

CE 815 – Secure Software Systems

Causal Analysis (Holmes)

Mehdi Kharrazi

Department of Computer Engineering
Sharif University of Technology



Acknowledgments: Some of the slides are fully or partially obtained from other sources. A reference is noted on the bottom of each slide, when the content is fully obtained from another source. Otherwise a full list of references is provided on the last slide. Thanks to Zahra Fazli for the help on the slides.



- globally-accessible knowledge base of adversary tactics and techniques based on real-world observations

The screenshot shows the MITRE ATT&CK website. The navigation bar includes links for Matrices, Tactics, Techniques, Defenses, CTI, Resources, Benefactors, and Blog, along with a search box. A banner below the navigation bar states: "ATT&CK v16 has been released! Check out the [blog post](#) for more information." The main content area features the "ATT&CK" logo and a navigation menu with links for "Get Started", "Take a Tour", "Contribute", "Blog", "FAQ", and "Random Page". To the right of the logo, there is introductory text: "MITRE ATT&CK® is a globally-accessible knowledge base of adversary tactics and techniques based on real-world observations. The ATT&CK knowledge base is used as a foundation for the development of specific threat models and methodologies in the private sector, in government, and in the cybersecurity product and service community. With the creation of ATT&CK, MITRE is fulfilling its mission to solve problems for a safer world – by bringing communities together to develop more effective cybersecurity. ATT&CK is open and available to any person or organization for use at no charge." Below this text is the "ATT&CK Matrix for Enterprise" section, which includes a "layout: side" dropdown and buttons for "show sub-techniques" and "hide sub-techniques". The matrix itself is a table with columns for different categories and their respective technique counts.

Reconnaissance	Resource Development	Initial Access	Execution	Persistence	Privilege Escalation	Defense Evasion	Credential Access	Discovery	Lateral Movement	Col...
10 techniques	8 techniques	10 techniques	14 techniques	20 techniques	14 techniques	44 techniques	17 techniques	32 techniques	9 techniques	17 te...
Active Scanning (3)	Acquire Access	Content Injection	Cloud Administration	Account Manipulation (27)	Abuse Elevation	Abuse Elevation Control Mechanism (46)	Adversary-in-the-Middle (46)	Account Discovery (4)	Exploitation of Remote	Advers...

Enterprise Tactics



ID	Name	Description
TA0043	Reconnaissance	The adversary is trying to gather information they can use to plan future operations.
TA0042	Resource Development	The adversary is trying to establish resources they can use to support operations.
TA0001	Initial Access	The adversary is trying to get into your network.
TA0002	Execution	The adversary is trying to run malicious code.
TA0003	Persistence	The adversary is trying to maintain their foothold.
TA0004	Privilege Escalation	The adversary is trying to gain higher-level permissions.
TA0005	Defense Evasion	The adversary is trying to avoid being detected.
TA0006	Credential Access	The adversary is trying to steal account names and passwords.
TA0007	Discovery	The adversary is trying to figure out your environment.
TA0008	Lateral Movement	The adversary is trying to move through your environment.
TA0009	Collection	The adversary is trying to gather data of interest to their goal.
TA0011	Command and Control	The adversary is trying to communicate with compromised systems to control them.
TA0010	Exfiltration	The adversary is trying to steal data.
TA0040	Impact	The adversary is trying to manipulate, interrupt, or destroy your systems and data.



Enterprise Techniques: Initial Access

ID	Name	Description
T1659	Content Injection	Adversaries may gain access and continuously communicate with victims by injecting malicious content into systems through online network traffic. Rather than luring victims to malicious payloads hosted on a compromised website (i.e., <i>Drive-by Target</i> followed by <i>Drive-by Compromise</i>), adversaries may initially access victims through compromised data-transfer channels where they can manipulate traffic and/or inject their own content. These compromised online network channels may also be used to deliver additional payloads (i.e., <i>Ingress Tool Transfer</i>) and other data to already compromised systems.
T1189	Drive-by Compromise	Adversaries may gain access to a system through a user visiting a website over the normal course of browsing. With this technique, the user's web browser is typically targeted for exploitation, but adversaries may also use compromised websites for non-exploitation behavior such as acquiring <i>Application Access Token</i> .
T1190	Exploit Public-Facing Application	Adversaries may attempt to exploit a weakness in an Internet-facing host or system to initially access a network. The weakness in the system can be a software bug, a temporary glitch, or a misconfiguration.
T1133	External Remote Services	Adversaries may leverage external-facing remote services to initially access and/or persist within a network. Remote services such as VPNs, Citrix, and other access mechanisms allow users to connect to internal enterprise network resources from external locations. There are often remote service gateways that manage connections and credential authentication for these services. Services such as <i>Windows Remote Management</i> and <i>VNC</i> can also be used externally.
T1200	Hardware Additions	Adversaries may introduce computer accessories, networking hardware, or other computing devices into a system or network that can be used as a vector to gain access. Rather than just connecting and distributing payloads via removable storage (i.e. <i>Replication Through Removable Media</i>), more robust hardware additions can be used to introduce new functionalities and/or features into a system that can then be abused.
T1566	Phishing	Adversaries may send phishing messages to gain access to victim systems. All forms of phishing are electronically delivered social engineering. Phishing can be targeted, known as spearphishing. In spearphishing, a specific individual, company, or industry will be targeted by the adversary. More generally, adversaries can conduct non-targeted phishing, such as in mass malware spam campaigns.
.001	Spearphishing Attachment	Adversaries may send spearphishing emails with a malicious attachment in an attempt to gain access to victim systems. Spearphishing attachment is a specific variant of spearphishing. Spearphishing attachment is different from other forms of spearphishing in that it employs the use of malware attached to an email. All forms of spearphishing are electronically delivered social engineering targeted at a specific individual, company, or industry. In this scenario, adversaries attach a file to the spearphishing email and usually rely upon <i>User Execution</i> to gain execution. Spearphishing may also involve social engineering techniques, such as posing as a trusted source.
.002	Spearphishing Link	Adversaries may send spearphishing emails with a malicious link in an attempt to gain access to victim systems. Spearphishing with a link is a specific variant of spearphishing. It is different from other forms of spearphishing in that it employs the use of links to download malware contained in email, instead of attaching malicious files to the email itself, to avoid defenses that may inspect email attachments. Spearphishing may also involve social engineering techniques, such as posing as a trusted source.
.003	Spearphishing via Service	Adversaries may send spearphishing messages via third-party services in an attempt to gain access to victim systems. Spearphishing via service is a specific variant of spearphishing. It is different from other forms of spearphishing in that it employs the use of third party services rather than directly via enterprise email channels.



Techniques :Content Injection

Content Injection

Adversaries may gain access and continuously communicate with victims by injecting malicious content into systems through online network traffic. Rather than luring victims to malicious payloads hosted on a compromised website (i.e., [Drive-by Target](#) followed by [Drive-by Compromise](#)), adversaries may initially access victims through compromised data-transfer channels where they can manipulate traffic and/or inject their own content. These compromised online network channels may also be used to deliver additional payloads (i.e., [Ingress Tool Transfer](#)) and other data to already compromised systems.^[1]

Mitigations

ID	Mitigation	Description
M1041	Encrypt Sensitive Information	Where possible, ensure that online traffic is appropriately encrypted through services such as trusted VPNs.
M1021	Restrict Web-Based Content	Consider blocking download/transfer and execution of potentially uncommon file types known to be used in adversary campaigns.

Detection

ID	Data Source	Data Component	Detects
DS0022	File	File Creation	Monitor for unexpected and abnormal file creations that may indicate malicious content injected through online network communications.
DS0029	Network Traffic	Network Traffic Content	Monitor for other unusual network traffic that may indicate additional malicious content transferred to the system. Use network intrusion detection systems, sometimes with SSL/TLS inspection, to look for known malicious payloads, content obfuscation, and exploit code.
DS0009	Process	Process Creation	Look for behaviors on the endpoint system that might indicate successful compromise, such as abnormal behaviors of browser processes. This could include suspicious files written to disk, evidence of Process Injection for attempts to hide execution, or evidence of Discovery .

