

Financial Lower Bounds of Online Advertising Abuse

ELECTRICAL **[+]** COMPUTER

E N G I N E E R I N G



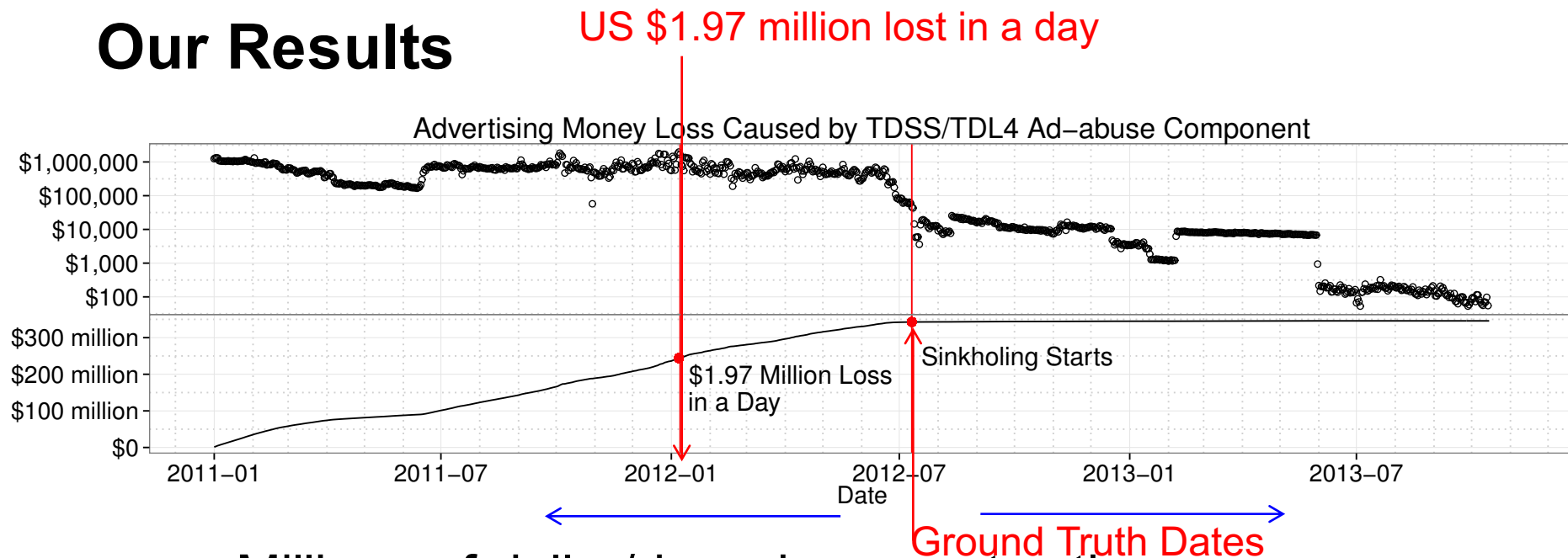
Yizheng Chen, Panagiotis Kintis,
Manos Antonakakis, Yacin Nadji, David Dagon,
Wenke Lee, and Michael Farrell



Monetization using Ads

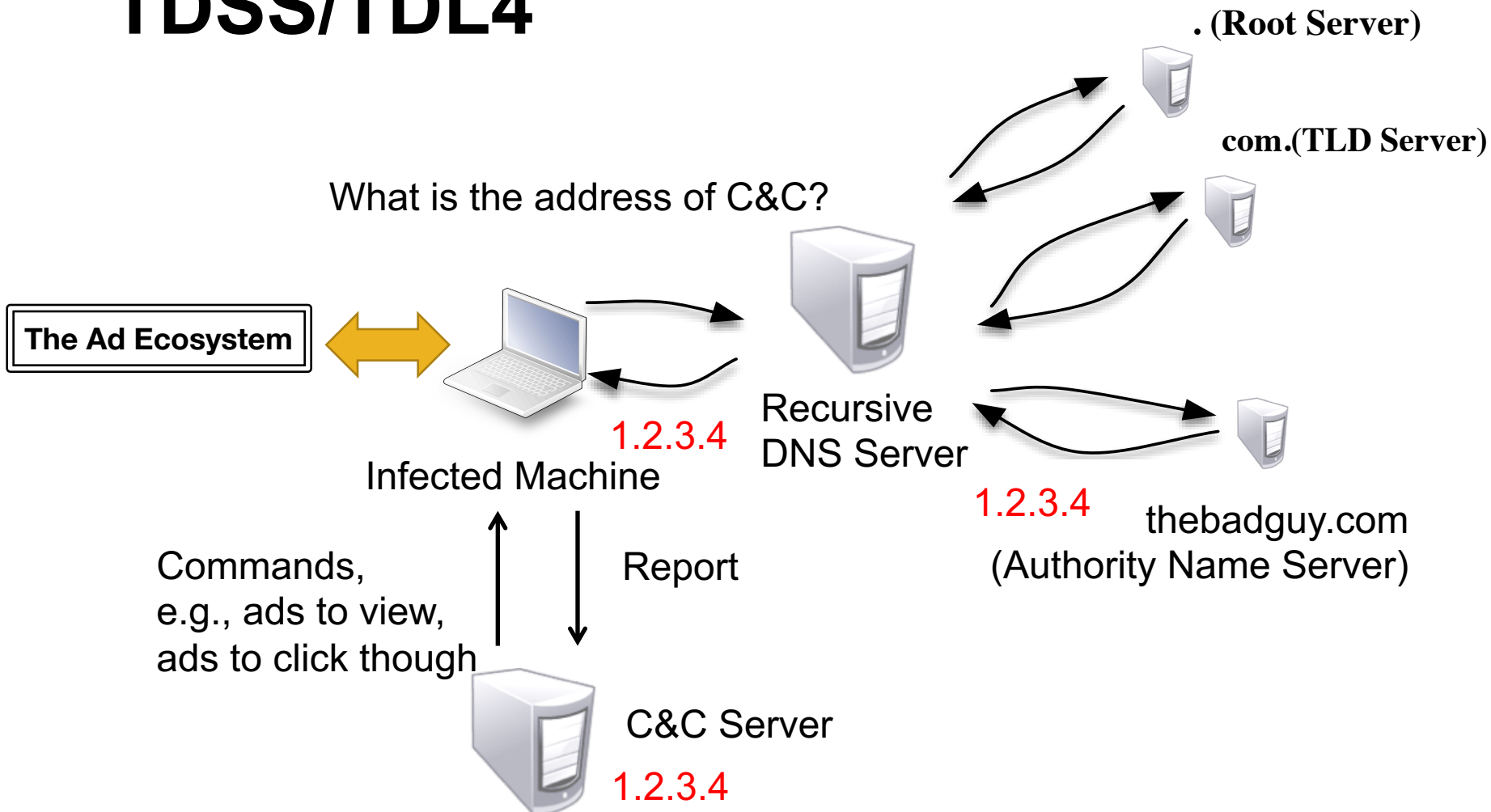
- Enable other malicious activities.
- Impression fraud.
 - Extremely hard to be detected.
- Methodology to track abuse in the long term.
 - We do not participate in the abuse.
- Financial lower bound.
 - Ad-abuse component of TDSS/TDL4.

