

Scalable Data Analytics Pipeline for Real-Time Attack Detection; Design, Validation, and Deployment in a Honeypot Environment

Eric Badger

Master's Student

Computer Engineering





Overview

- Introduction/Motivation
- Challenges
- Pipeline Design
- Pipeline Deployment
- Validation of Alerts and Attack Detection Tools
- Future Work
- Conclusion



Research Problem

Our goal is to detect potential attacks as early as possible. Security analysts attempt to detect and prevent attacks, but they can't analyze everything in their infrastructure by hand. They need tools to automate the analysis for early detection of attacks.

- How do we transition attack detection models from theory to practice?
- How do we validate that the alerts we are using are useful?
 - Does combining alerts from different monitors make the attack detection better?
 - Is the extra performance overhead worth it?
- How do we validate that our attack detection model is adequate and better than others?
 - What models are suitable for real-time attack detection in practical deployment?

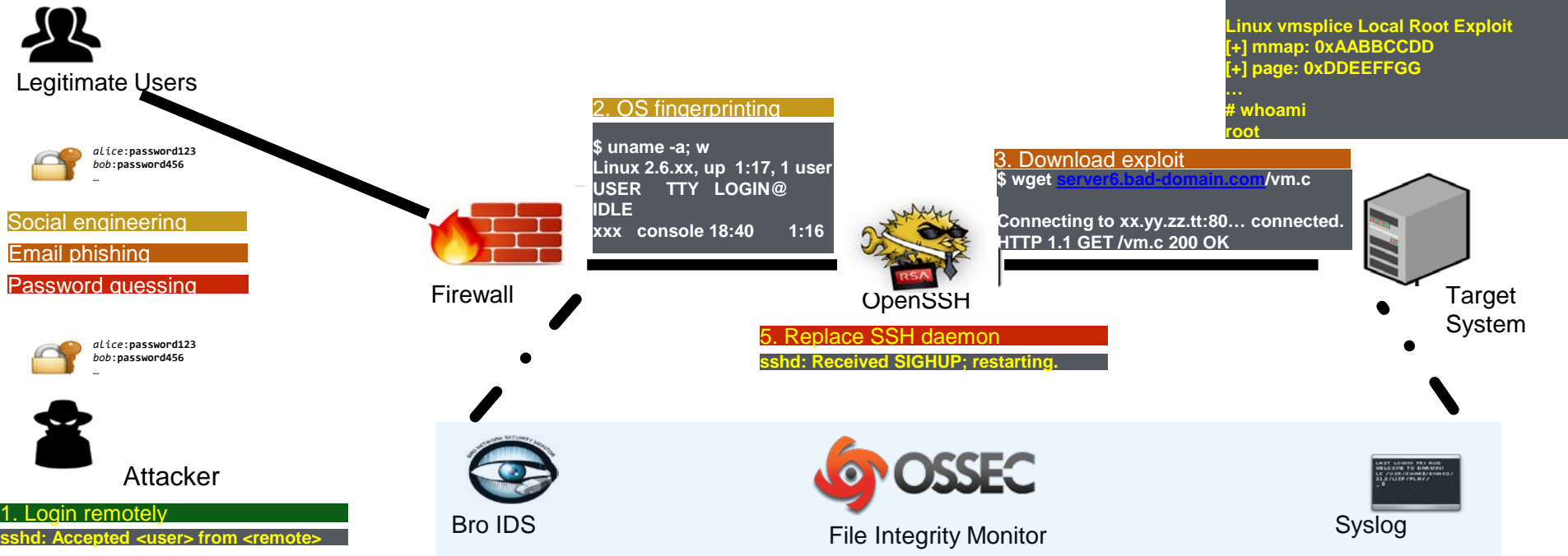


Research Challenges: Transitioning Attack Detection from Theory to Practice

- Identifying which alerts are useful for attack detection
- Normalizing all logs to a common format
- Achieving both high-accuracy and real-time attack detection
- Achieving high-accuracy attack detection in the face of alert randomness, noise, and imperfect monitors
- Scaling the data pipeline

The chain of tools used for data-driven attack detection

Example Attack Scenario





How to Extract Important Alerts

- **Network Monitors**

- Bro

- Network IDS used for packet analysis

- CriticalStack Intel Feed

- **Host Monitors**

- OSSEC

- Runs periodic system checks and file integrity monitoring

- Aggregates and correlates all other host alerts

- Snoopy Logger

- Logs all execv system calls

- RKHunter

- Searches for rootkits, hidden folders/files/ports, and other system issues

- Syslogs

- Normal GNU/Linux “/var/log” logs, such as auth.log, kern.log, dpkg.log, and others

- Bash Logs

- Logs Bash history as the commands are executed

Log Normalization and Aggregation

Auth Logs

```
118 2015-09-29T16:54:39.425503-05:00 whitacre sudo: pam_unix(sudo:session): session opened for user root by eric(uid=1000)
119 2015-09-29T16:54:39.499108-05:00 whitacre su[10044]: Successful su for root by root
120 2015-09-29T16:54:39.499118-05:00 whitacre su[10044]: + /dev/pts/6 root:root
121 2015-09-29T16:54:39.499296-05:00 whitacre su[10044]: pam_unix(su:session): session opened for user root by eric(uid=0)
122 2015-09-29T17:00:00.000000-05:00 whitacre pam_unix(cron:session): session opened for user root by eric(uid=0)
123 2015-09-29T17:00:00.000000-05:00 whitacre pam_unix(cron:session): session closed for user root
```