HOLMES: Real-time APT Detection through Correlation of Suspicious Information Flows, Sadegh M. Milajerdi, Rigel Gjomemo, Birhanu Eshetey, R. Sekarz, V.N. Venkatakrishnan, IEEE Symposium on Security and Privacy, 2019.



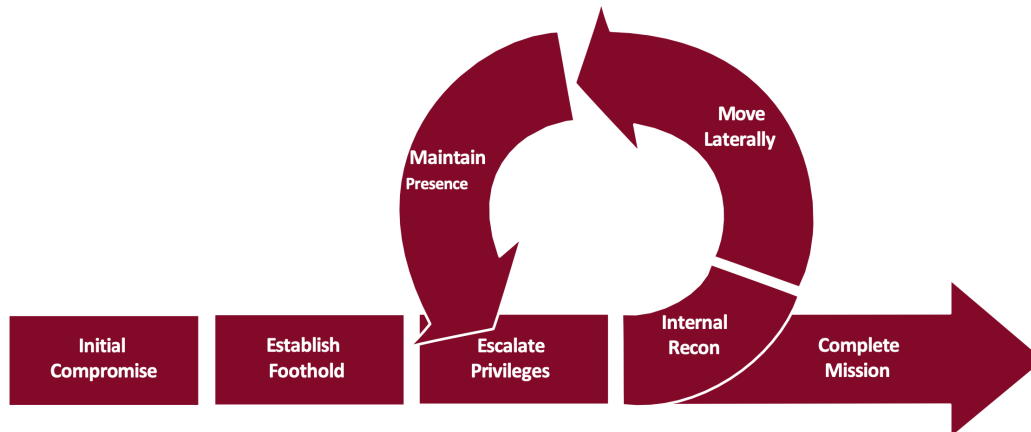
Advanced Persistent Threat (APT) and its challenges

- Targeted cyber attacks on organizations getting more sophisticated and stealthy.
- Goal: to steal data, disrupt operations or destroy infrastructure.
- APTs combine many different attack vectors
- Each appearing in some log sources
- Firewall, IDS/IPS, Netflow, DNS logs, Identity and access management tools
- Might occur over a long duration
- Correlating heterogeneous alarms using heuristics like timestamp is not so effective
- Lacking the full picture (root cause, affected entities, etc.).
- Significant manual effort and expertise are needed to piece together numerous alarms emanated by multiple security tools.



Intuition

- APT behaviors often conform to the kill-chain [MANDIANT-APT1]

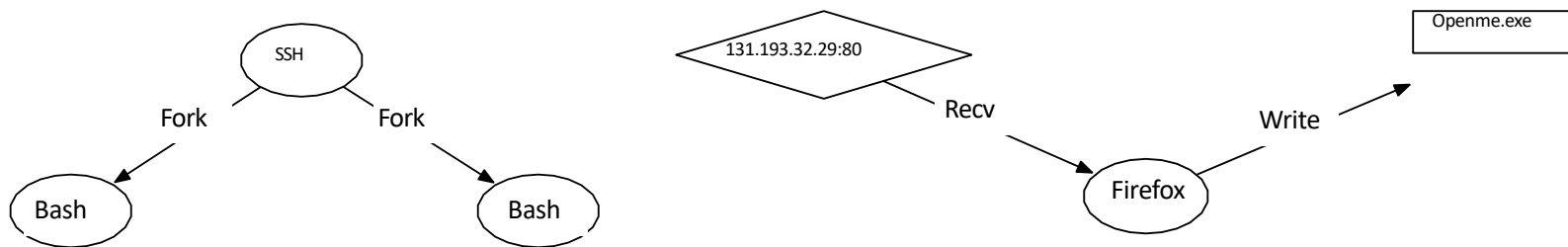


- Our analysis of over 300 APT whitepapers confirms that most APTs follow this kill-chain
- In particular, high-level steps of APTs need to be causally connected
- Use connectedness of high-level steps as a basis for campaign detection



Approach

- Use Provenance Graph to enable alert correlation for attack campaign detection
- vertices: system entities (socket, process, file, memory, etc.), and agents (user, groups,...)
- edges: system calls (causal dependencies or information flow)

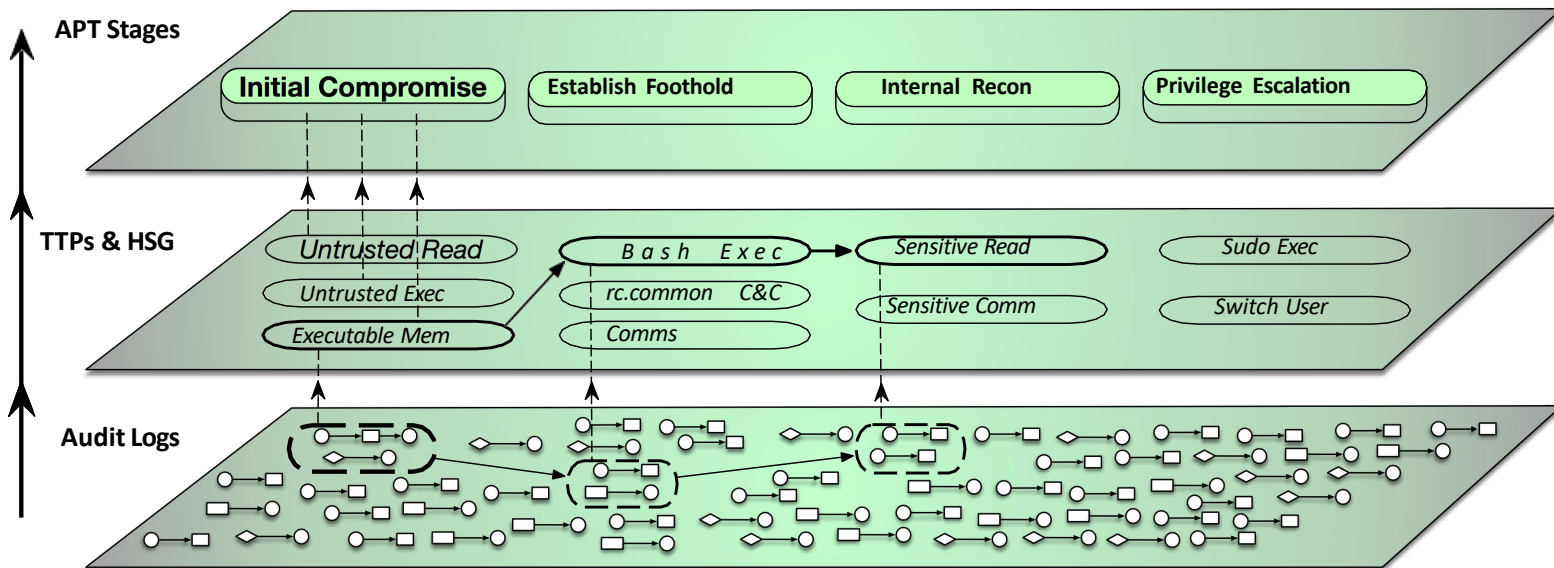


- Leverage the full historical context of a system
- Reason about interrelationships between different events and objects
- Key challenge: How to bridge semantic gap between low-level records and high-level activities in kill-chain?



