

# DNS Monitoring, looking out for anomalies using the time frame of Name - IP association

*Lautaro Dolberg*

- 1 Introduction
- 2 Background
- 3 MAM
- 4 Analysis
- 5 Evaluation
- 6 Conclusions

- 1 Introduction
- 2 Background
- 3 MAM
- 4 Analysis
- 5 Evaluation
- 6 Conclusions

# Personal Information

Lautaro Dolberg, Born in Buenos Aires - Argentina. 29  
Years old

- ▶ PhD Studies at Universite de Luxembourg.
- ▶ Started **15/1/2012** - Expected Defense **2X/3/2015**
- ▶ Secan Lab - Vehicular Lab. J Francois, R State, R. Frank & T. Engel

Lautaro Dolberg, Born in Buenos Aires - Argentina. 29  
Years old

- ▶ PhD Studies at Universite de Luxembourg.
- ▶ Started **15/1/2012** - Expected Defense **2X/3/2015**
- ▶ Secan Lab - Vehicular Lab. J Francois, R State, R. Frank & T. Engel
- ▶ APR 14 - Prof Raouf Boutaba research team at Waterloo University, Ontario, Canada.

Lautaro Dolberg, Born in Buenos Aires - Argentina. 29 Years old

- ▶ PhD Studies at Universite de Luxembourg.
- ▶ Started **15/1/2012** - Expected Defense **2X/3/2015**
- ▶ Secan Lab - Vehicular Lab. J Francois, R State, R. Frank & T. Engel
- ▶ APR 14 - Prof Raouf Boutaba research team at Waterloo University, Ontario, Canada.
- ▶ "Licenciatura en Ciencias de la Computacion" at DC - FCEN UBA (2011)

Lautaro Dolberg, Born in Buenos Aires - Argentina. 29 Years old

- ▶ PhD Studies at Universite de Luxembourg.
- ▶ Started **15/1/2012** - Expected Defense **2X/3/2015**
- ▶ Secan Lab - Vehicular Lab. J Francois, R State, R. Frank & T. Engel
- ▶ APR 14 - Prof Raouf Boutaba research team at Waterloo University, Ontario, Canada.
- ▶ "Licenciatura en Ciencias de la Computacion" at DC - FCEN UBA (2011)
- ▶ Research Keywords: Security, Networks, VANET, Monitoring, Big Data, Mobile Devices, SDN

- ▶ Unrestricted access networks
  - ▶ Internet Applications Mobile (Crowd-sourcing)
  - ▶ Traffic Sensing Applications (LuxTraffic)
  - ▶ Internet Name Resolution (DNS)
- ▶ Large Data Collections
- ▶ Monitoring for detecting misbehavior
- ▶ Mobile Security (Android + IOS)
- ▶ Software Defined Networks
- ▶ Network Awareness

*Security & Management of distributed applications*



- 1 Introduction
- 2 Background**
- 3 MAM
- 4 Analysis
- 5 Evaluation
- 6 Conclusions

*DNS traffic reflects Internet activities and behaviors*

- ▶ Internet Threats Growing: Phishing, Malware, Spoofed Domains.
- ▶ Identify malware behavior by assessing association time between names and networks.
- ▶ Helps for contextualizing other data such as Netflow, etc.
- ▶ Available resources such as Passive DNS.
- ▶ As both DNS and IP space follow hierarchical organization, MAM can be used.

- ▶ *Can we use DNS records and its changes over time to trace Internet activities?*

- ▶ *Do malicious domains behave different from others in terms of name - ip association?*

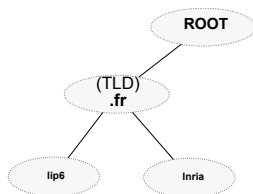
- ▶ *Do malicious domains behave different from others in terms of name - ip association?*

## Example DNS Records

Type	Name	IPV4	TTL
CNAME	www.lip6.fr	ww.lip6.fr	X
A	ww.lip6.fr	132.227.104.15	X

*DNS is an essential service for Internet*

- ▶ Emerged in 1987 (RFC 1035, 1123, 2181)
- ▶ Domain names are labels separated with dots.
- ▶ Strictly hierarchical pattern
- ▶ TLD (eg: .com, .fr, .etc)