1443126106.661021

Bro Notice Logs

```
1 #separator \\\n2 #set_separator ,\n3 #empty_field (empty)\n4 #unset_field -\n5 #path notice\n6 #open 2015-09-24-15-20-50\n7 #fields ts uid id.orig_h id.orig_p id.resp_h id.resp_p fuid file_mime_type file_desc proto note msg src dst p n peer_descr actions suppress_for_dropped remote_location.country_code remote_location.region remote_location.city remote_location.latitude remote_location.longitude\n8 #types time string addr port addr port string string string enum enum string string addr addr port count string set[enum] interval boolean string string string double double\n9 1443126106.661021 Ceu4Bv3Kd6Pr6MiUa6 130.126.137.23 58190 54.231.133.204 443 - - - tcp SSL::Invalid Server Cert SSL certificate validation failed with (unable to get local issuer certificate) CN=*.s3-eu-west-1.amazonaws.com, OU=S3-B, O=Amazon.com\\, Inc., L=Seattle, ST=Washington, C=US 130.126.137.23 54.231.133.204 443 - bro Notice::ACTION LOG 3600.000000 F - - - -
```

OSSEC Logs

```
3763 ** Alert 1443563681.247829: - syslog,sudo\n3764 2015 Sep 29 16:54:41 whitacre->/var/log/auth.log\n3765 Rule: 5402 (level 3) -> 'Successful sudo to ROOT executed'\n3766 User: eric\n3767 2015-09-29T16:54:39.425299-05:00 whitacre sudo: eric : TTY=pts/6 ; PWD=/var/log ; USER=root ; COMMAND=/bin/su\n3768\n3769 ** Alert 1443563681.248105: mail - local,syslog,\n3770 2015 Sep 29 16:54:41 whitacre->/var/log/auth.log\n3771 Rule: 105501 (level 11) -> 'User successfully changed UID to root.'\n3772 User: eric
```

2015-09-29T08:00:06.257580-05:00

RKHunter Logs

```
5582 [16:02:07] System checks summary\n5583 [16:02:07] =====\n5584 [16:02:07]\n5585 [16:02:07] File properties checks...\n5586 [16:02:07] Files checked: 142\n5587 [16:02:07] Suspect files: 2\n5588 [16:02:07]\n5589 [16:02:07] Rootkit checks...\n5590 [16:02:07] Rootkits checked : 380\n5591 [16:02:07] Possible rootkits: 0
```

Snoopy Logs

```
4141 2015-09-29T08:00:06.252345-05:00 whitacre snoopy[32190]: [username:root tty username:(none) uid:0 sid:26590 tty:(none) cwd:/root filename:/bin/uname]: uname -r\n4142 2015-09-29T08:00:06.254930-05:00 whitacre snoopy[32194]: [username:root tty username:(none) uid:0 sid:26590 tty:(none) cwd:/root filename:/bin/grep nl: grep ^\n4143 2015-09-29T08:00:06.257580-05:00 whitacre snoopy[32197]: [username:root tty username:(none) uid:0 sid:26590 tty:(none) cwd:/root filename:/bin/egrep]: egrep (^[[^\]])[!]*}
```



Log Normalization and Aggregation (2)

- Since the logs are all in different formats, they need to be normalized to a common format

Normalized Log

```
209 1443457145.717,143.219.0.11:root,ALERT_FAILED_PASSWORD,NaN,NaN,1443461619.48
3
210 1443457147.510,143.219.0.11:root,ALERT_FAILED_PASSWORD,NaN,NaN,1443461619.49
0
211 1443457661.505,143.219.0.11:LOGIN,login,NaN,NaN,1443461619.516
212 1443457661.469,143.219.0.11:root,read_host_configuration,NaN,NaN,1443461619.
520
213 1443457662.963,143.219.0.11:ubuntu,ALERT_GET_LOGGEDIN_USERS,NaN,NaN,14434616
19.536
214 1443461754.305,,:,ALERT_INTERNAL_ADDRESS_SCAN,NaN,NaN,1443461764.436
```

All logs needed to be centralized so that we can act on them

TLS/SSL encryption is necessary to secure the movement of logs through the pipeline