Bridging the Semantic Gap

- Use Tactics, Techniques, and Procedures (TTPs) from MITRE's ATT&CK framework as an intermediate layer to bridge low-level audit records to high-level steps





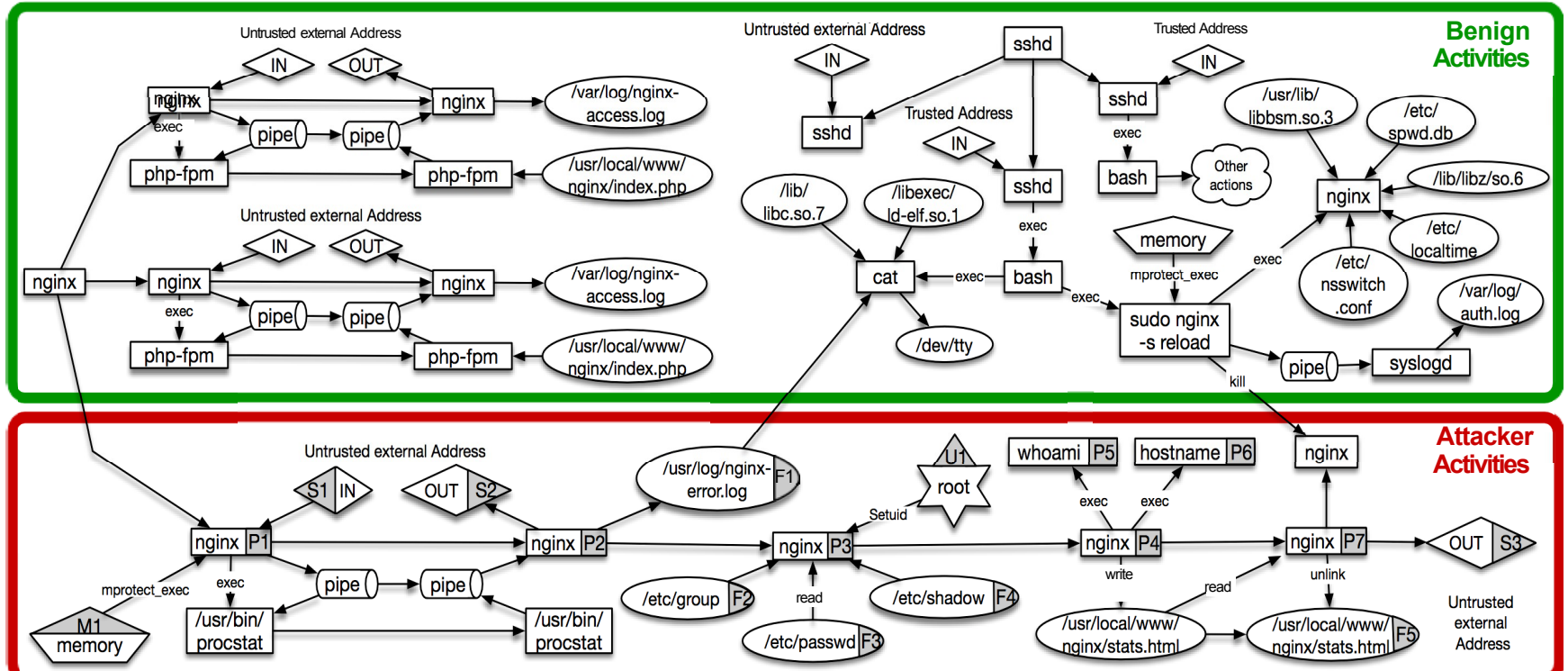
Bridging the Semantic Gap

Hardware Additions	Scheduled Task		Binary Padding	Credentials in Registry	Browser Bookmark Discovery	Exploitation of Remote Services	Data from Information Repositories	Exfiltration Over Physical Medium	Remote Access Tools
Trusted Relationship	LSASS Driver		Extra Window Memory Injection	Exploitation for Credential Access	Network Share Discovery	Distributed Component Object Model	Video Capture	Exfiltration Over Command and Control Channel	Port Knocking
Supply Chain Compromise	Local Job Scheduling		Access Token Manipulation	Forced Authentication	Peripheral Device Discovery	Pass the Ticket	Audio Capture	Automated Exfiltration	Multi-hop Proxy
Spearphishing Attachment	Trap		Bypass User Account Control	Hooking	File and Directory Discovery	Replication Through Removable Media	Clipboard Data	Exfiltration Over Other Network Medium	Domain Fronting
Exploit Public-Facing Application	Launch Daemon		Process Injection	Password Filter DLL	Permission Groups Discovery	Windows Admin Shares	Automated Collection	Standard	Data Encoding
Replication Through Removable Media	Scheduled Binary		Image File Execution Options Injection	LLMNR/NBT-NS Poisoning	System Network Connections Discovery	Third-party Software	Email Collection	Non-Application Layer Protocol	Remote File Copy
Spearphishing via Service	Proxy Execution		Valid Accounts	Private Keys	System Owner/User Discovery	Shared Web boot	Screen Capture	Multi-Stage Channels	Web Service
Drive-by Compromise	User Execution		Valid Accounts	Keychain	System Network Connections Discovery	Logon Scripts	Data Staged	Standard Application Layer Protocol	Remote File Copy
Valid Accounts	CMSTP		Hooking	Input Prompt	System Network Connections Discovery	Logon Scripts	Input Capture	Standard Application Layer Protocol	Remote File Copy
Valid Accounts	Dynamic Data Exchange		Startup Items	Bash History	System Network Connections Discovery	Logon Scripts	Data from Network Shared Drive	Exfiltration Over Alternative Protocol	Remote File Copy
Valid Accounts	Mshta		Launch Daemon	Port Knocking	System Network Connections Discovery	Logon Scripts	Data from Local System	Data Transfer	Remote File Copy
Valid Accounts	AppleScript		Dylib Hijacking	Indirect Command Execution	System Network Connections Discovery	Logon Scripts	Man in the Browser	Size Limits	Remote File Copy
Valid Accounts	Source		Application Shimming	Indirect Command Execution	System Network Connections Discovery	Logon Scripts	Data from Removable Media	Scheduled Transfer	Remote File Copy
Valid Accounts	Space after Filename		AppInit DLLs	BITS Jobs	System Network Configuration Discovery	Application Deployment Software			Remote File Copy
Valid Accounts	Execution through Module Load		Web Shell	Control Panel Items	Application Window Discovery	SSH Hijacking			Remote File Copy
Valid Accounts	Regsvcs/Regasm		Service Registry Permissions Weakness	CMSTP	Application Window Discovery	AppleScript			Remote File Copy
Valid Accounts	InstallUtil		New Service	Process Doppelganging	Application Window Discovery	AppleScript			Remote File Copy
Valid Accounts	File System Permissions Weakness		Mshta	Credential Dumping	Password Policy Discovery	Target Shared Content			Remote File Copy
Valid Accounts	Regsvr32		Hidden Files and Directories	Securityd Memory	System Time Discovery	Remote Desktop Protocol			Remote File Copy
Valid Accounts	PowerShell		Space after Filename	Brute Force	Account Discovery	Remote Services			Remote File Copy