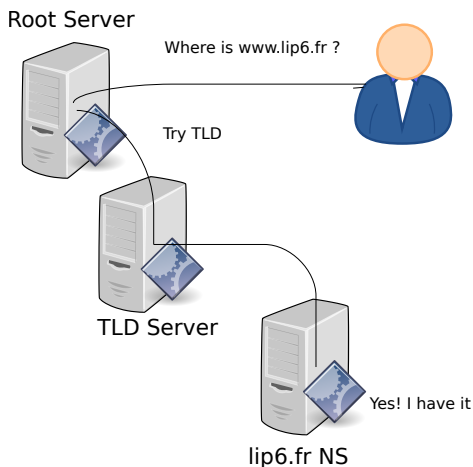


Max depth 127. Limited to 253 characters, label limit 63 characters.

# DNS Structure & Procedure

Resolving `www.lip6.fr`

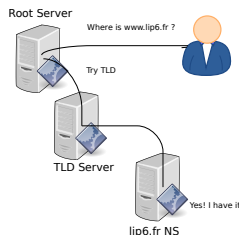
# DNS Structure & Procedure





# DNS Structure & Procedure

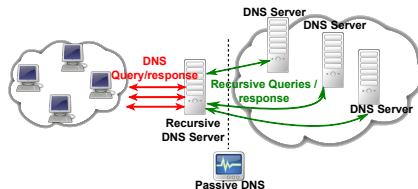
## Resolving www.lip6.fr



## Relevant DNS Registers

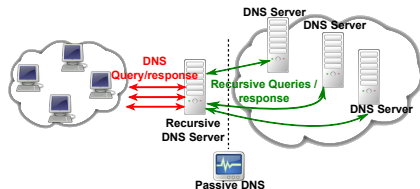
- ▶ A: Association IPV4 or/and IPV6.
- ▶ CNAME: Redirect.
- ▶ PTR: Reverse DNS search.

*A Passive DNS DB contains:*



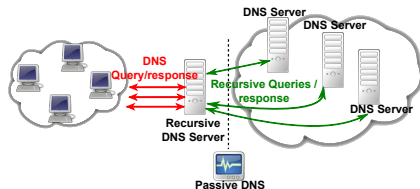
- ▶ Where did this domain name point to in the past?

*A Passive DNS DB contains:*



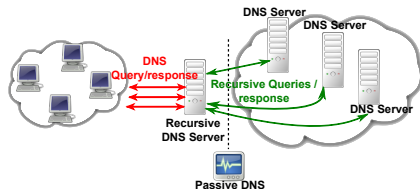
- ▶ Where did this domain name point to in the past?
- ▶ What domain names are hosted by a given nameserver?

*A Passive DNS DB contains:*



- ▶ Where did this domain name point to in the past?
- ▶ What domain names are hosted by a given nameserver?
- ▶ What domain names point into a given IP network?

*A Passive DNS DB contains:*



- ▶ Where did this domain name point to in the past?
- ▶ What domain names are hosted by a given nameserver?
- ▶ What domain names point into a given IP network?
- ▶ What subdomains exist below a certain domain name?

## Previous Work

- ▶ Exposure: Finding malicious domains using passive dns analysis, in *Network and Distributed System Security Symposium - NDSS, 2011*.
- ▶ DNSSM: A large-scale Passive DNS Security Monitoring Framework, in *IEEE/IFIP Network Operations and Management Symposium, 2012*.
- ▶ SDBF: Smart DNS Brute-Forcer, in *IEEE/IFIP Network Operations and Management Symposium - NOMS, 2012*.

- 1 Introduction
- 2 Background
- 3 MAM**
- 4 Analysis
- 5 Evaluation
- 6 Conclusions

*How do we organize all the DNS data that we have?*



*How do we organize all the DNS data that we have?*

- ▶ Preserve the hierarchy of Data (DNS & IP).
- ▶ Reduce the scale

*How do we organize all the DNS data that we have?*

- ▶ Preserve the hierarchy of Data (DNS & IP).
- ▶ Reduce the scale
- ▶ Minimize information loss due to aggregation.