If not, the logs could be added, deleted, or changed by a MITM attack



Pipeline Design

**Example
Tools**

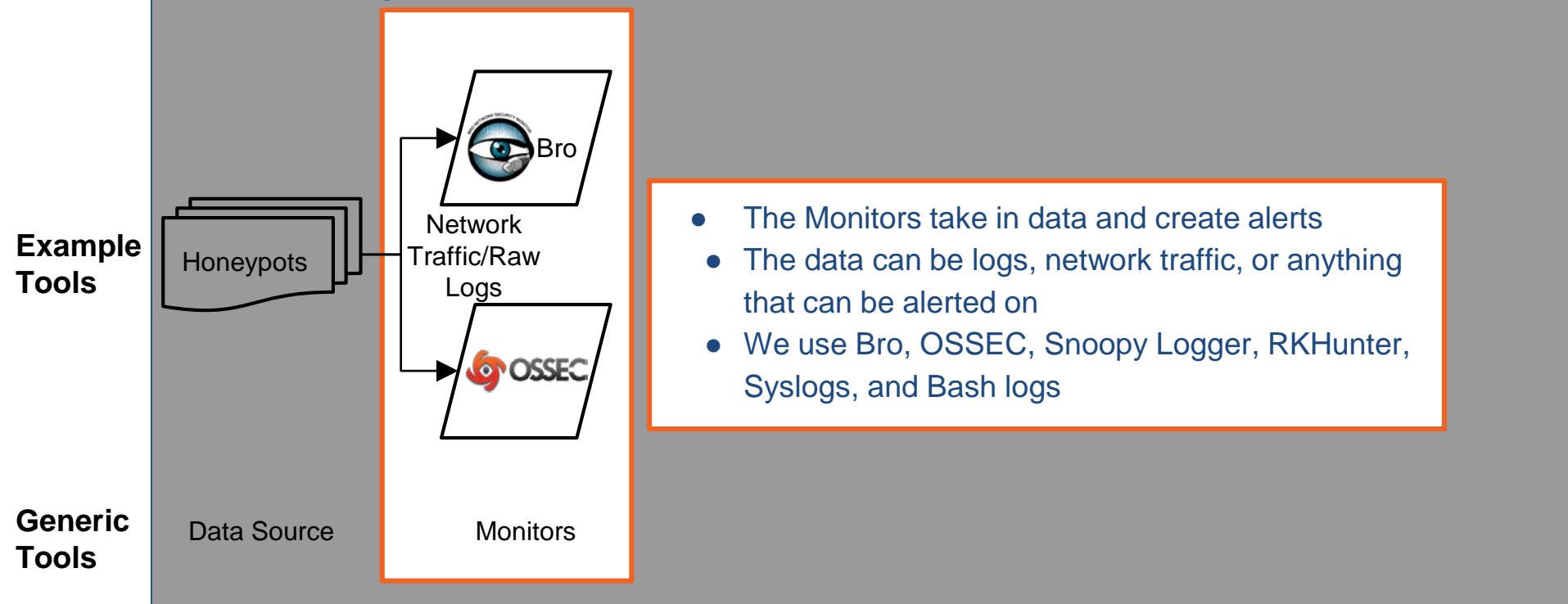


**Generic
Tools**

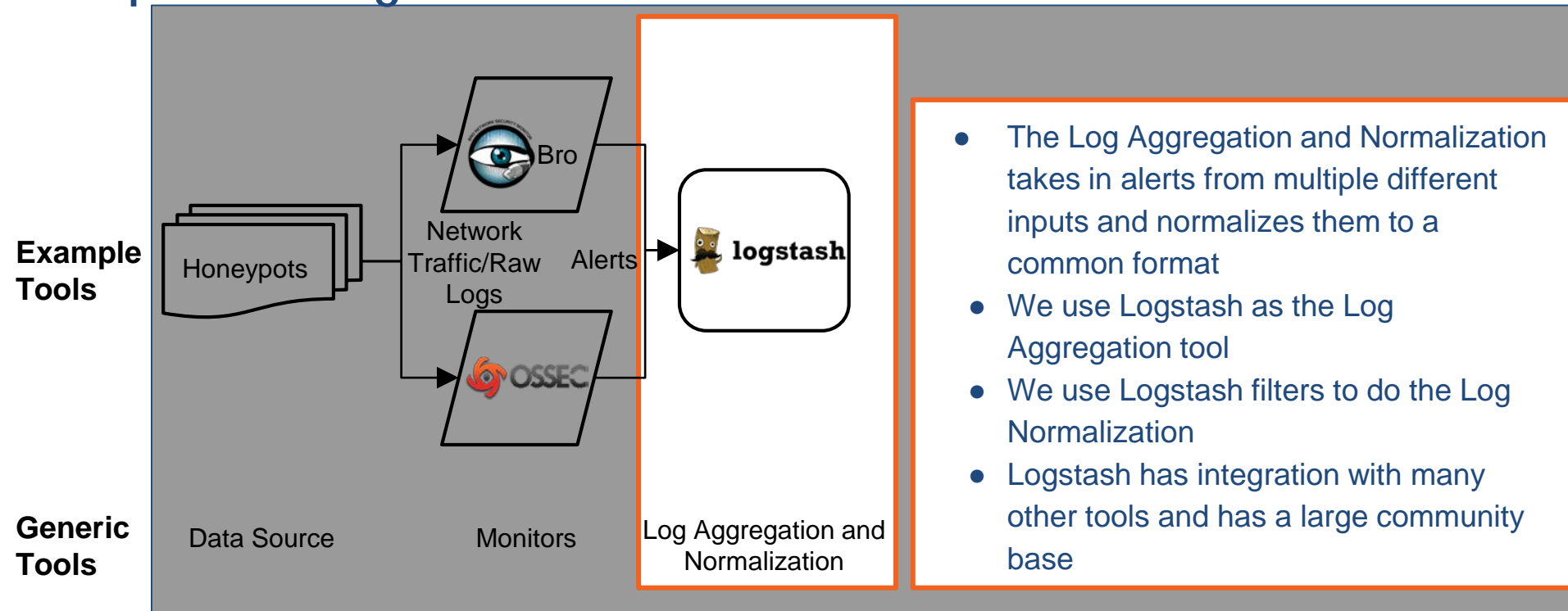
Data Source

- The Data Source can be any sort of online or offline computer or device
 - Online
 - Honeyd, servers, workstations, phones
 - Offline
 - Security testbed, old logs
- We use customized honeypots deployed at the NCSA