Our Results

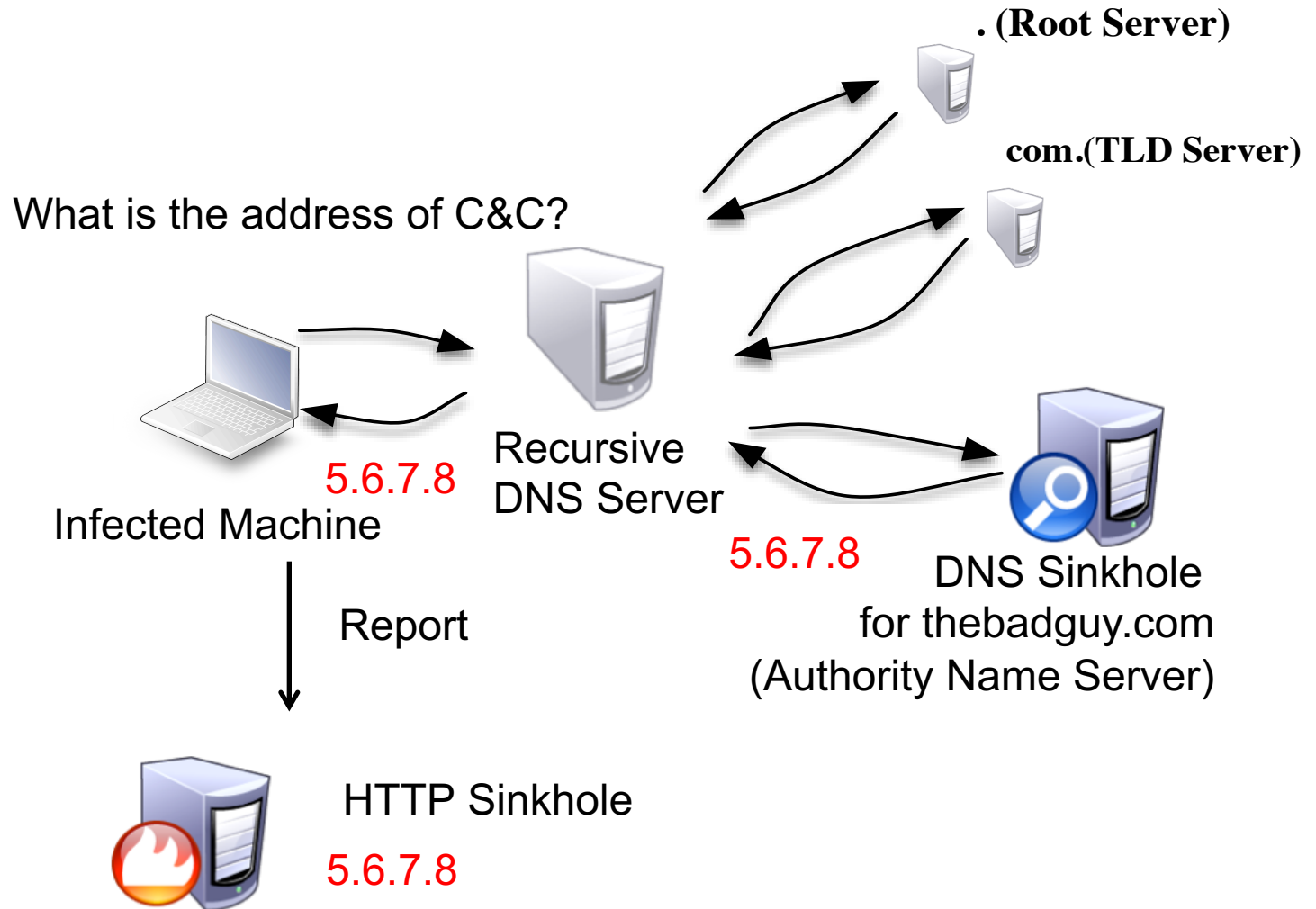


- Millions of dollar/day when most active
- 228 IP addresses and 863 resolved domain names
- \$24.22 million gain for botmasters