*How do we organize all the DNS data that we have?*

- ▶ Preserve the hierarchy of Data (DNS & IP).
- ▶ Reduce the scale
- ▶ Minimize information loss due to aggregation.
- ▶ Fine control granularity for data analysis.

MAM is an enabler for aggregation and data retrieval.

# SNT Multidimensional Aggregation Monitoring + Applications

## *Aggregation*

- ▶ **Scalable** way to represent information
  - ▶ **Outline** relevant correlated facts
  - ▶ Flexible granularity

# SNT Multidimensional Aggregation Monitoring

## + Applications

### *Aggregation*

- ▶ **Scalable** way to represent information
  - ▶ **Outline** relevant correlated facts
  - ▶ Flexible granularity
- ▶ **Features:**
  - ▶ Custom Units (e.g. traffic packets, vehicle, traffic units)
  - ▶ choose **criteria** for aggregation

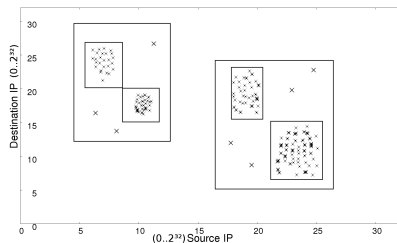
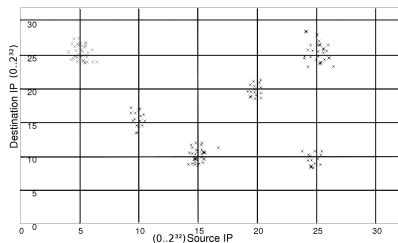
# SNT Multidimensional Aggregation Monitoring

## + Applications

### Aggregation

- ▶ **Scalable** way to represent information
  - ▶ **Outline** relevant correlated facts
  - ▶ Flexible granularity
- ▶ Features:
  - ▶ Custom Units (e.g. traffic packets, vehicle, traffic units)
  - ▶ choose **criteria** for aggregation
- ▶ **Temporal and Spatial aggregation**
  - ▶ Temporal: time windows split ( $\beta$ )
  - ▶ Spatial: keep nodes with activity  $> \alpha$  e.g. *traffic volume*, aggregate the others into their parents  $\rightarrow$  needs **hierarchical relationships**

*Example: IPv4 space partitioning, static vs dynamic*



With MAM is possible to generate time aggregated combining multiple data.

- ▶ Two dimensions

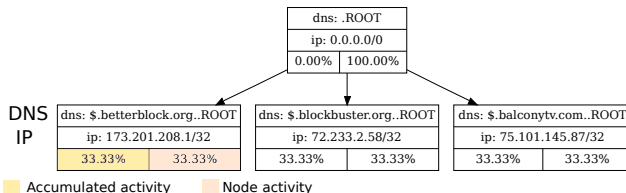


With MAM is possible to generate time aggregated combining multiple data.

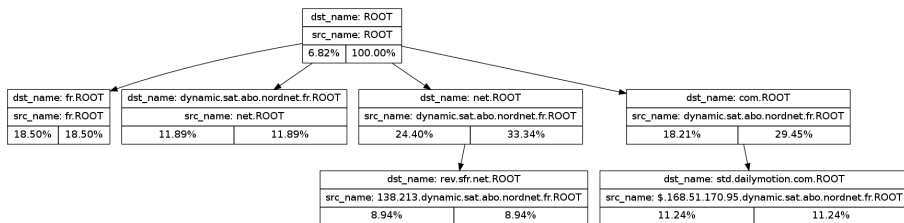
- ▶ Two dimensions
- ▶ Hierarchically derived from data model (IPV4 & DNS Data Space)

With MAM is possible to generate time aggregated combining multiple data.