Pipeline Design

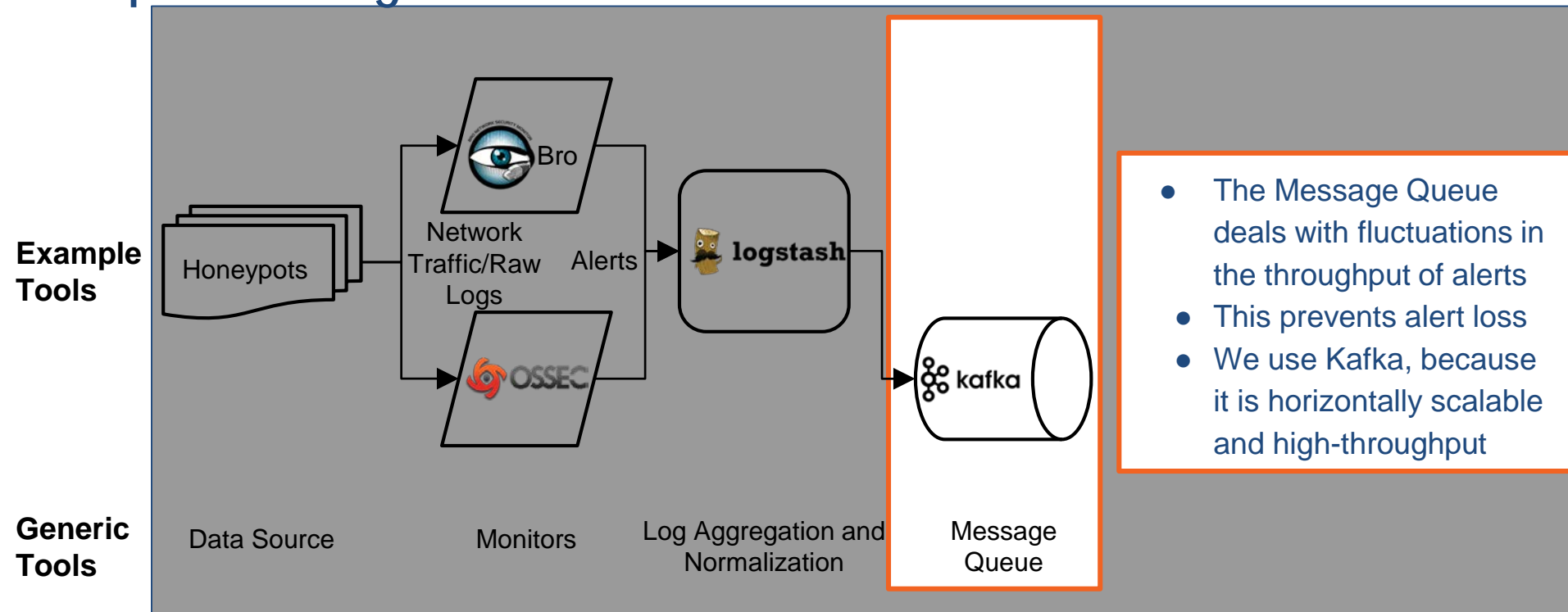


Pipeline Design





Pipeline Design

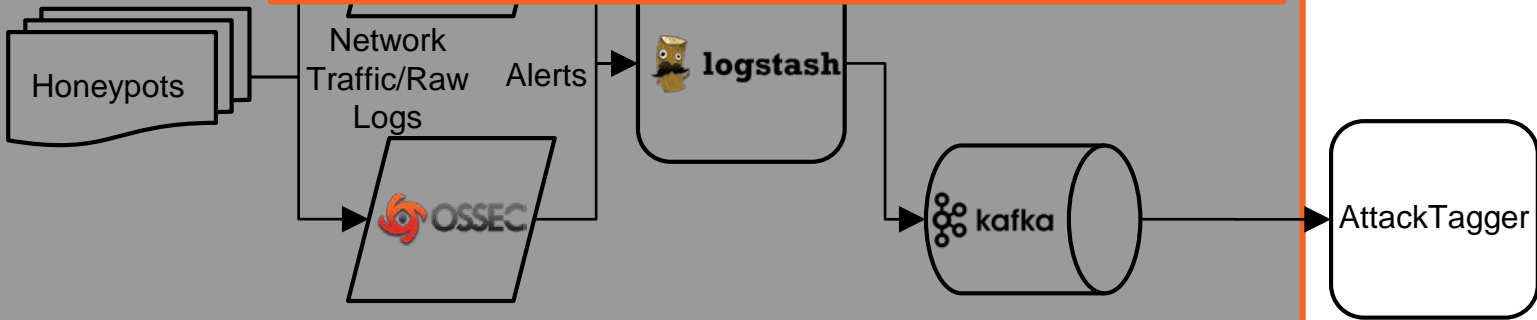




Pipeline Design

- The Attack Detection tool takes in alerts from the Message Queue and does analysis to detect attacks
- We use AttackTagger, which is an attack detection tool based on factor graphs