ATT&CK™

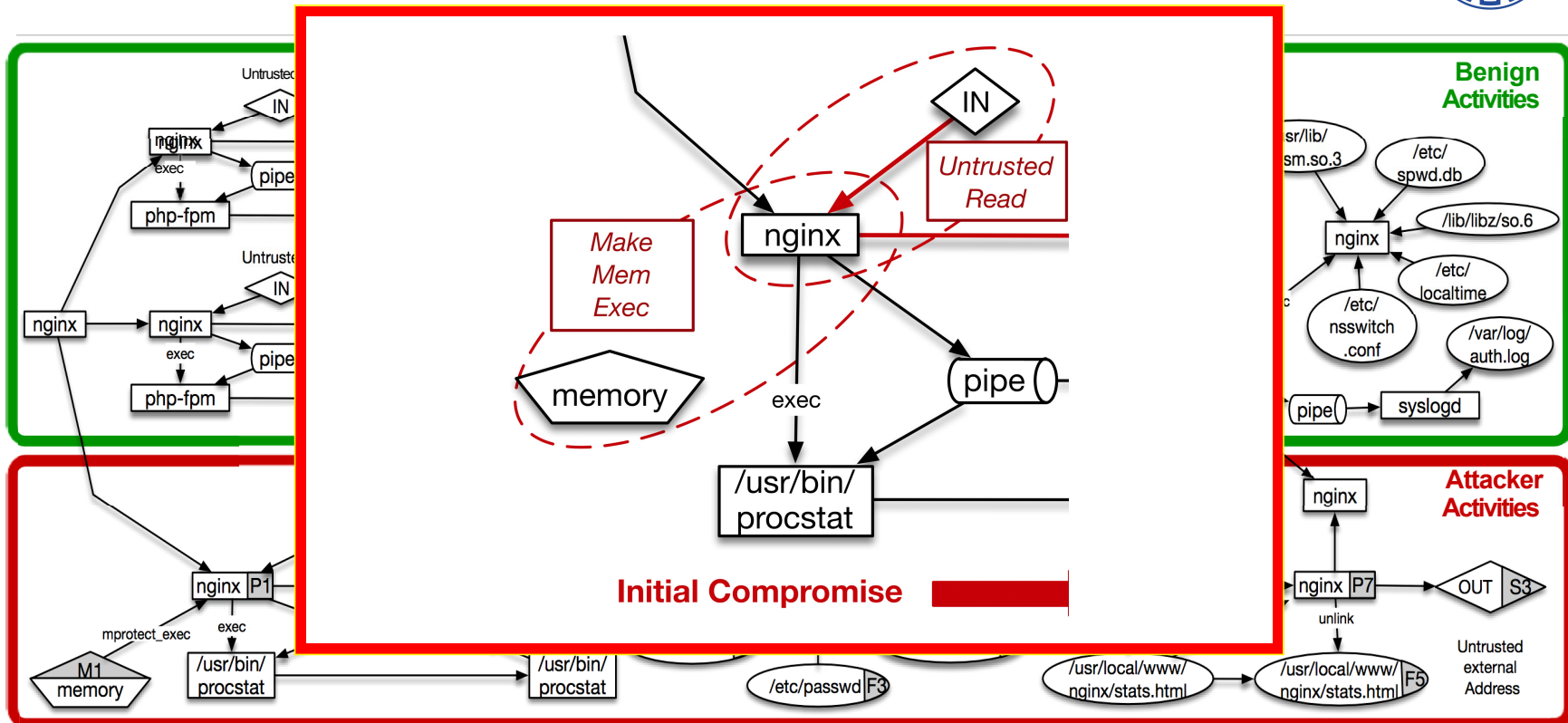


Illustrative Example



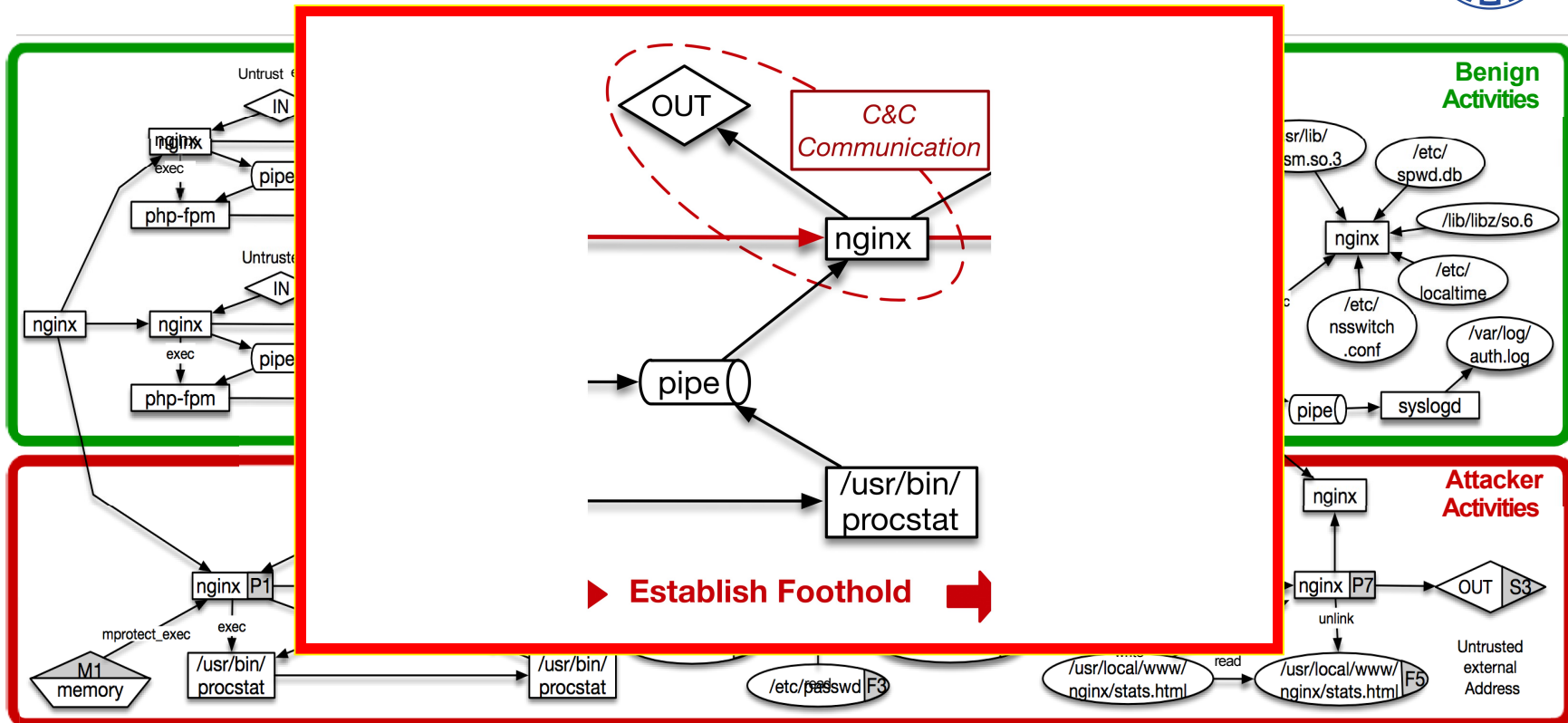


Illustrative Example



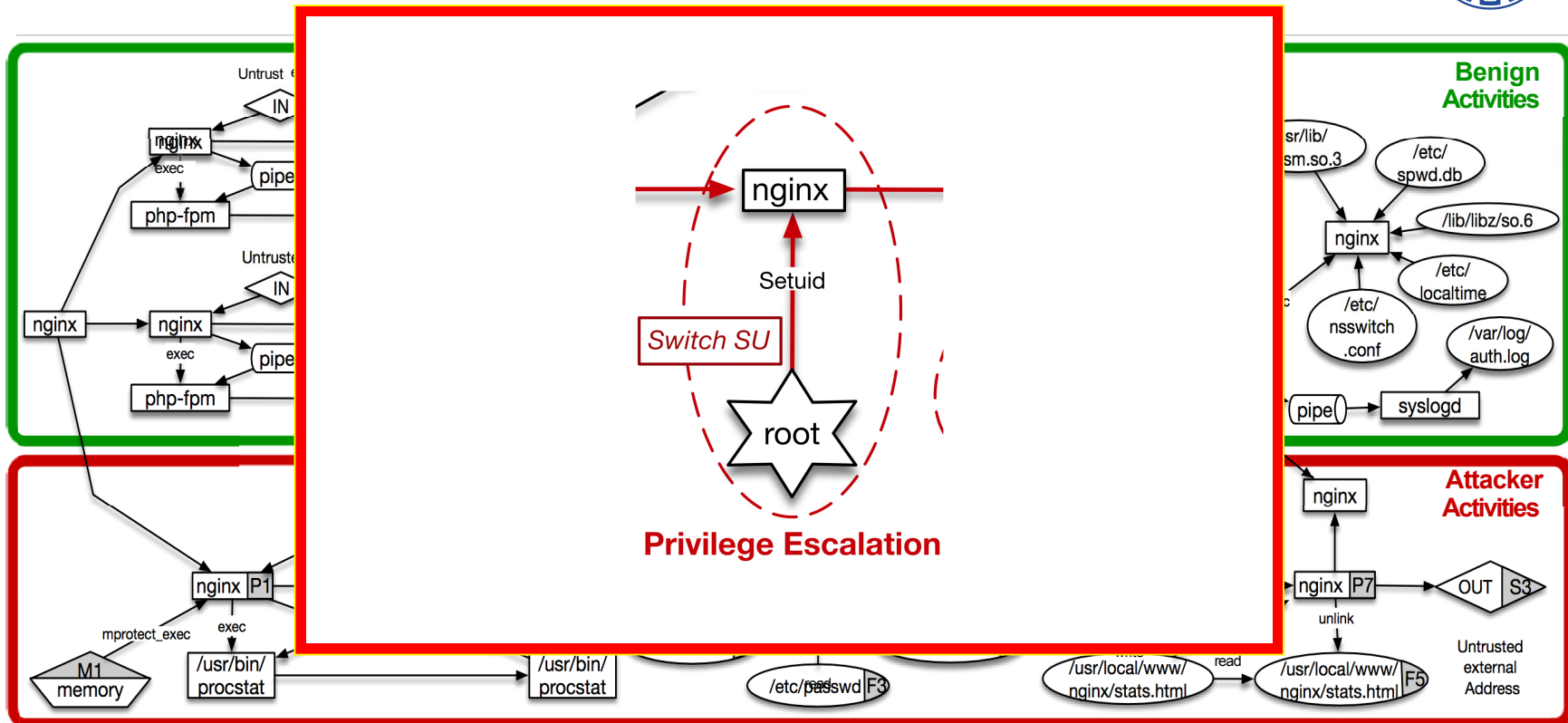


Illustrative Example



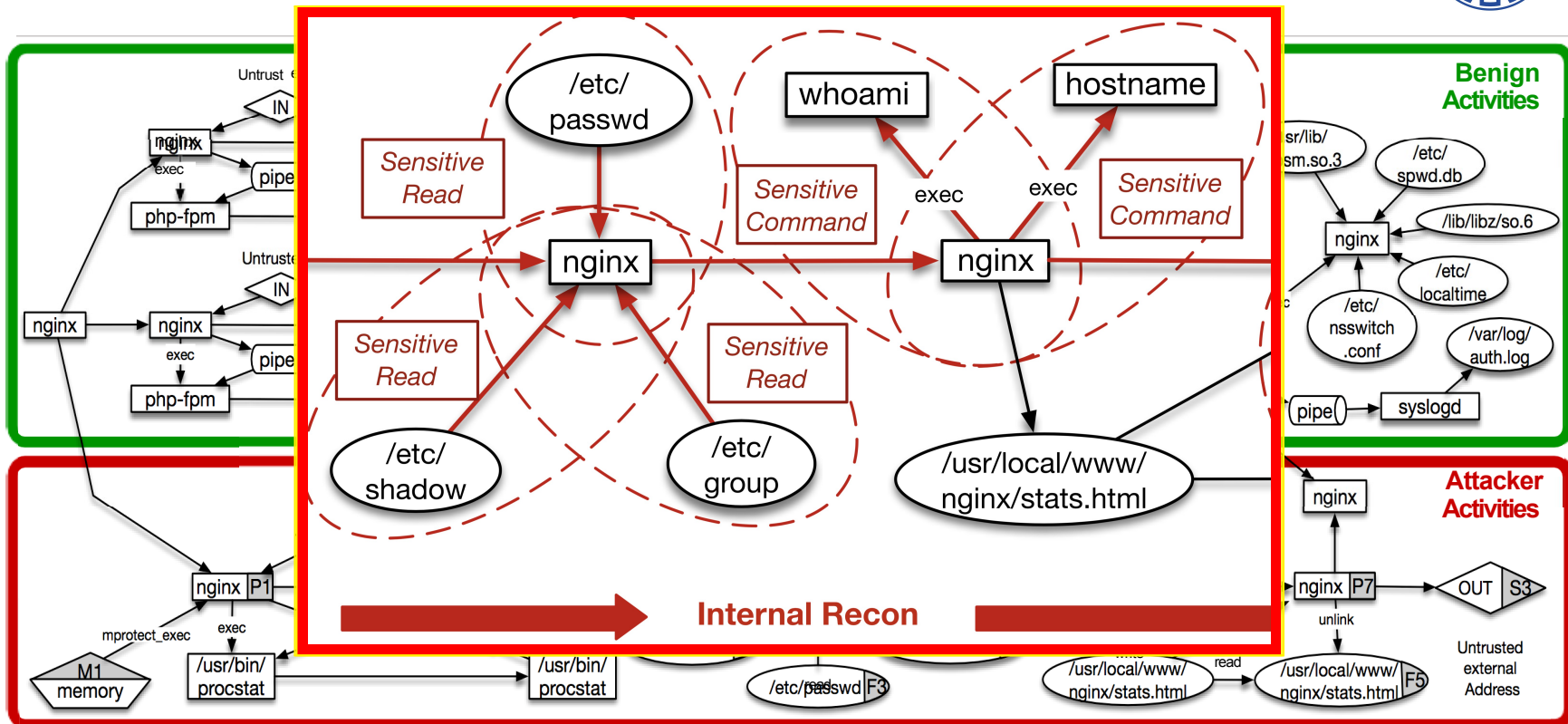


Illustrative Example



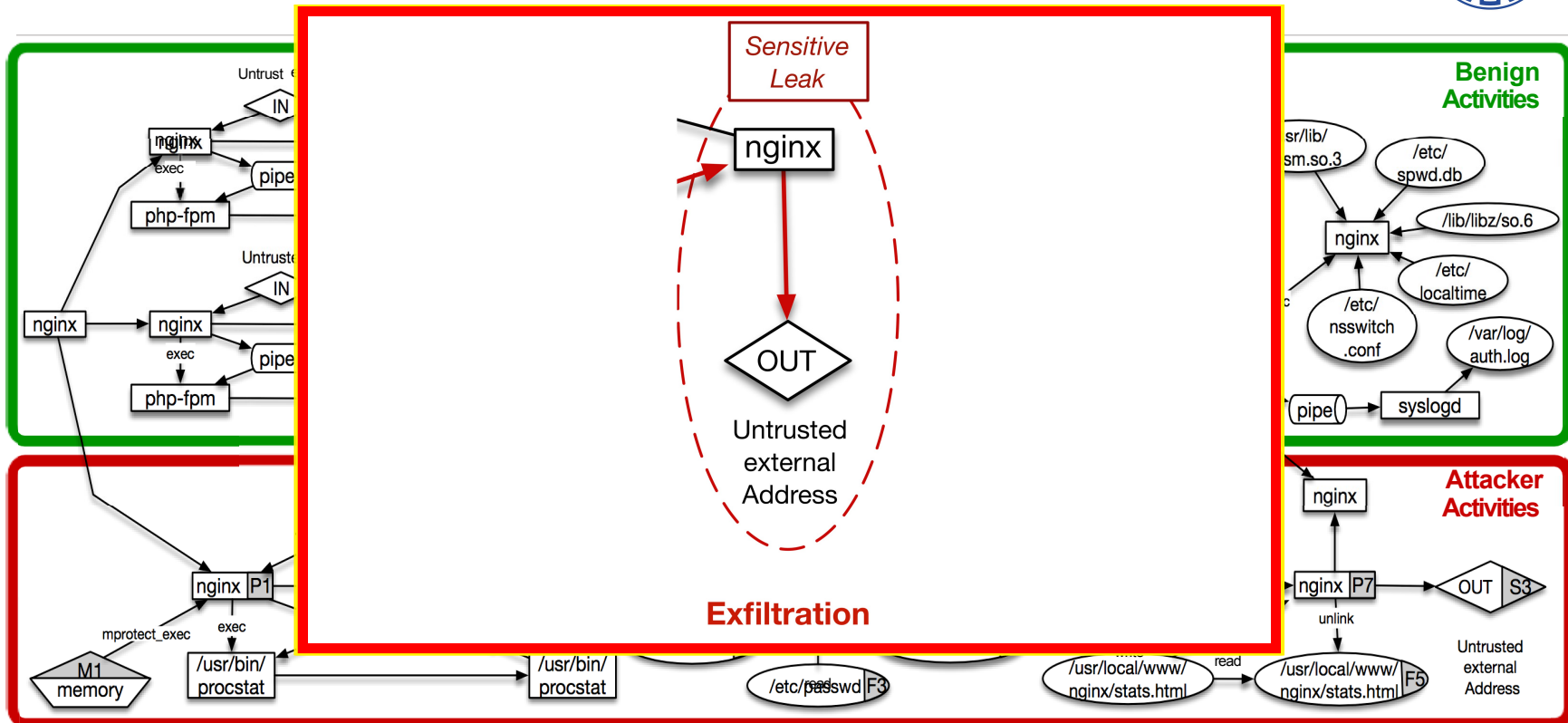


Illustrative Example



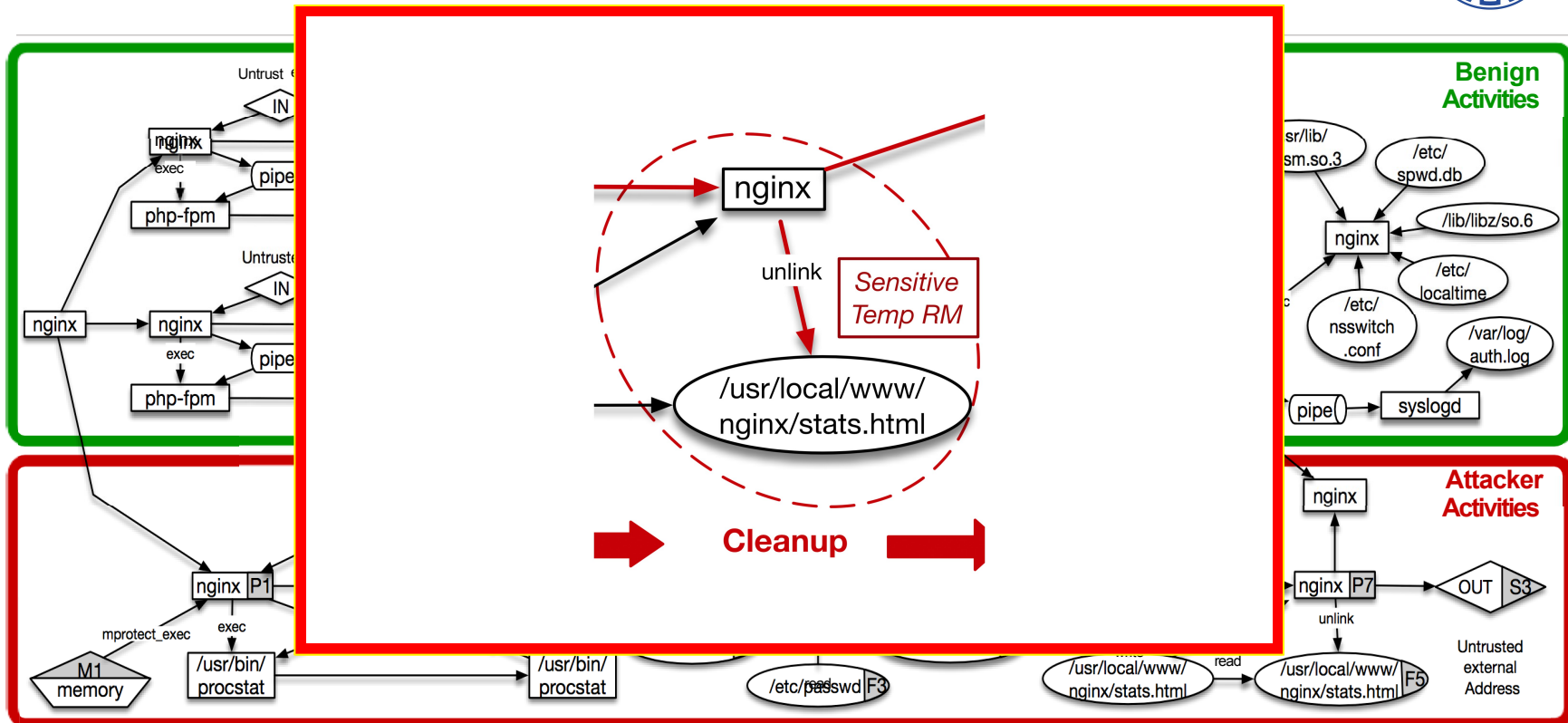


Illustrative Example

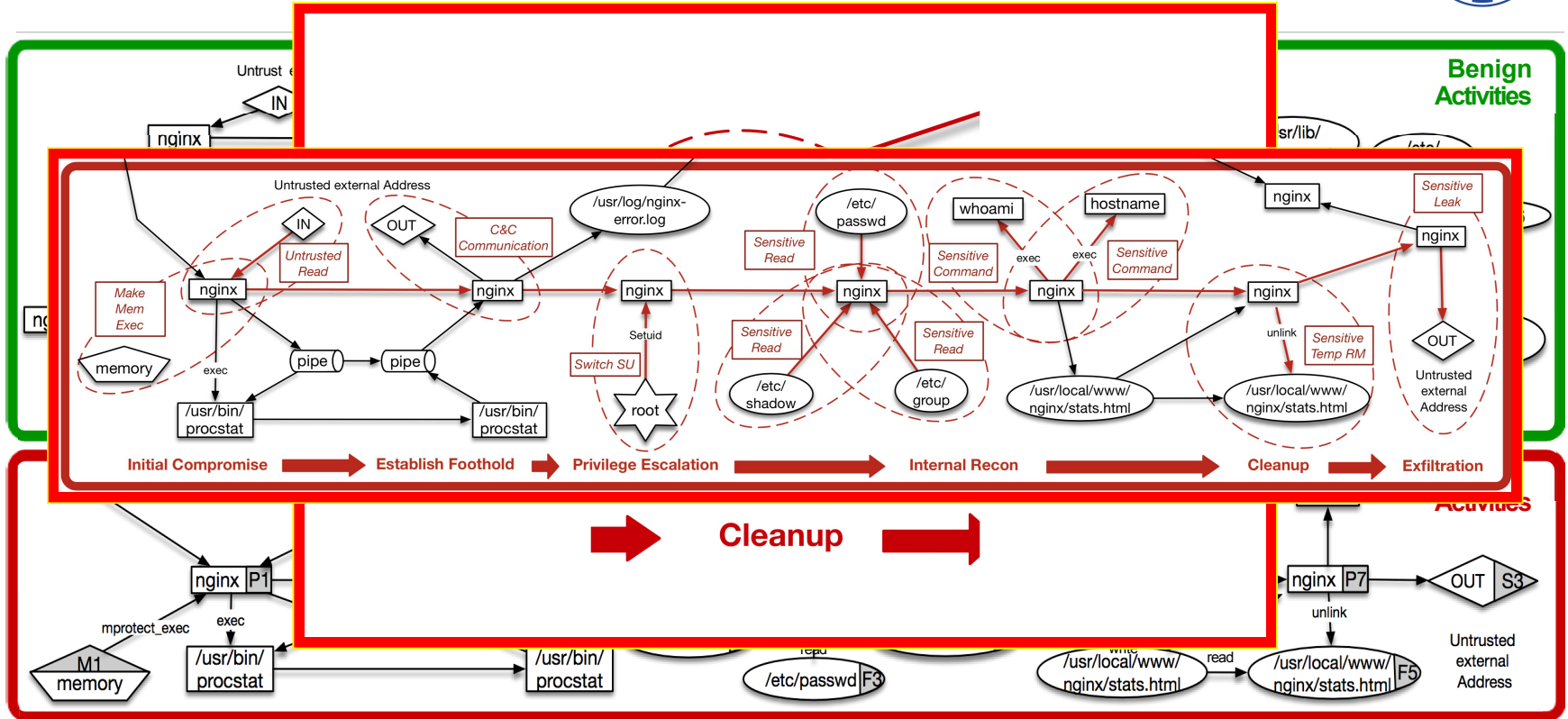




Illustrative Example



Illustrative Example

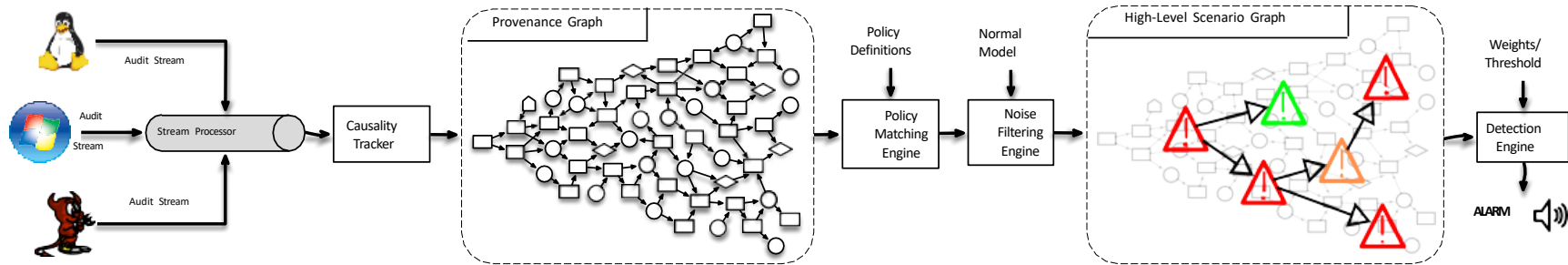


Benign Activities

Activities



Holmes Architecture



- Develop TTP specifications over audit logs
- Use specifications to detect TTPs
- Filter noise based on data quantities of benign information flows, measured in bytes transferred