- ▶ Two dimensions
- ▶ Hierarchically derived from data model (IPV4 & DNS Data Space)
- ▶ Example



## *Multidimensional tree for source and destination names*



*So far...*

- ▶ DNS Interest and Background

*So far...*

- ▶ DNS Interest and Background
- ▶ Multidimensional Aggregation for hierarchical data

*So far...*

- ▶ DNS Interest and Background
- ▶ Multidimensional Aggregation for hierarchical data
- ▶ Coming next ← Leverage DNS activity analysis with MAM

- 1 Introduction
- 2 Background
- 3 MAM
- 4 Analysis**
- 5 Evaluation
- 6 Conclusions

*Asses the duration of the association (IP-NAME) over time*



*Asses the duration of the association (IP-NAME) over time*

We want to keep the notion of subnet and subdomain.

This should be stable over time



www.lip6.fr

132.227.104.15/32

*Asses the duration of the association (IP-NAME) over time*  
We want to keep the notion of subnet and subdomain.

lets asume bmp.fr is fishing for bnp



www.bmp.fr

192.168.1.2/32

Assuming a sequence of K Trees

- ▶  $S = \{T_1 \dots T_K\}$  representing DNS record association over time split in K

Assuming a sequence of  $K$  Trees

- ▶  $S = \{T_1 \dots T_K\}$  representing DNS record association over time split in  $K$
- ▶  $n \in \text{Nodes}(T_i)$

Assuming a sequence of  $K$  Trees

- ▶  $S = \{T_1 \dots T_K\}$  representing DNS record association over time split in  $K$
- ▶  $n \in \text{Nodes}(T_i)$ 
  - ▶  $n.dns$  represents the DNS name (using 'dollar' as terminal symbol)

Assuming a sequence of K Trees

- ▶  $S = \{T_1 \dots T_K\}$  representing DNS record association over time split in K
- ▶  $n \in \text{Nodes}(T_i)$ 
  - ▶  $n.dns$  represents the DNS name (using 'dollar' as terminal symbol)
  - ▶  $n.ip$  represents the IPv4 address as a tuple (address, prefix\_length)

Assuming a sequence of K Trees

- ▶  $S = \{T_1 \dots T_K\}$  representing DNS record association over time split in K
- ▶  $n \in \text{Nodes}(T_i)$ 
  - ▶  $n.dns$  represents the DNS name (using 'dollar' as terminal symbol)
  - ▶  $n.ip$  represents the IPv4 address as a tuple (address, prefix\_length)
  - ▶  $n.accum$  is the accumulated value representing the number of A Records.

Assuming a sequence of  $K$  Trees

- ▶  $S = \{T_1 \dots T_K\}$  representing DNS record association over time split in  $K$
- ▶  $n \in \text{Nodes}(T_i)$ 
  - ▶  $n.dns$  represents the DNS name (using 'dollar' as terminal symbol)
  - ▶  $n.ip$  represents the IPv4 address as a tuple (address, prefix\_length)
  - ▶  $n.accum$  is the accumulated value representing the number of A Records.
- ▶ The trees from  $S$  are aggregated according to a given *alpha*



# Steadiness Metric I (Similarity)

The similarity of nodes is positive if:

- ▶  $n_1$ ,  $n_2$  Tree nodes are from the same domain (Total overlap)

# Steadiness Metric I (Similarity)

The similarity of nodes is positive if:

- ▶  $n1, n2$  Tree nodes are from the same domain (Total overlap)
- ▶  $n1.dns \subset n2.dns$  OR  $n1.ip \subset n2.ip$  (Partial overlap)

# Steadiness Metric I (Similarity)

The similarity of nodes is positive if:

- ▶  $n1, n2$  Tree nodes are from the same domain (Total overlap)
- ▶  $n1.dns \subset n2.dns$  OR  $n1.ip \subset n2.ip$  (Partial overlap)
- ▶  $n2.dns \subset n1.dns$  OR  $n2.ip \subset n1.ip$  (Partial overlap)

Otherwise similarity is 0.