Example Tools



Generic Tools

Data Source

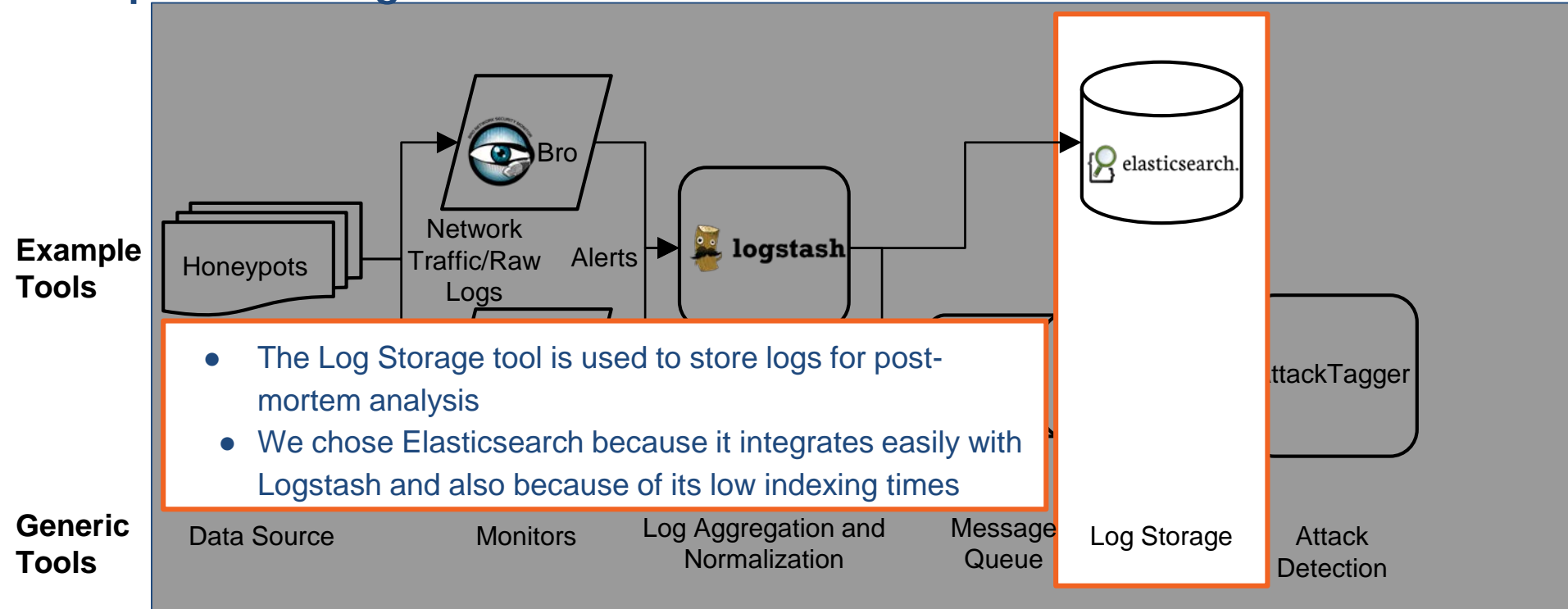
Monitors

Log Aggregation and Normalization

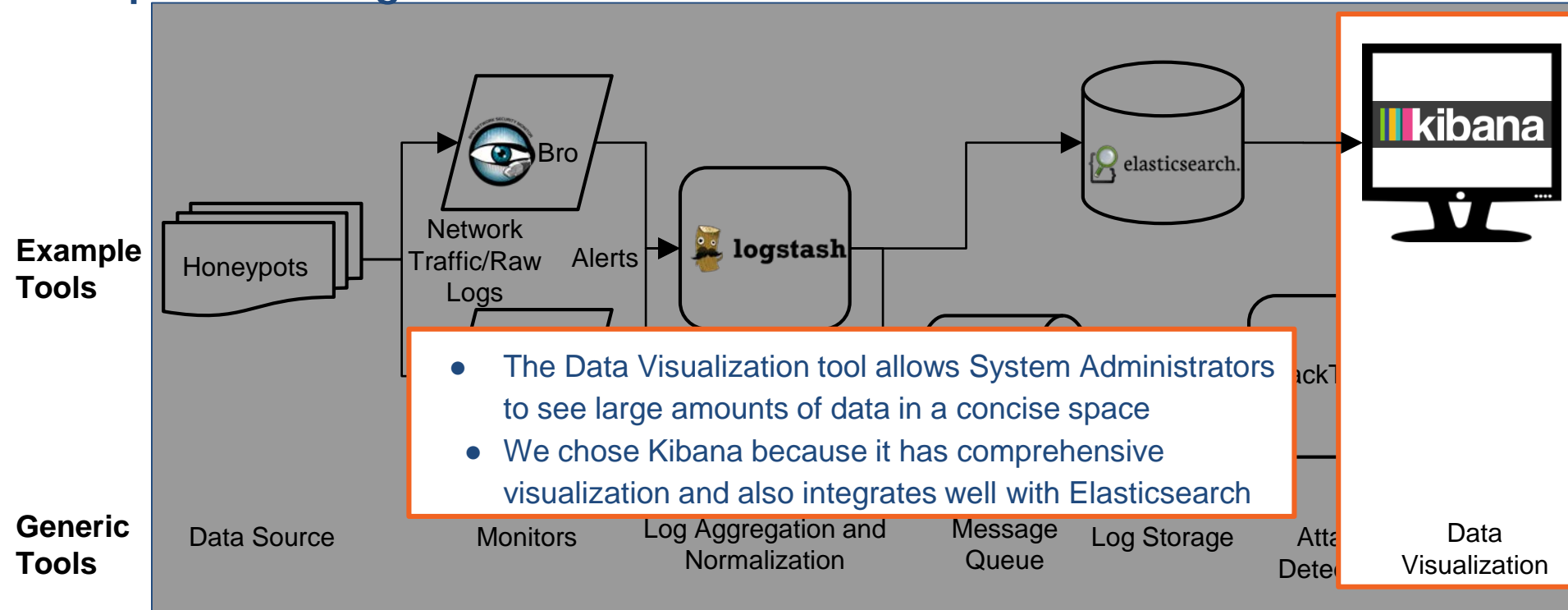
Message Queue

Attack Detection

Pipeline Design



Pipeline Design

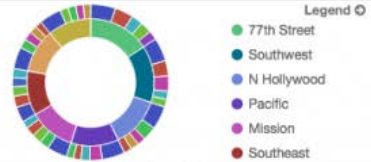


LAPD: Crime Dashboard

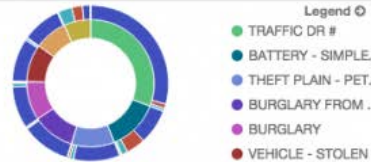
LAPD: Tile map