TDSS/TDL4

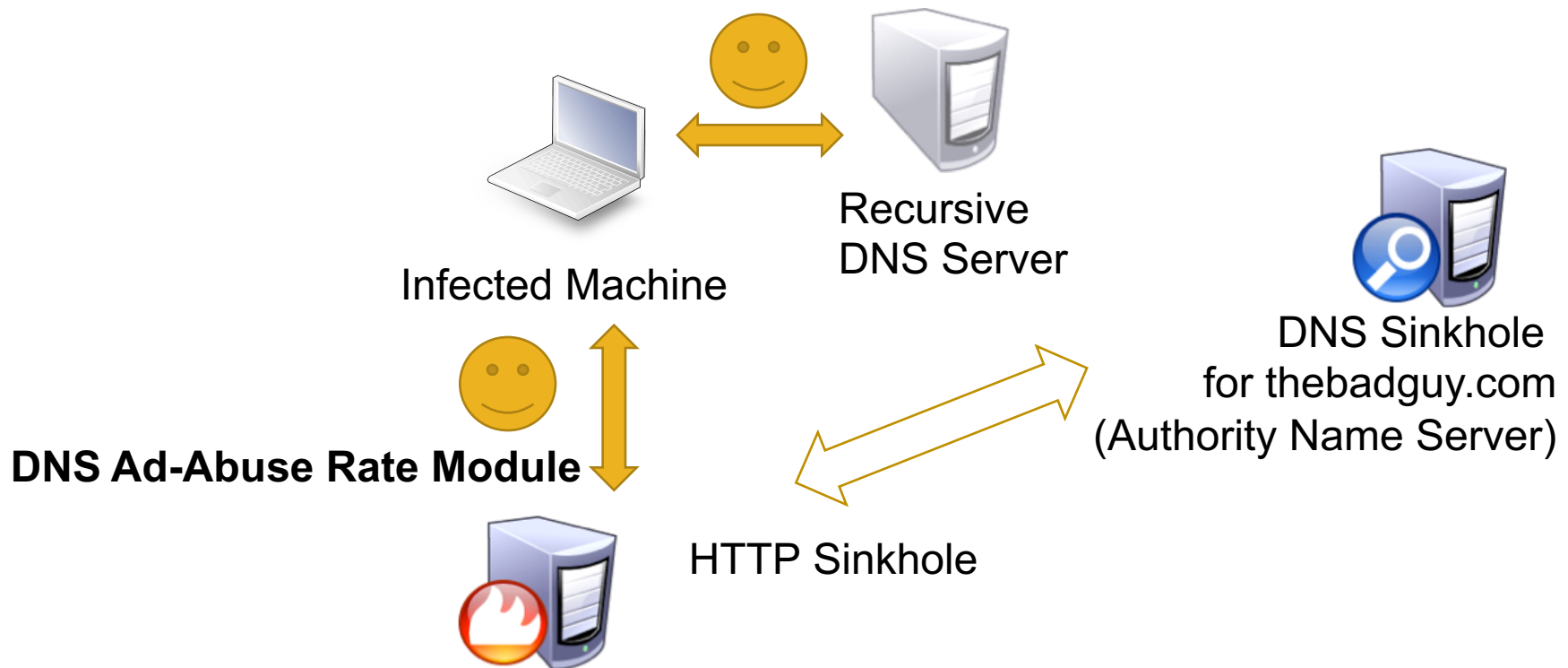


Ad-abuse Analysis System (A²S)