## Steadiness Metric I (Similarity)

The similarity of nodes is positive if:

- ▶  $n1, n2$  Tree nodes are from the same domain (Total overlap)
- ▶  $n1.dns \subset n2.dns$  OR  $n1.ip \subset n2.ip$  (Partial overlap)
- ▶  $n2.dns \subset n1.dns$  OR  $n2.ip \subset n1.ip$  (Partial overlap)

Otherwise similarity is 0.

Given  $n$  from  $T$ , we look for  $m \in Nodes(T_x)$  where

$$\forall m_x \in Nodes(T_x) : sim(n, m) \geq sim(n, m_x)$$

# Steadiness Metric II (Similarity)

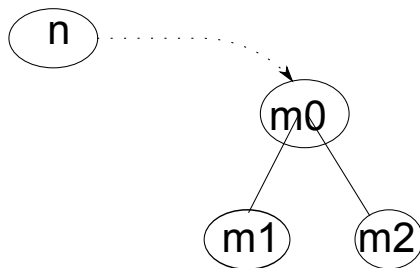
In other words:

- ▶ It's almost as inserting the node into the tree, and looking for a possible parent

# Steadiness Metric II (Similarity)

In other words:

- ▶ It's almost as inserting the node into the tree, and looking for a possible parent



# Steadiness Metrics III

*Tree comparison, how to establish a similarity criteria?*

## Steadiness Metrics III

*Tree comparison, how to establish a similarity criteria?*

$$\text{sim}(n1, n2) = \alpha \times \text{IP\_sim}(n1, n2) + \beta \times \text{DNS\_sim}(n1, n2) + \gamma \times \text{vol\_sim}(n1, n2)$$



*Tree comparison, how to establish a similarity criteria?*

$$\text{sim}(n1, n2) = \alpha \times \text{IP\_sim}(n1, n2) + \beta \times \text{DNS\_sim}(n1, n2) + \gamma \times \text{vol\_sim}(n1, n2)$$

$$\text{IP\_sim}(n1, n2) = 1 - \frac{|n1_{\text{prefix\_len}} - n2_{\text{prefix\_len}}|}{32}$$

*Tree comparison, how to establish a similarity criteria?*

$$\text{sim}(n1, n2) = \alpha \times \text{IP\_sim}(n1, n2) + \beta \times \text{DNS\_sim}(n1, n2) + \gamma \times \text{vol\_sim}(n1, n2)$$

$$\text{IP\_sim}(n1, n2) = 1 - \frac{|n1_{\text{prefix\_len}} - n2_{\text{prefix\_len}}|}{32}$$

$$\text{DNS\_sim}(n1, n2) = \frac{|n1_{\text{dns}} \cap n2_{\text{dns}}|}{|n1_{\text{dns}} \cup n2_{\text{dns}}|}$$

## Steadiness Metrics III

*Tree comparison, how to establish a similarity criteria?*

$$sim(n1, n2) = \alpha \times IP\_sim(n1, n2) + \beta \times$$

$$DNS\_sim(n1, n2) + \gamma \times vol\_sim(n1, n2)$$

$$IP\_sim(n1, n2) = 1 - \frac{|n1_{prefix\_len} - n2_{prefix\_len}|}{32}$$

$$DNS\_sim(n1, n2) = \frac{|n1_{dns} \cap n2_{dns}|}{|n1_{dns} \cup n2_{dns}|}$$

$$vol\_sim(n1, n2) = 1 - 0.01 \times |n1_{acc\_vol} - n2_{acc\_vol}|$$

*Smoothing helps considering early time windows*

- ▶  $n1 \in \text{Nodes}(T_1)$ ,  $n2 \in \text{Nodes}(T_2)$
- ▶ We compute  $n2 = \text{most\_sim}(n1)$  is
- ▶  $\text{stead}(n1) = \text{sim}(n1, n2) + \mu \times \text{stead}(n2)$ . With  $T_0$  as base case and  $\mu \in \mathbb{R}$ .