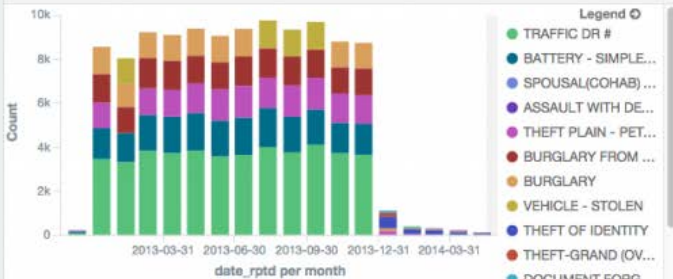
LAPD: Crime Type by Area



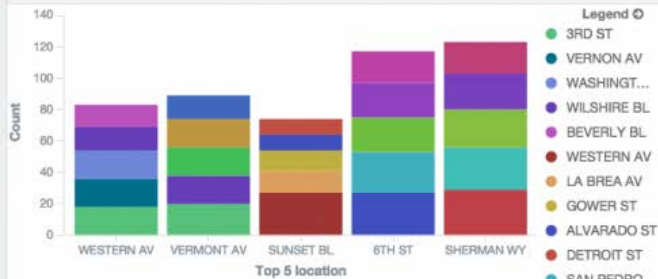
LAPD: Crime Status by Crime Type



LAPD: Monthly histogram by Crime Type

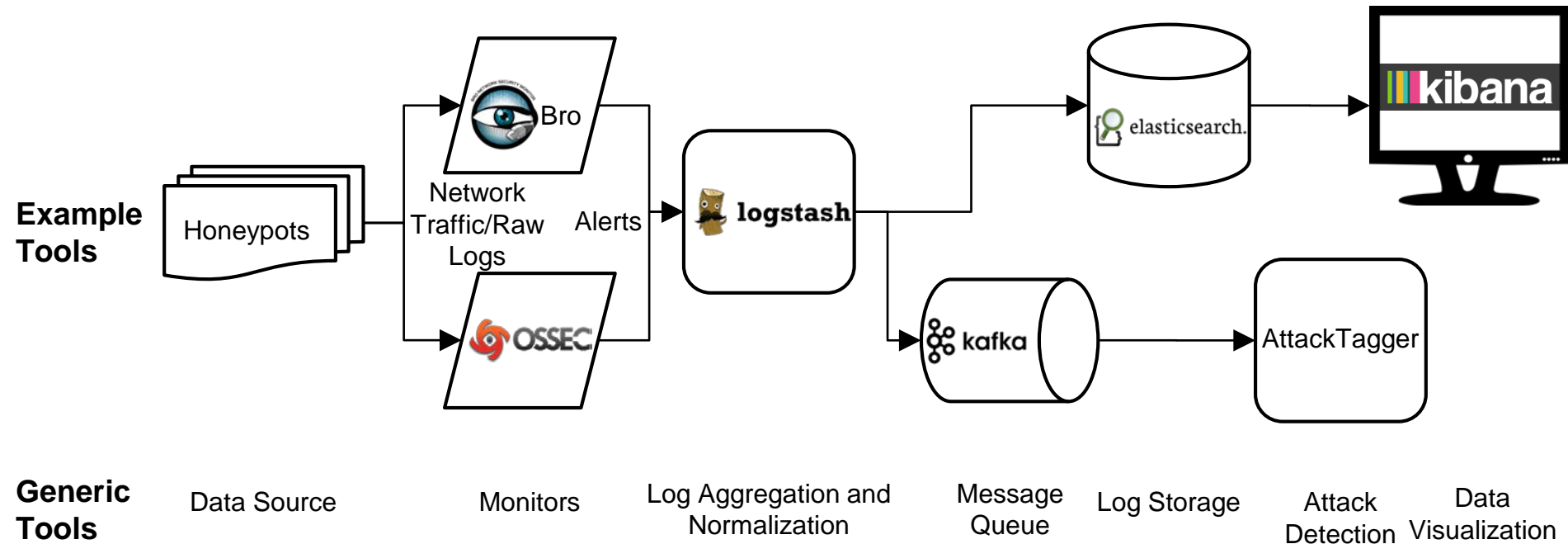


LAPD: Intersections



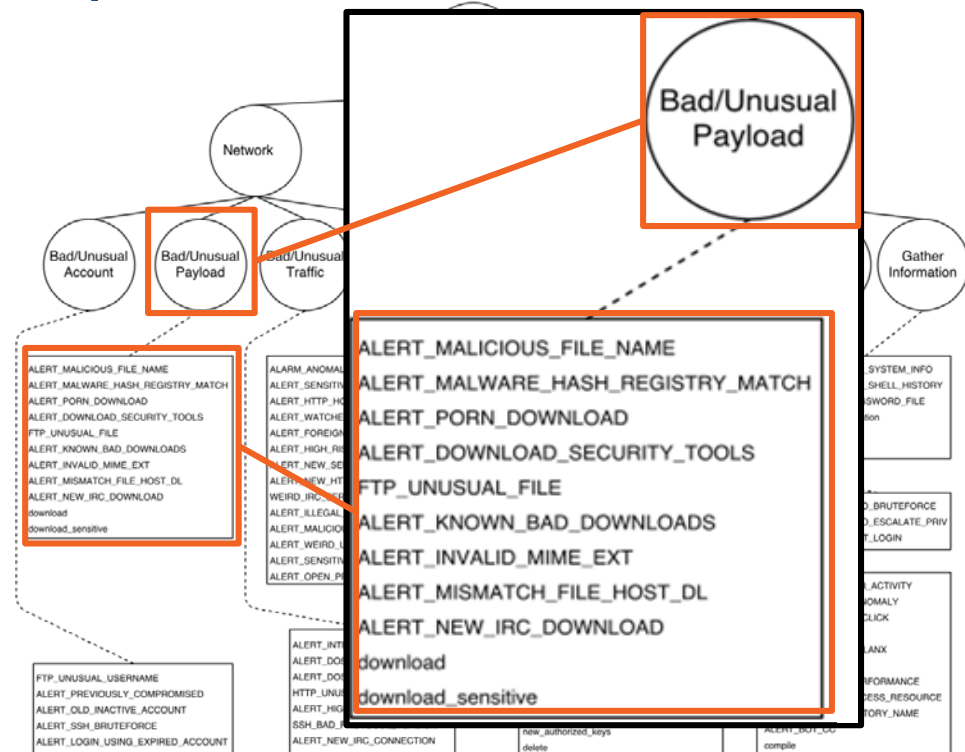


Pipeline Design



How Do I Know What Alerts Are Bad/Important?