- Construct high-level graph (HSG) that correlates individual alerts/TTPs
- Derive campaign detection signal from graph



Example TTP specifications

APT Stage	TTP	Event Family	Events	Severity	Prerequisites
<i>Initial_Compromise(P)</i>	<i>Untrusted_Read(S, P)</i>	READ	FileIoRead (Windows), read/pread/readv/preadv (Linux,BSD)	L	$S.ip \notin \{\text{Trusted_IP_Addresses}\}$
	<i>Make_Mem_Exec(P, M)</i>	MPROTECT	VirtualAlloc (Windows), mprotect (Linux,BSD)	M	$\$PROT_EXEC\$ \in M.flags$ $\wedge \exists Untrusted_Read(?, P') :$ $path_factor(P', P) \leq path_thres$
<i>Establish_Foothold(P)</i>	<i>Shell_Exec(F, P)</i>	EXEC	ProcessStart (Windows), execve/fexecve (Linux,BSD)	M	$F.path \in \{\text{Command_Line_Utilities}\}$ $\wedge \exists Initial_Compromise(P') :$ $path_factor(P', P) \leq path_thres$

TABLE 4. Example TTPs. In the Severity column, L=Low, M=Moderate, H=High, C=Critical. Entity types are shown by the characters: P=Process, F=File, S=Socket, M=Memory, U=User.

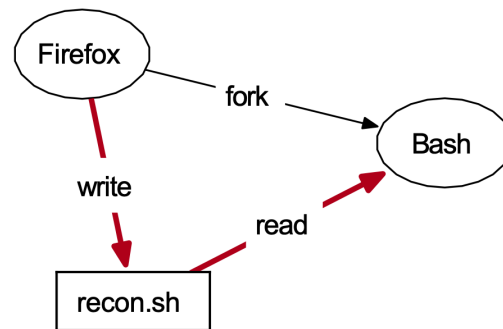
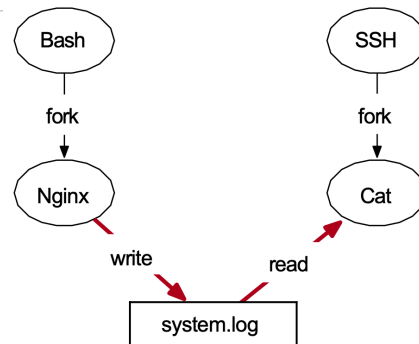


Avoiding spurious dependencies

- Spurious dependencies can result in dependence explosion
- Addressed by asking a key question: what is the influence that attacker had in creating a dependency?
- Key notion Ancestor cover for f : set of all processes that influence a dependency f .

$\forall p \in f \exists a \in AC(f) \ a = p$ or a is an ancestor of p

- Minimal Ancestor cover for f - corresponds to the minimum number of processes attacker should exploit to influence a dependency f .