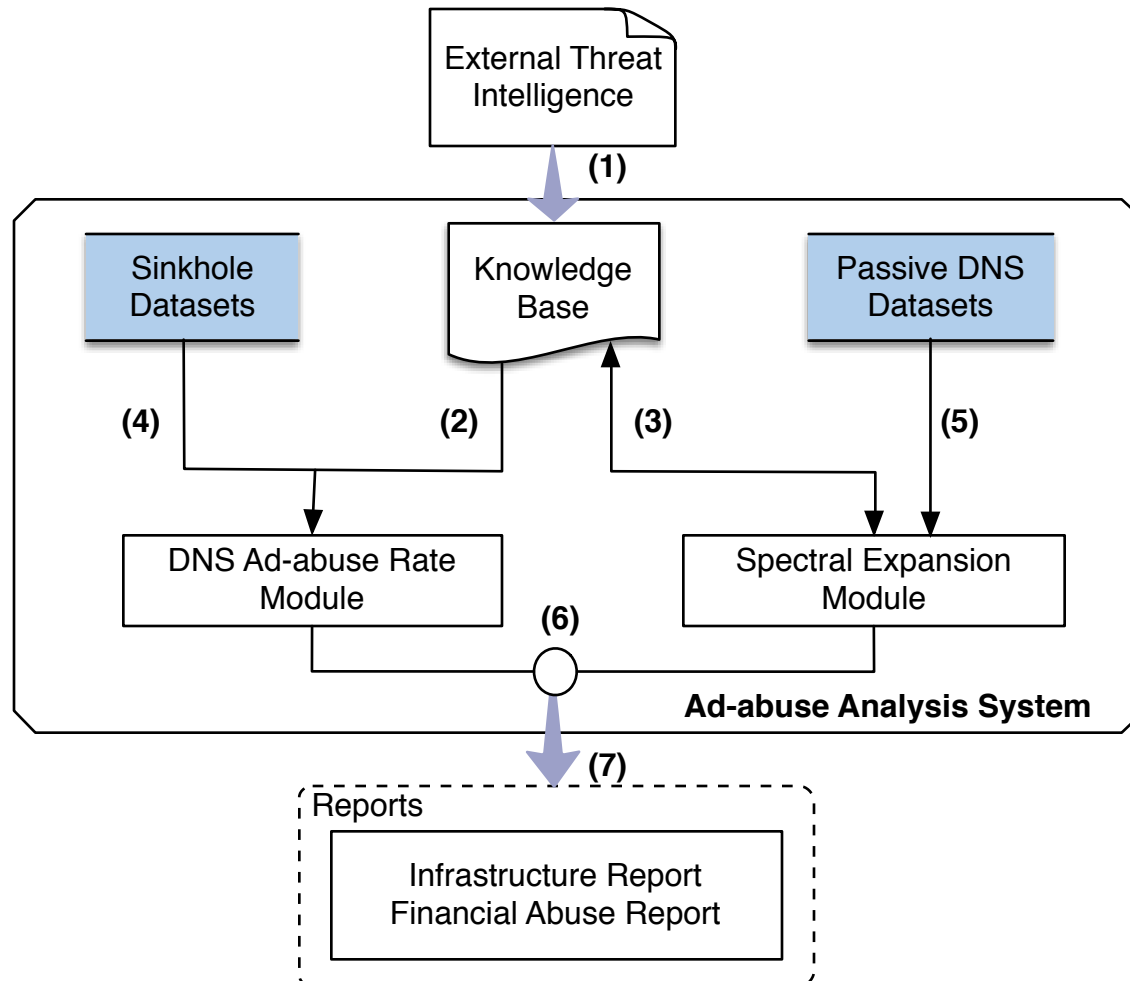


Ad-abuse Analysis System (A²S)

Spectral Expansion Module

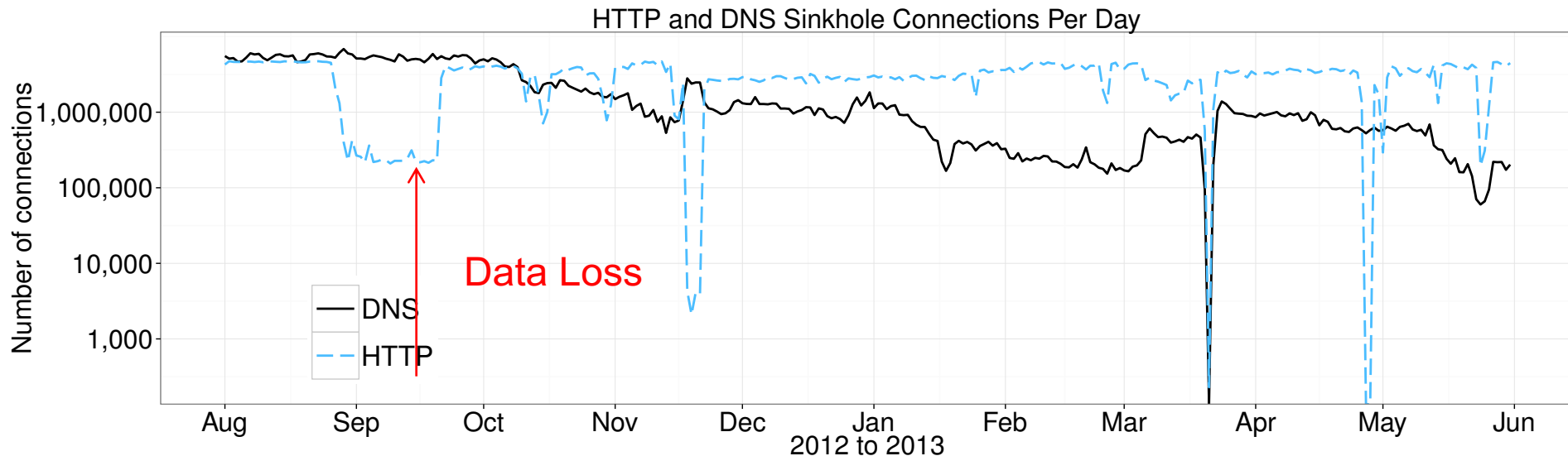


Ad-abuse Analysis System (A²S)