*Smoothing helps considering early time windows*

- ▶  $n1 \in \text{Nodes}(T_1)$ ,  $n2 \in \text{Nodes}(T_2)$
- ▶ We compute  $n2 = \text{most\_sim}(n1)$  is
- ▶  $\text{stead}(n1) = \text{sim}(n1, n2) + \mu \times \text{stead}(n2)$ . With  $T_0$  as base case and  $\mu \in \mathbb{R}$ .

So we compute the global steadiness of a Tree  $T$  by:

$$\text{Persistence}(T) = \frac{\sum_{n \in \text{Nodes}(T)} \text{stead}(n)}{|\text{Nodes}(T)|}$$

- 1 Introduction
- 2 Background
- 3 MAM
- 4 Analysis
- 5 Evaluation**
- 6 Conclusions

## *Aggregation Window: 1 Week Time Length*

- ▶ Macro: Up to 52 weeks from 2011-04-23 to 2012-06-30 (662 K)
- ▶ Micro: 10 weeks maximum

### *Aggregation Window: 1 Week Time Length*

- ▶ Macro: Up to 52 weeks from 2011-04-23 to 2012-06-30 (662 K)
- ▶ Micro: 10 weeks maximum

### *Malicious data*

- ▶ Time: Periodically, Steady



### *Aggregation Window: 1 Week Time Length*

- ▶ Macro: Up to 52 weeks from 2011-04-23 to 2012-06-30 (662 K)
- ▶ Micro: 10 weeks maximum

### *Malicious data*

- ▶ Time: Periodically, Steady
- ▶ Proportion: 0.1%, 1% and 10%

### *Aggregation Window: 1 Week Time Length*

- ▶ Macro: Up to 52 weeks from 2011-04-23 to 2012-06-30 (662 K)
- ▶ Micro: 10 weeks maximum

### *Malicious data*

- ▶ Time: Periodically, Steady
- ▶ Proportion: 0.1%, 1% and 10%
- ▶ Source: Blacklists (Exposure, WOT) 175K

### *Aggregation Granularity: 2%*

### *Distribution of Local Steadiness (Leafs)*

- ▶ More than 50% of malicious nodes have less than 0.7 of stability

# Steadiness Distribution

## *Distribution of Local Steadiness (Leafs)*

- ▶ More than 50% of malicious nodes have less than 0.7 of stability
- ▶ Less than 20% of malicious nodes have more than 0.85 of stability

### *Distribution of Local Steadiness (Leafs)*

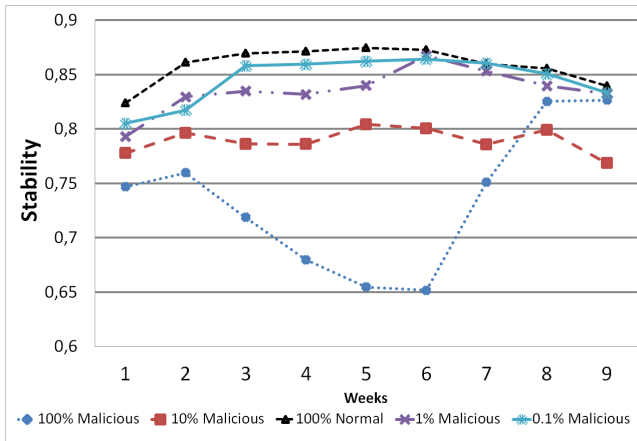
- ▶ More than 50% of malicious nodes have less than 0.7 of stability
- ▶ Less than 20% of malicious nodes have more than 0.85 of stability
- ▶ Less than 40% of normal data have a steadiness of 0.8% or less.