- Research was done in [1] and [2] that studied attacks over a six-year period at NCSA. This research identified important alerts related to these attacks and developed the AttackTagger detection tool
- We utilized and extended a custom set of monitors to create alerts corresponding to the inputs that were used in AttackTagger
- In essence, we brought AttackTagger from a theoretical tool to actual deployment



[1] Phuong Cao, Key-whan Chung, Zbigniew Kalbarczyk, Ravishankar Iyer, and Adam J. Slagell. 2014. Preemptive intrusion detection. HotSoS '14.
[2] Phuong Cao, Eric Badger, Zbigniew Kalbarczyk, Ravishankar Iyer, and Adam Slagell. 2015. Preemptive intrusion detection: theoretical framework and real-world measurements. HotSoS '15.

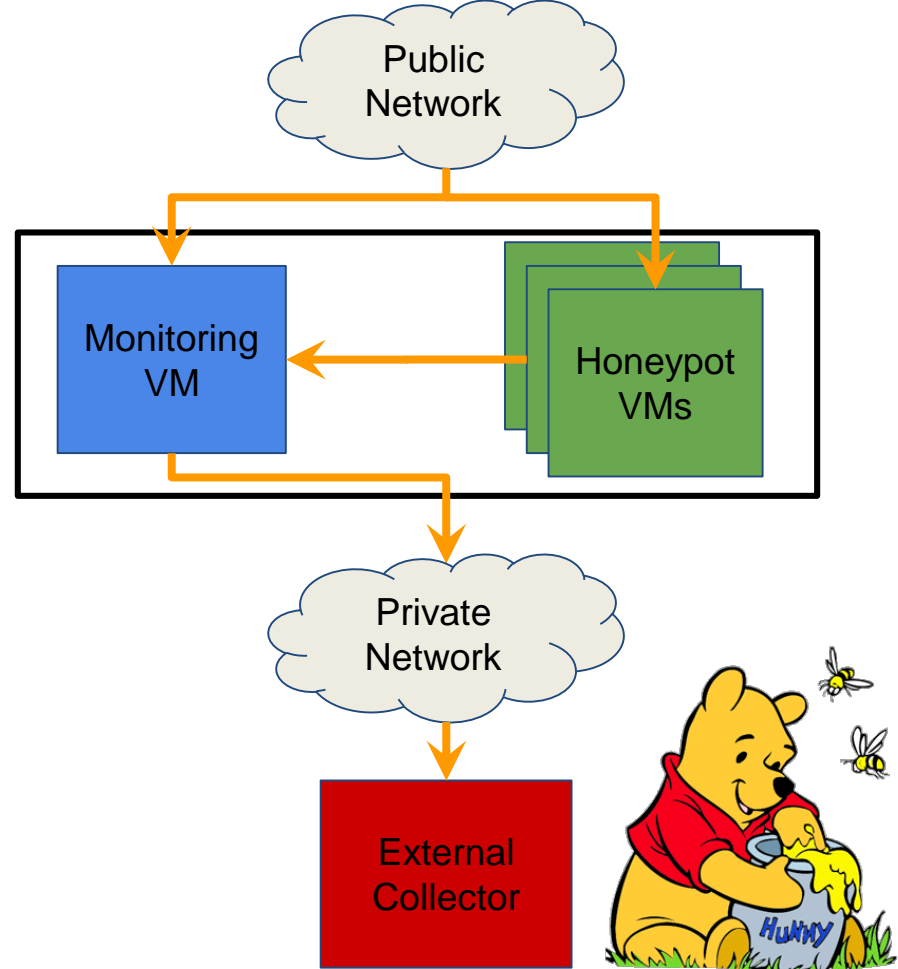


What Can We Do with This Pipeline?