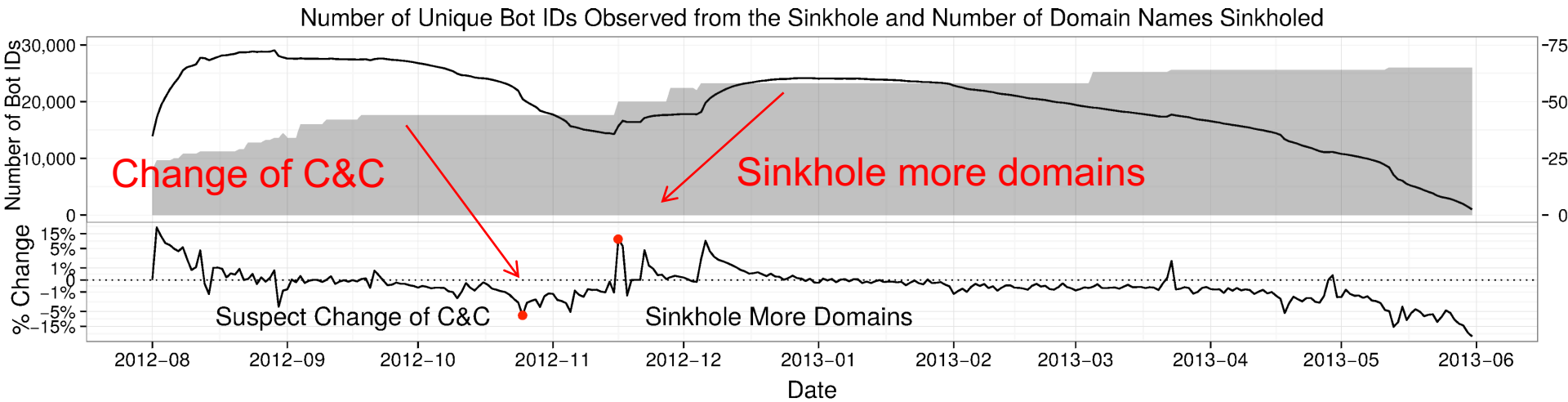


Datasets

	Date Range	Size	Records (millions)
DNS Sinkhole	8/1/2012 - 5/31/2013	6.9G	565
HTTP Sinkhole	8/1/2012 - 5/31/2013	248.6G	919
NXDOMAIN	6/27/2010 - 9/15/2014	133.5G	13,557
pDNS-DB	1/1/2011 - 11/6/2014	17.9T	10,209

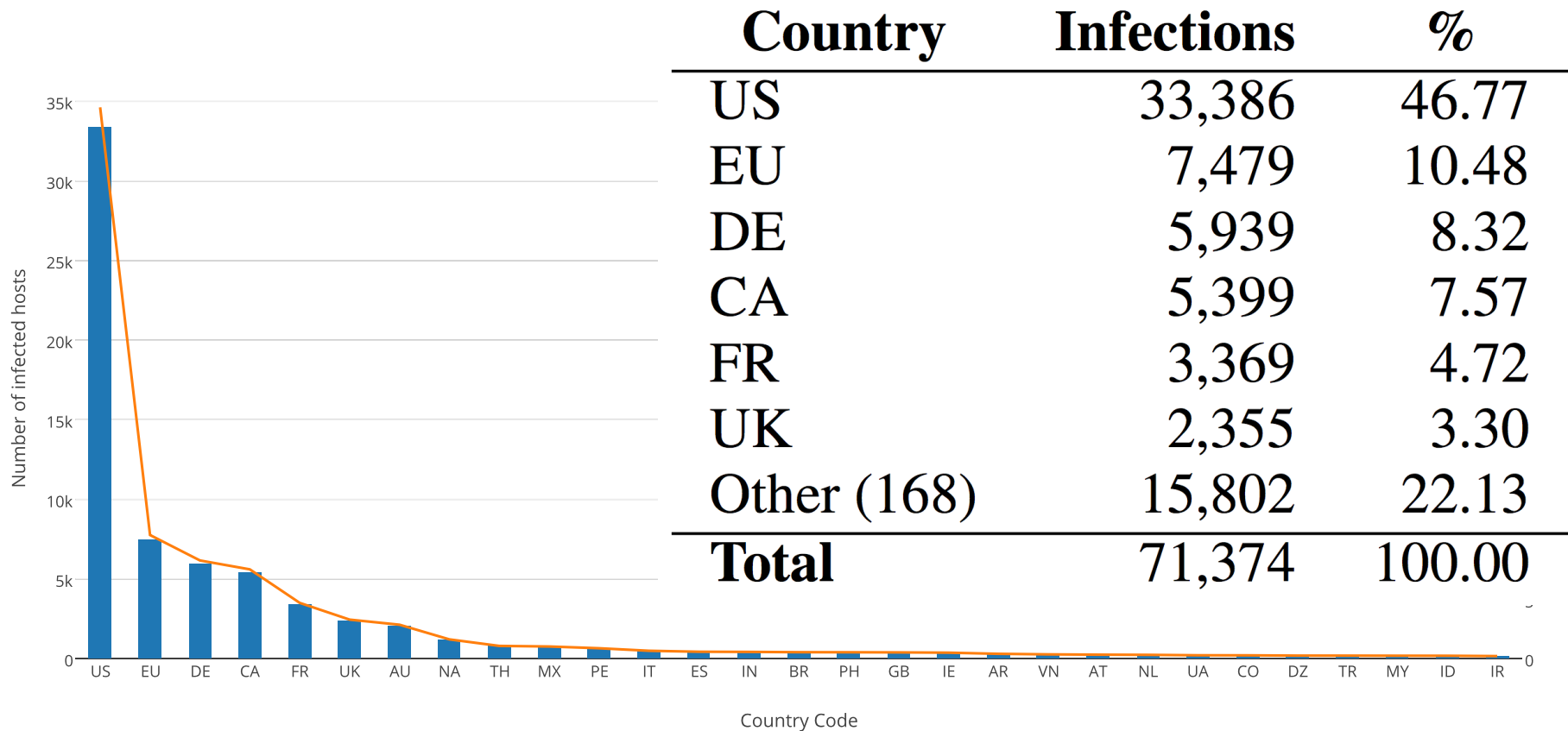


Sinkhole Datasets



- C&C domains change a lot.
- We need to track change.

