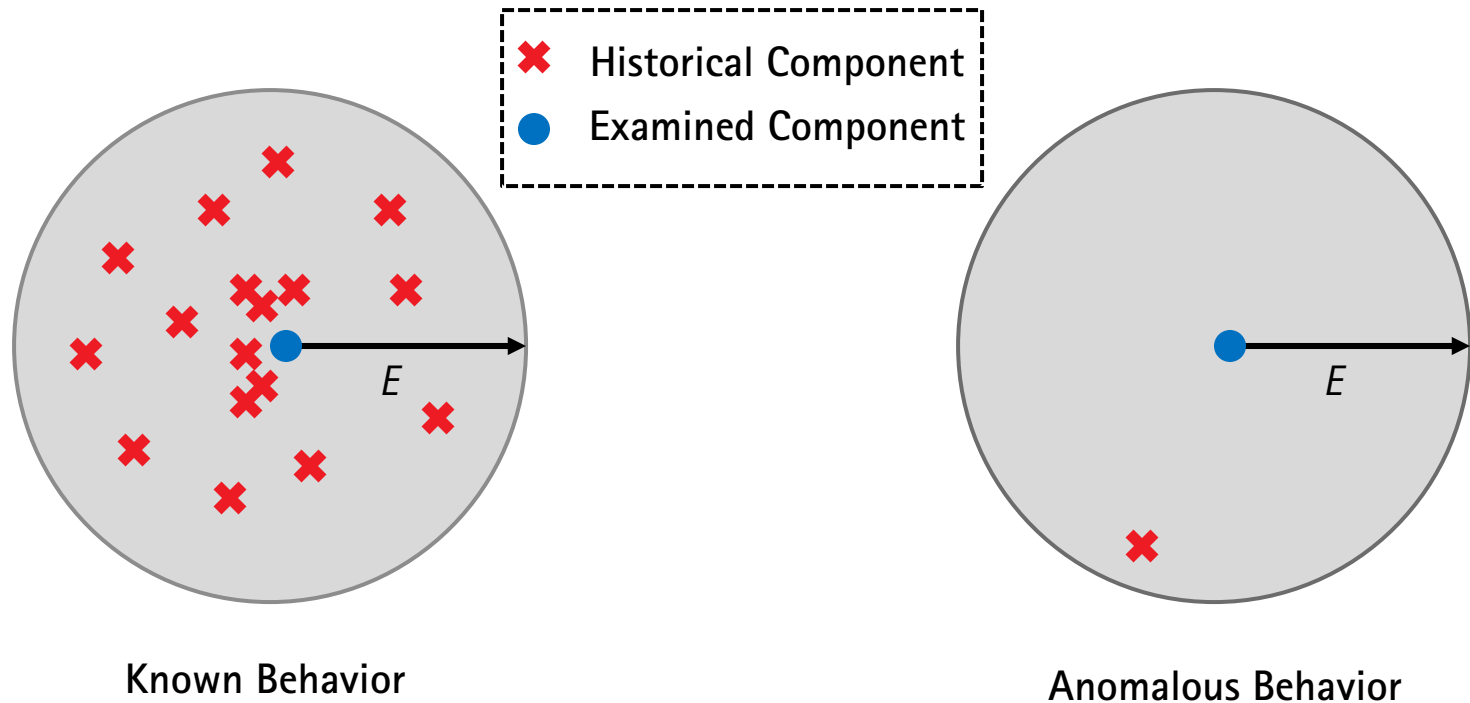


Epsilon Ball Approach

Treat component as vector, compare to historical components

Count components inside a hypersphere of radius E

Rank by count of components inside hypersphere



Comparison of Anomaly Detection Approaches

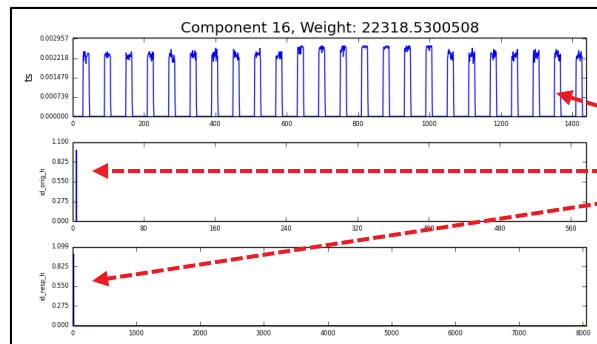
	Execution Time	Parametric	Detects Anomalous Variations of Historical Behavior	Detects Anomalous Behavior Unrelated to Historical Behavior
LDA – Dom Topic	High	Yes	Yes	No
HLDA – Dom Topic	High	No	Yes	No
LDA – Component Reconstruct	High	Yes	Yes	Yes
HLDA – Component Reconstruct	High	No	Yes	Yes
Decomp Diff	Low	Yes	Somewhat	Yes
Approximate Convex Hull	Low	No	No	Yes
Epsilon Ball	Low	Yes	Somewhat	Yes

PART 3: GRAPHS AND DATABASES

Graphs and Databases in Context

Components only tell a small part of the story

- E.g., Timestamp, Source IP, Destination IP



Component represents beaconing behavior between two IP addresses. Is it C2 traffic? Hourly batch jobs? Hourly log transfers?