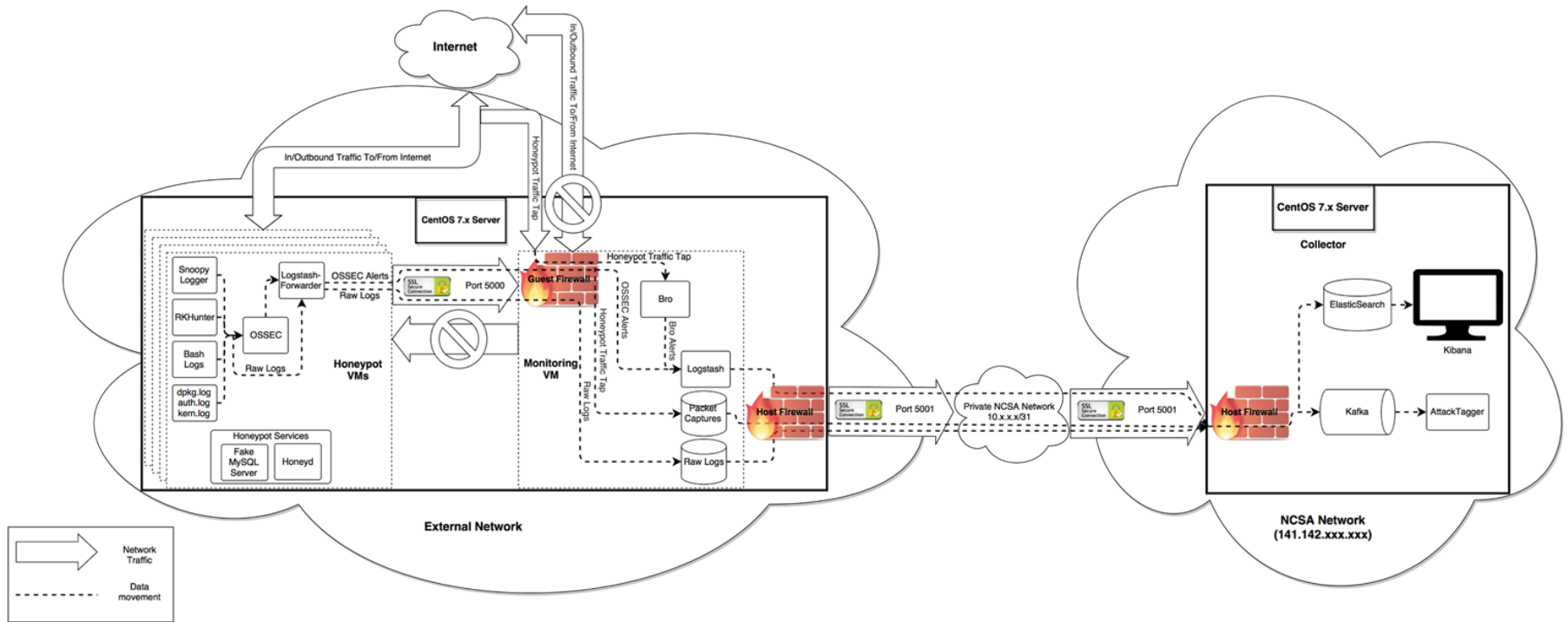
- Both online and offline deployment
- Online
 - Analysis of attacks happening on the infrastructure
 - Analysis of attack detection tools on live data
- Offline
 - Post-mortem log analysis (via Elasticsearch/Kibana)
 - Analysis of old attacks
 - Development of attack detection tools
 - Validation of alerts

Honeypots at NCSA

- NCSA server running several VMs
 - Honeypot VMs
 - Open to public
 - Monitoring VM
 - Allows TCP Port 5000 (Logstash) from honeypots
 - Allows TCP Port 22 from NCSA, UI, and UI wireless
 - Sends logs to Collector via Private Network
- Collector
 - Allows TCP Port 5001 (Logstash) from private network
 - Allows TCP Port 22 from NCSA, UI, and UI wireless



Honeypots at NCSA (2)





Preliminary Honeypot Results

- 3 separate SSH brute-force attacks successfully compromised one of the honeypots in the first 3 days
- Appeared to download and execute either an open proxy or a DDoS attack through the program “/tmp/squid64”
- They beat my monitors! (Well, sort of...)
 - They pushed their malware from the anomalous host instead of pulling it from the honeypot
 - They deleted the malware immediately after running it, so it was not seen by OSSEC’s file-integrity monitoring



How to Validate Importance of Inputs (Alerts)

- Mix and match which monitors/alerts that we use in our attack detection
- Evaluate the difference in attack detection coverage and accuracy
 - Adding monitors/alerts will likely add detection coverage because of extra data
 - Adding monitors/alerts could possibly decrease detection accuracy because of additional noise
- Determine whether the difference in detection coverage is worth the additional overhead



How to Validate Accuracy of Outputs (Detection Tools)

- Compare and contrast different attack detection tools
e.g. Factor Graphs, Bayesian Networks, Markov Random Fields,
Signature Detection, etc.
Which are most accurate?
Which are least complex?
- In the pipeline, attack detection tools are plug and play as long as they can read the normalized alert format
If they can't, a translation filter can be added

Future Work

- Validate data pipeline inputs (alerts) and outputs (attack detection tools)
- Add additional data types to data pipeline

Netflows

Full file-integrity monitoring (e.g. Tripwire)

Administrator-generated alerts/profiles

Keystroke data (e.g. iSSHD)

- Convert detection model into a stream-processing system
 - Detectors such as AttackTagger are currently batch processing detectors
 - We need to process the data in real-time
- Transition entire pipeline into practice at NCSA production system

Conclusion

- Showed how to transition attack detection software from theory to practice
- Showed how to evaluate the effectiveness of the inputs (alerts) and outputs (attack detection tools) of the pipeline
- Identified challenges and how to overcome them

Special Thanks

- Phuong Cao
- Alex Withers
- Adam Slagell
- NCSA
- DEPEND Research Group

Questions?

Citations

- [1] Phuong Cao, Key-whan Chung, Zbigniew Kalbarczyk, Ravishankar Iyer, and Adam J. Slagell. 2014. Preemptive intrusion detection. In *Proceedings of the 2014 Symposium and Bootcamp on the Science of Security (HotSoS '14)*. ACM, New York, NY, USA, , Article 21 , 2 pages. DOI=10.1145/2600176.2600197 <http://doi.acm.org/10.1145/2600176.2600197>
- [2] Phuong Cao, Eric Badger, Zbigniew Kalbarczyk, Ravishankar Iyer, and Adam Slagell. 2015. Preemptive intrusion detection: theoretical framework and real-world measurements. In *Proceedings of the 2015 Symposium and Bootcamp on the Science of Security (HotSoS '15)*. ACM, New York, NY, USA, , Article 5 , 12 pages. DOI=10.1145/2746194.2746199 <http://doi.acm.org/10.1145/2746194.2746199>

The Honey Pot and The Honey Badger