Countries Affected



Observation: The measurement in ISP includes less than 15% of the botnet population.

DNS Ad-abuse Rate Module

- $\zeta = (\# \text{ of HTTP C\&C}) / (\# \text{ of DNS resolutions})$
- # of DNS requests to # of ad-abuse C&C connections
- “short-term” sinkhole observations to “long-term” passive DNS observations.

Spectral Expansion Module

- Unknown C&C domains can be
 - queried by known infected hosts
 - share infrastructure with known C&C domains.
- Iterative algorithm
 - Tripartite Graph: infected hosts - domains they queried - resolved data (A, CNAME).
 - Spectral clustering
 - Singular Value Decomposition + XMeans
 - Update Ad-abuse domains by cluster analysis
- Sanitize resulting domains
 - Financial Analysis: known email addresses or TDSS name servers.

Algorithm 1 Spectral Expansion Algorithm

Infected Hosts

- 1: $H \leftarrow \{h | \exists q \in D_A : h \text{ queried } q \text{ on day } d_i\}$
- 2: $D \leftarrow \{q | \exists h \in H : h \text{ queried } q \text{ on } d_i\}$
- 3: $Rdata \leftarrow \{ip | \exists q \in D : q \text{ resolved to } ip \text{ historically}\} \cup \{cname | \exists q \in D : q \text{ resolved to } cname \text{ historically}\}$
- 4: Apply thresholds α and β to the sets of $Rdata$ and H , respectively, to remove noisy IPs and hosts.
- 5: $M \leftarrow$ relationship between D and $(Rdata, H)$. Normalize by IPs, CNAMEs and Hosts.
- 6: $S \leftarrow M \times M^T$
- 7: $U\Sigma V^* \leftarrow SVD(S)$
- 8: $clusters \leftarrow XMeans(U)$
- 9: $D_A \leftarrow$ Analyze $clusters$.
- 10: $i = i + \delta$, Go to line 1.

Algorithm 1 Spectral Expansion Algorithm

Require: δ

Domains Queried by Known Infected Hosts

- 2: $D \leftarrow \{q | \exists h \in H : h \text{ queried } q \text{ on } d_i\}$
- 3: $Rdata \leftarrow \{ip | \exists q \in D : q \text{ resolved to } ip \text{ historically}\} \cup \{cname | \exists q \in D : q \text{ resolved to } cname \text{ historically}\}$
- 4: Apply thresholds α and β to the sets of $Rdata$ and H , respectively, to remove noisy IPs and hosts.
- 5: $M \leftarrow$ relationship between D and $(Rdata, H)$. Normalize by IPs, CNAMEs and Hosts.
- 6: $S \leftarrow M \times M^T$
- 7: $U\Sigma V^* \leftarrow SVD(S)$
- 8: $clusters \leftarrow XMeans(U)$
- 9: $D_A \leftarrow$ Analyze $clusters$.
- 10: $i = i + \delta$, Go to line 1.

Algorithm 1 Spectral Expansion Algorithm

Require: δ

1: $H \leftarrow \{h | \exists q \in D_A : h \text{ queried } q \text{ on day } d_i\}$

Resolved IP addresses and Canonical Names

3: $Rdata \leftarrow \{ip | \exists q \in D : q \text{ resolved to } ip \text{ historically}\} \cup \{cname | \exists q \in D : q \text{ resolved to } cname \text{ historically}\}$

4: Apply thresholds α and β to the sets of $Rdata$ and H , respectively, to remove noisy IPs and hosts.

5: $M \leftarrow$ relationship between D and $(Rdata, H)$. Normalize by IPs, CNAMEs and Hosts.

6: $S \leftarrow M \times M^T$

7: $U\Sigma V^* \leftarrow SVD(S)$

8: $clusters \leftarrow XMeans(U)$

9: $D_A \leftarrow$ Analyze $clusters$.

10: $i = i + \delta$, Go to line 1.

Algorithm 1 Spectral Expansion Algorithm

Require: δ

- 1: $H \leftarrow \{h | \exists q \in D_A : h \text{ queried } q \text{ on day } d_i\}$
 - 2: $D \leftarrow \{q | \exists h \in H : h \text{ queried } q \text{ on } d_i\}$
 - 3: $Rdata \leftarrow \{ip | \exists q \in D : q \text{ resolved to } ip \text{ historically}\} \cup$
Remove (more than) parking IPs, sinkholes, and large gateways.
 - 4: Apply thresholds α and β to the sets of $Rdata$ and H , respectively, to remove noisy IPs and hosts.
 - 5: $M \leftarrow$ relationship between D and $(Rdata, H)$. Normalize by IPs, CNAMEs and Hosts.
 - 6: $S \leftarrow M \times M^T$
 - 7: $U \Sigma V^* \leftarrow SVD(S)$
 - 8: $clusters \leftarrow XMeans(U)$
 - 9: $D_A \leftarrow$ Analyze $clusters$.
 - 10: $i = i + \delta$, Go to line 1.
-

Algorithm 1 Spectral Expansion Algorithm

Require: δ

- 1: $H \leftarrow \{h | \exists q \in D_A : h \text{ queried } q \text{ on day } d_i\}$
- 2: $D \leftarrow \{q | \exists h \in H : h \text{ queried } q \text{ on } d_i\}$
- 3: $Rdata \leftarrow \{ip | \exists q \in D : q \text{ resolved to } ip \text{ historically}\} \cup \{cname | \exists q \in D : q \text{ resolved to } cname \text{ historically}\}$
- 4: Apply thresholds α and β to the sets of $Rdata$ and H ,