### *Distribution of Local Steadiness (Leafs)*

- ▶ More than 50% of malicious nodes have less than 0.7 of stability
- ▶ Less than 20% of malicious nodes have more than 0.85 of stability
- ▶ Less than 40% of normal data have a steadiness of 0.8% or less.
- ▶ Only 10% of normal data have a steadiness of less than 0.5

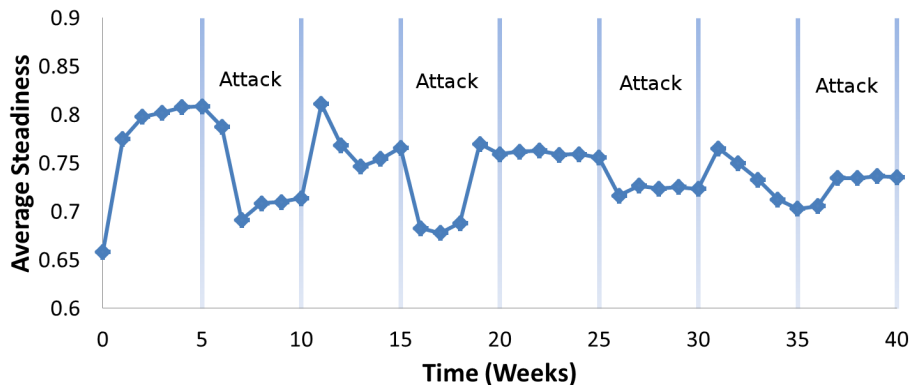
# Microscopic observation

*Malicious domains causes a drop on average steadiness*



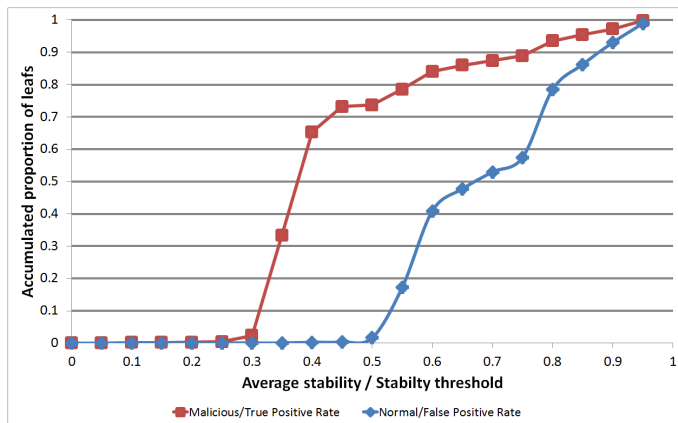
# Macroscopic observation

*Malicious domains causes a drop on average steadiness:  
Macro*





*Accuracy: Steadiness as metric for filtering malicious domains*



- 1 Introduction
- 2 Background
- 3 MAM
- 4 Analysis
- 5 Evaluation
- 6 Conclusions**

- ▶ A methodology for assessing DNS - IP association time frame was proposed.

- ▶ A methodology for assessing DNS - IP association time frame was proposed.
- ▶ Reduced the scale of data, helpful in the context of network security. (From 80K nodes to 2K with  $\alpha = 2\%$ )

- ▶ A methodology for assessing DNS - IP association time frame was proposed.
- ▶ Reduced the scale of data, helpful in the context of network security. (From 80K nodes to 2K with  $\alpha = 2\%$ )
- ▶ Definition of steadiness metrics for a local and global scope was introduced.

- ▶ A methodology for assessing DNS - IP association time frame was proposed.
- ▶ Reduced the scale of data, helpful in the context of network security. (From 80K nodes to 2K with  $\alpha = 2\%$ )
- ▶ Definition of steadiness metrics for a local and global scope was introduced.
- ▶ Evaluation using real data and during several time frames. Validation of the metrics

- ▶ A methodology for assessing DNS - IP association time frame was proposed.
- ▶ Reduced the scale of data, helpful in the context of network security. (From 80K nodes to 2K with  $\alpha = 2\%$ )
- ▶ Definition of steadiness metrics for a local and global scope was introduced.
- ▶ Evaluation using real data and during several time frames. Validation of the metrics
- ▶ Scalability: It can be implemented using dynamic programming / distributed computing.

### *Relevant Publications*

- TCP** 2012: L. Dolberg, J. Francois, T. Engel. "Multidimensional Aggregation Monitoring", Usenix LISA 2012
- DNS** 2013: L. Dolberg, J. Francois, T. Engel. "DNS Malware Detection using Stability Metrics", IEEE LCN 2013



*Thanks!*