More information necessary to make a malicious / benign decision

- E.g., user, asset type, network topology, known behaviors, threat intel, ...
- Needed info stored in external DB / graph / ... or enriched data in SIEM

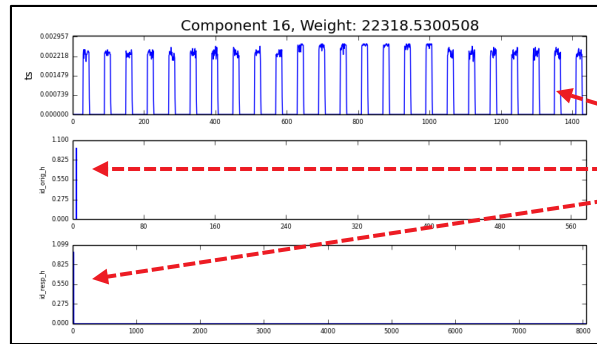
Use anomalous component as trailhead into investigation

- Generate targeted queries to provide context and assist decision making
- Massively reduces scope of graph / database analysis

Generating Targeted Queries

Use component labels with nonzero scores to generate "WHERE" clause

- E.g., "SELECT * WHERE ts=(00:00, 01:00, ...), src_ip=1.2.3.4, dst_ip=5.6.7.8"



Component represents beaconing behavior between two IP addresses. Is it C2 traffic? Hourly batch jobs? Hourly log transfers?

Problem: Data was binned before conversion to tensor

Solution Part 1: Generate backtracking data when building tensor

- Map tensor entries to lines in original log

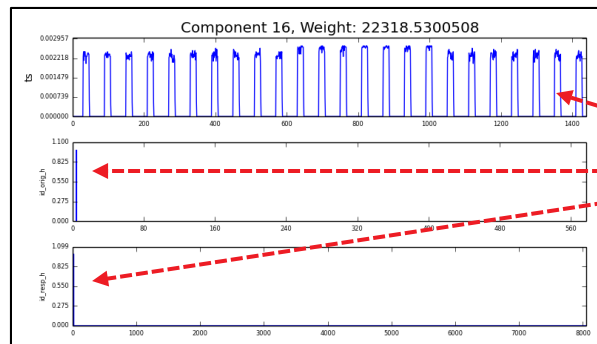
Solution Part 2: Reconstruct into tensor, get subset of relevant log entries

- Original entries provide more context – exact timestamps, flow IDs, ...

Generating Targeted Queries

Use enriched data to filter false positives

- E.g., "SELECT * WHERE ts=(00:00, 01:00, ...), src_ip=1.2.3.4, dst_ip=5.6.7.8"
AND src_ip NOT "batch_server" AND src_ip NOT "log_transfer_hourly"



Component represents beaconing behavior between two IP addresses. Is it C2 traffic? Hourly batch jobs? Hourly log transfers?

Further queries based on results of targeted query

- Query within the returned data or use as guide for further focused queries

Targeted query massively reduces size of graph / DB / SIEM data to investigate

- Not "boiling the ocean" by running analytics over entire graph / DB / SIEM
- Tensor decompositions highly optimized and run on ten-billion scale logs in reasonable time (high minutes / low hours)

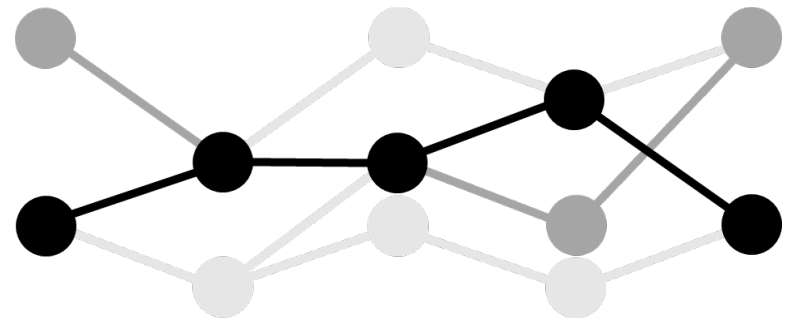
Conclusion

Contact the Speaker

- Thomas Henretty
henretty@reservoir.com

Recent Papers

- *Combining Tensor Decompositions and Graph Analytics to Provide Cyber Situational Awareness at HPC Scale*
HPEC, Sep 2019
- *Fast and Scalable Distributed Tensor Decompositions*
HPEC, Sep 2019
- *Enhancing Network Visibility and Security through Tensor Analysis*
Future Generation Computer Systems, July 2019



Pattern Discovery

Tensor decomposition provides a model for Zeek log data that allows behaviors to be separated as coherent patterns