Matrix normalization

- 5: $M \leftarrow$ relationship between D and $(Rdata, H)$. Normalize by IPs, CNAMEs and Hosts.
 - 6: $S \leftarrow M \times M^T$
 - 7: $U\Sigma V^* \leftarrow SVD(S)$
 - 8: $clusters \leftarrow XMeans(U)$
 - 9: $D_A \leftarrow$ Analyze $clusters$.
 - 10: $i = i + \delta$, Go to line 1.
-

Algorithm 1 Spectral Expansion Algorithm

Require: δ

- 1: $H \leftarrow \{h | \exists q \in D_A : h \text{ queried } q \text{ on day } d_i\}$
- 2: $D \leftarrow \{q | \exists h \in H : h \text{ queried } q \text{ on } d_i\}$
- 3: $Rdata \leftarrow \{ip | \exists q \in D : q \text{ resolved to } ip \text{ historically}\} \cup \{cname | \exists q \in D : q \text{ resolved to } cname \text{ historically}\}$
- 4: Apply thresholds α and β to the sets of $Rdata$ and H , respectively, to remove noisy IPs and hosts.
- 5: $M \leftarrow$ relationship between D and $(Rdata, H)$. Normalize

Association matrix \rightarrow similarity matrix

- 6: $S \leftarrow M \times M^T$
 - 7: $U\Sigma V^* \leftarrow SVD(S)$
 - 8: $clusters \leftarrow XMeans(U)$
 - 9: $D_A \leftarrow$ Analyze $clusters$.
 - 10: $i = i + \delta$, Go to line 1.
-

Algorithm 1 Spectral Expansion Algorithm

Require: δ

- 1: $H \leftarrow \{h | \exists q \in D_A : h \text{ queried } q \text{ on day } d_i\}$
- 2: $D \leftarrow \{q | \exists h \in H : h \text{ queried } q \text{ on } d_i\}$
- 3: $Rdata \leftarrow \{ip | \exists q \in D : q \text{ resolved to } ip \text{ historically}\} \cup \{cname | \exists q \in D : q \text{ resolved to } cname \text{ historically}\}$
- 4: Apply thresholds α and β to the sets of $Rdata$ and H , respectively, to remove noisy IPs and hosts.
- 5: $M \leftarrow$ relationship between D and $(Rdata, H)$. Normalize by IPs, CNAMEs and Hosts.

Singular Value Decomposition

- 7: $U\Sigma V^* \leftarrow SVD(S)$
- 8: $clusters \leftarrow XMeans(U)$
- 9: $D_A \leftarrow$ Analyze $clusters$.
- 10: $i = i + \delta$, Go to line 1.

Algorithm 1 Spectral Expansion Algorithm

Require: δ

- 1: $H \leftarrow \{h | \exists q \in D_A : h \text{ queried } q \text{ on day } d_i\}$
- 2: $D \leftarrow \{q | \exists h \in H : h \text{ queried } q \text{ on } d_i\}$
- 3: $Rdata \leftarrow \{ip | \exists q \in D : q \text{ resolved to } ip \text{ historically}\} \cup \{cname | \exists q \in D : q \text{ resolved to } cname \text{ historically}\}$
- 4: Apply thresholds α and β to the sets of $Rdata$ and H , respectively, to remove noisy IPs and hosts.
- 5: $M \leftarrow$ relationship between D and $(Rdata, H)$. Normalize by IPs, CNAMEs and Hosts.
- 6: $S \leftarrow M \times M^T$

Cluster the left Eigen vectors

- 8: $clusters \leftarrow XMeans(U)$
- 9: $D_A \leftarrow$ Analyze $clusters$.
- 10: $i = i + \delta$, Go to line 1.

Algorithm 1 Spectral Expansion Algorithm

Require: δ

- 1: $H \leftarrow \{h | \exists q \in D_A : h \text{ queried } q \text{ on day } d_i\}$
- 2: $D \leftarrow \{q | \exists h \in H : h \text{ queried } q \text{ on } d_i\}$
- 3: $Rdata \leftarrow \{ip | \exists q \in D : q \text{ resolved to } ip \text{ historically}\} \cup \{cname | \exists q \in D : q \text{ resolved to } cname \text{ historically}\}$
- 4: Apply thresholds α and β to the sets of $Rdata$ and H , respectively, to remove noisy IPs and hosts.
- 5: $M \leftarrow$ relationship between D and $(Rdata, H)$. Normalize by IPs, CNAMEs and Hosts.
- 6: $S \leftarrow M \times M^T$
- 7: $U\Sigma V^* \leftarrow SVD(S)$

Cluster analysis and label propagation

- 9: $D_A \leftarrow$ Analyze clusters.
- 10: $i = i + \delta$, Go to line 1.