Avoiding spurious dependencies (Cont.)

$$path_factor(N_1, N_2) = \min_{\forall f: f.src=N_1, f.dst=N_2} AC_{min}(f)$$

- *path_factor* value computed incrementally in real-time

APT Stage	TTP	Event Family	Severity	Prerequisites
<i>Complete_Mission(P)</i>	<i>Sensitive_Leak(P, S)</i>	SEND	H	$S.ip \notin \{\text{Trusted_IP_Addresses}\}$ $\wedge \exists \text{Internal_Reconnaissance}(P') :$ $path_factor(P', P) \leq path_thres$ $\wedge \exists \text{Initial_Compromise}(P'') :$ $path_factor(P'', P) \leq path_thres$

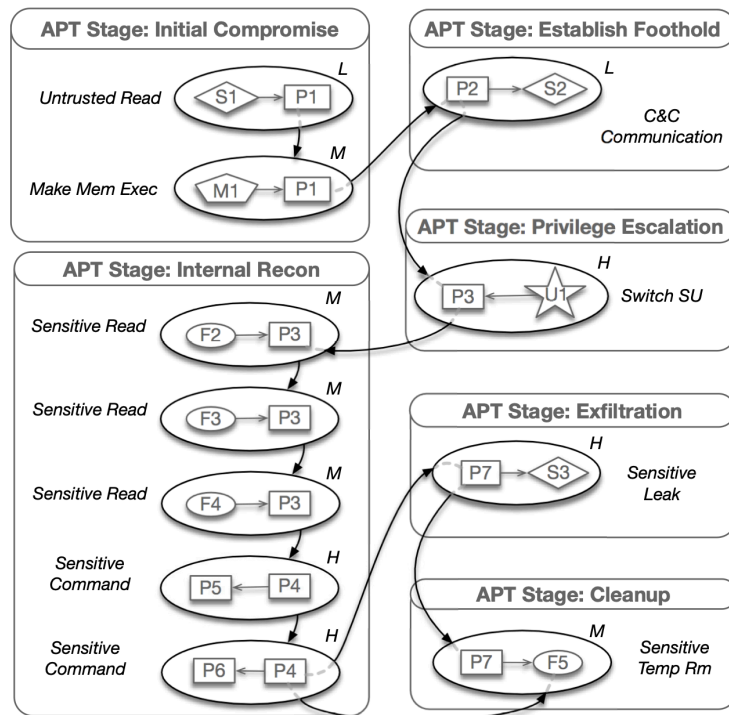
Value of *path_thres* could be set based on the threat an organization is preventing from

- we assume attacker is not willing or capable to compromise more than 3 exploits.