Algorithm 1 Spectral Expansion Algorithm

Require: δ

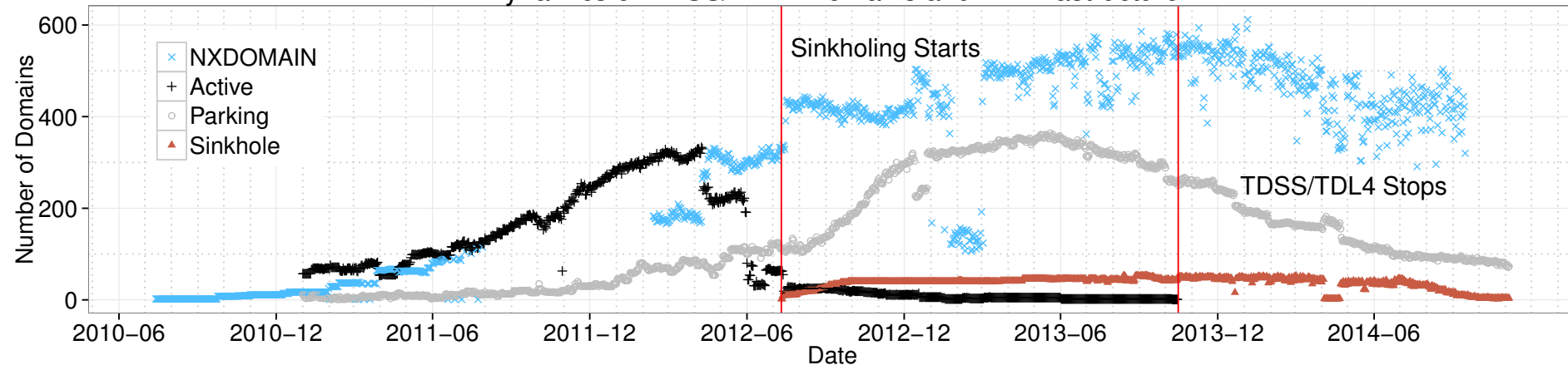
- 1: $H \leftarrow \{h | \exists q \in D_A : h \text{ queried } q \text{ on day } d_i\}$
- 2: $D \leftarrow \{q | \exists h \in H : h \text{ queried } q \text{ on } d_i\}$
- 3: $Rdata \leftarrow \{ip | \exists q \in D : q \text{ resolved to } ip \text{ historically}\} \cup \{cname | \exists q \in D : q \text{ resolved to } cname \text{ historically}\}$
- 4: Apply thresholds α and β to the sets of $Rdata$ and H , respectively, to remove noisy IPs and hosts.
- 5: $M \leftarrow$ relationship between D and $(Rdata, H)$. Normalize by IPs, CNAMEs and Hosts.
- 6: $S \leftarrow M \times M^T$
- 7: $U\Sigma V^* \leftarrow SVD(S)$
- 8: $clusters \leftarrow XMeans(U)$

Next iteration

- 10: $i = i + \delta$, Go to line 1.

Infrastructure Report

Dynamics of TDSS/TDL4 Domains and IP Infrastructure



- 3 False Positives out of 1,134
- <http://tinyurl.com/dimva16-tdss-domains>

Observation: The number of active domain names daily increased from 2010, and reached the maximum (333) on 4/9/2012.

Observation: None of the domains resolved to any active IP after 10/15/2013.

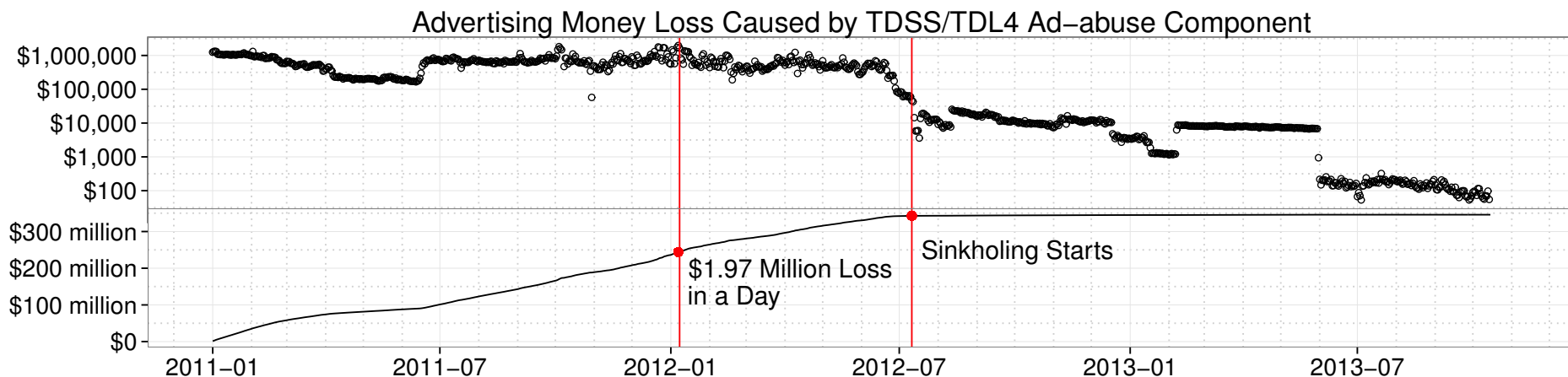
Financial Model

$$M_{impression} = \sum_i \zeta * R_i * \underbrace{\left(p_{im} * \frac{\mu_{im}}{1000} * CPM \right)}_{\text{Impression Fraud Component}}$$

Number of DNS Requests

- ζ : DNS Ad-abuse Rate
- p_{im} : impression percentage
- μ_{im} : number of impressions per C&C HTTP connection
- CPM: Cost Per Millie (Cost per thousand impressions)

Financial Report



Note: \$14 million were found in DNSChanger operators' bank accounts.

Stakeholders		Money (millions)
Advertisers' Capital		346.00
DSP	45%	155.70
Ad Exchange (inbound)	8%	27.68
Ad Exchange (outbound)	8%	27.68
Ad Networks	32%	110.72
Ad Server/Publisher (Affiliates)	7%	24.22

Questions?