Signal Correlation, HSG, and Threat Triples

- A TTP is matched and added to the HSG if all its prerequisites are satisfied.
- HSG \rightarrow Threat Tuple: represents various stages of an APT campaign.
- Each element in tuple takes on severity levels $\langle M, L, H, H, -, H, M \rangle$
- HSG provides a compact, visual summary of the campaign at any moment.
- cyber-analyst can quickly infer the big picture of the attack (scope and magnitude)





HSG Ranking and Prioritization

- Severity level transformed to a number based on NIST severity score mappings

Qualitative level	Quantitative Range	Rounded up Average Value
Low	0.1 - 3.9	2.0
Medium	4.0 - 6.9	6.0
High	7.0 - 8.9	8.0
Critical	9.0 - 10.0	10.0

- Tuple transformed into numeric value as weighted product

$$\prod_{i=1}^n (S_i)^{w_i} \geq \mathcal{T} \quad w_i = \frac{10 + i}{10}$$

- Alert raised based on threshold learned from benign activity data

- $\langle C, M, -, H, -, H, M \rangle \rightarrow \langle 10, 6, 1, 8, 1, 8, 6 \rangle \rightarrow 1163881$



Evaluation Datasets

Dataset 1: Using this dataset, we measure the optimal threshold value

Stream No.	Duration	Platform	Scenario No.	Scenario Name	Attack Surface
1	0d1h17m	Ubuntu 14.04 (64bit)	1	Drive-by Download	Firefox 42.0
2	2d5h8m	Ubuntu 12.04 (64bit)	2	Trojan	Firefox 20.0
3	1d7h25m	Ubuntu 12.04 (64bit)	3	Trojan	Firefox 20.0
4	0d1h39m	Windows 7 Pro (64bit)	4	Spyware	Firefox 44.0
5	5d5h17m	Windows 7 Pro (64bit)	5.1	Eternal Blue	Vulnerable SMB
			5.2	RAT	Firefox 44.0
6	2d5h17m	FreeBSD 11.0 (64bit)	6	Web-Shell	Backdoored Nginx
7	8d7h15m	FreeBSD 11.0 (64bit)	7.1	RAT	Backdoored Nginx
			7.2	Password Hijacking	Backdoored Nginx

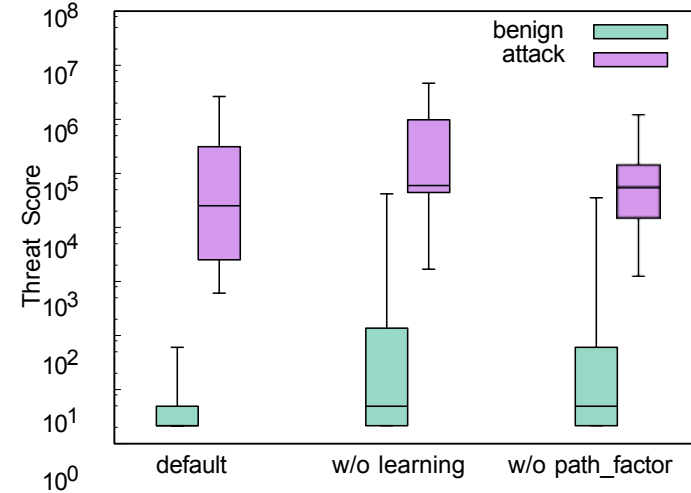
Dataset 2: live detection in a setting that we have no prior knowledge of when or how red-team is conducting the attacks.

- After this experiment, dataset has been released publicly.

Evaluation

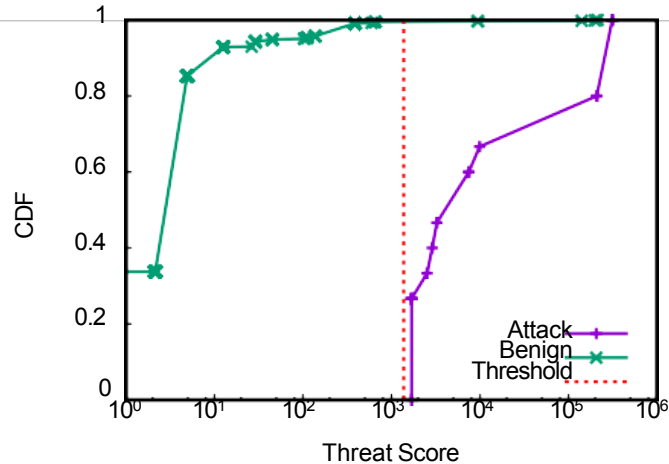
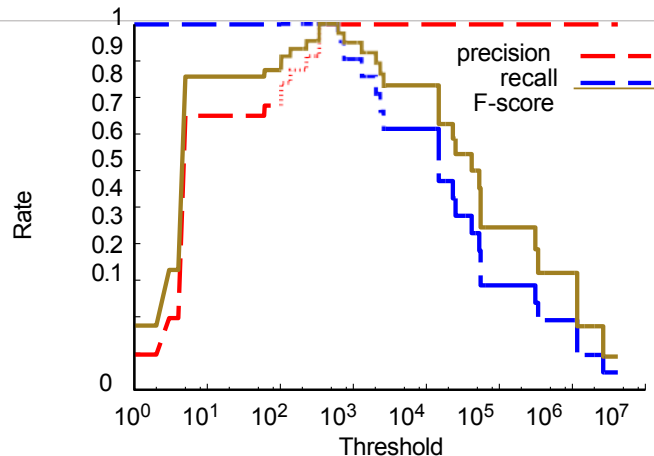


Scenario No.	Threat Tuple	Threat Score	Highest Benign Score in Dataset
1	$\langle C, M, -, H, -, H, M \rangle$	1163881	61
2	$\langle C, M, -, H, -, H, - \rangle$	55342	226
3	$\langle C, M, -, H, -, H, M \rangle$	1163881	338
4	$\langle C, M, -, H, -, -, M \rangle$	41780	5
5.1	$\langle C, L, -, M, -, H, H \rangle$	339504	104
5.2	$\langle C, L, -, -, -, -, M \rangle$	608	
6	$\langle L, L, H, M, -, H, - \rangle$	25162	137
7.1	$\langle C, L, H, H, -, H, M \rangle$	4649220	133
7.2	$\langle M, L, H, H, -, H, M \rangle$	2650614	





Optimal threshold Value & Live Experiment Results



- F-score maximum at [338.25, 608.26] for 6 APT stages
- Average severity of each APT step = 2.09
- Threshold set for Live experiment (7 APT stages): $2.09^{9.8} = 1378$
- A few false positive: system administrator connecting via SSH



Summary

- Presented a real-time APT detection system that correlates TTPs that might be used to carry out each APT stage.
- visualize high-level APT behavior in real time.
- Dependence explosion mitigation by using the concept of minimum ancestral cover
- Benign system activities pruning based on data quantities in the flow of information
- Experiments show high accuracy and performance for Holmes
- Effectiveness evaluated using a live experiment w/o having prior knowledge of attacks.



Acknowledgments

- [MITRE]<https://attack.mitre.org/>
- [Holmes] HOLMES: Real-Time APT Detection through Correlation of Suspicious Information Flows, S. Momeni Milajerdi, R. Gjomemo, B. Eshete, R. Sekar, V. N. Venkatakrishnan, IEEE Symposium on Security and Privacy 2019.