

# **Modern Botnets**

#### and the Rise of Automatically Generated Domains

Joint work with

Stefano Schiavoni (POLIMI & Google, MSc),

Edoardo Colombo (POLIMI)

Lorenzo Cavallaro (RHUL, PhD),

Stefano Zanero (POLIMI, PhD)



Federico Maggi federico@maggi.cc Politecnico di Milano

### Who I am

Federico Maggi, PhD

**Post-doctoral Researcher** 



POLITECNICO DI MILANO



#### Topics

Android malware, malware analysis, web measurements

#### Background

Intrusion detection, anomaly detection



#### www.red-book.eu



#### The RED BOOK

A Roadmap for Systems Security Research

#### Audience

Policy makers Researchers Journalists

#### Content

Vulnerabilities Social Networks Critical Infrastructure Mobile Devices Malware







### Roadmap

- 1. Botnets
- 2. Communication channels
- 3. Domain generation algorithms (DGAs)
- 4. Detecting DGA-based botnets
- 5. Results



### Roadmap

#### 1. Botnets

- 2. Communication channels
- 3. Domain generation algorithms (DGAs)
- 4. Detecting DGA-based botnets
- 5. Results



### **Botnets: from malware to service**

#### Botnet

- Network of (malware infected) computers
- Controlled by an external entity (e.g., cybercriminal)

#### Bot

- Computer member of a botnet
- Infected with malicious software

#### **Botmaster**

Person or group managing the botnet



### **Centralized topology example**





# Infected machines = \$\$\$

#### **Steal sensitive information**

- harvest contacts
- online banking credentials

#### **Run malicious activities**

- send spam, phishing emails, click fraud
- denial of service

#### Make money

• rent the infrastructure as a service

#### Maintenance

• update the malware



# **Command & control flow**





# Administration dashboard (spyeye)

ccess, place your bookmarks here on the bookmarks bar. Import bookmarks now				
	Bots	Full Statistic 😽 Create Task 🕕	Tasks	
2011	VIRTEST	Plugins backconnect	SOCKS 5 0 642	
06:07:52			5561	
			609 C+	
	0	Loge Diler Blar		
		Logs Files Settings	J	
		GEO info		
Flag	Country	Online Bots/All Bots	Detail State	
	Austria	(11/228)	0	
	Belgium	(1/5)	0	
Bos	inia and Herzegovina	(0/9)	0	
	Brazil	(0/4)	Ŵ	
	Bulgaria	(6/14)	<b>U</b>	
	Canada	(0/8)	•	
	China	(0/1)	•	
<u> </u>	Cyprus	(0/2)	<u> </u>	
	Denmark	(0/2)	<u> </u>	
	Estonia	(0/1)	<u> </u>	
	Europe	(0/2)	0	
	Finland	(6/13)	<u> </u>	
	France	(11/32)	<u> </u>	
N 📕 ——	French Gulana	(0/1)	<u> </u>	
	Germany	(0/241)		
	Greece	(0/2)		
×	Hong Kong	(11/26)	<b>U</b>	
	Hungary	(22/79)	<b>U</b>	
	India	(8/19)	<b>U</b>	
Irar	n, Islamic Republic of	(0/2)	V	

Source (webroot.com)



### Some notable examples

#### Flashback (2012–today)

- 600K compromised Macs (so, it's not just Windows)
- credentials stealing

#### Grum (2008–2012)

- 840K compromised devices,
- 40bln/mo spam emails

#### TDL-4 (2011–today)

- 4,5M compromised machines (first 3 months)
- known as "indestructible".

Cryptolocker (October 2013–today) NEW



### Roadmap

1. Botnets

#### 2. Communication channels

- 3. Domain generation algorithms (DGAs)
- 4. Detecting DGA-based botnets
- 5. Results



### Where is the my C&C server?

- 1. Where is my C&C server located?
- 2. Contact the C&C server
- 3. Receive command





### **C&C** channel: single point of failure





#### P2P is the natural answer.

# We focus on **centralized botnets** because they're still a **majority**.



### **Centralized C&C mechanisms**

#### Hardcoded IPs (e.g., 123.123.123.123)

- Bot software (malware) ships with the IPs
- Botmaster can update IPs regularly
- Knowing the IP makes takedown easy

#### Hardcoded domain names (e.g., cnc.example.com)

- Decouple IP from domain
- Botmaster free to change domain names and IPs
- Frequently changing IPs make takedown harder
- Botmaster must own many IPs



### Hardcoded domain names (2)





#### Hardcoded domain names (1)





### Roadmap

- 1. Botnets
- 2. Communication channels

#### 3. Domain generation algorithms (DGAs)

- 4. Detecting DGA-based botnets
- 5. Results



#### **Game-changing approach**

#### **Goals of the botmaster**

- Make the C&C server harder to locate
- Make the C&C channel resilient to hijacking

Reversing the malware binary should not reveal the location of the C&C nor any useful information toward that.



### Single domain vs. Domain flux

	vljiic.org	yxipat.cn	
	f0938772fb.co.cc	rboed.info	BOTS
	jyzirvf.info	79ec8f57ef.cc	
	hughfgh142.tk	gkeqr.org	
cnc.example.com	fyivbrl3b0dyf.cn	xtknjczaafo.biz	
	vitgyyizzz.biz	yxzje.info	
	nlgie.org	ukujhjg11.tk	DGA
1	aawrqv.biz		I
SINGLE DOMAIN predictable easy to leak	THOUSANDS OF DOMAINS PER DAY unpredictable impossible to leak		-



### **Domain of the day**





### Where is my C&C server?





# **Leveraging DNS**

- Only the botmaster knows the active domain
- The DNS protocol does the rest
- The **DGA** can be made more **unpredictable** (e.g., Twitter trending topic)

Reversing the malware binary only reveals the generation algorithm not the active domain of the day!



#### **Message in a bottle**







### Roadmap

- 1. Botnets
- 2. Communication channels
- 3. Domain generation algorithms (DGAs)
- 4. Detecting DGA-based botnets
  - 5. Results



## **Natural observation point: DNS**





### **Domain reputation systems**

#### Notos

• [Antonakakis et al., 2010]

#### KOPIS

• [Antonakakis et al., 2011]

#### **EXPOSURE**

- [Bilge et al., 2011]
- http://exposure.iseclab.org



#### They tell malicious vs. benign domains apart

#### No insights on what is the purpose of the domain

- C&C of what botnet?
- Could the same C&C be used for multiple botnets?
- Is the domain malicious for other reasons?
  - Phishing
  - Spam
  - Drive-by download



### More insights needed





### **NXDOMAINs**





### **Finding distinct DGAs**





#### **Drawbacks**

#### Needs an unpractical observation point

- No global view
- Hard to deploy

#### **Needs the IP of the clients**

• Privacy of the clients is not enforced



### **Lower level DNS servers**





### **OUR SOLUTION**



#### **Overview of our solution**





# Step 1: Linguistic analysis

We measure the "randomness" of the strings with respect to non-DGA-generated domains









### Linguistic features (2D PCA)



First principal component



### **Step 2: IP analysis**





# **Step 2: DBSCAN Clustering**



#### Cluster 1

Domains that, in their lifetime, have resolved to the very same IPs

#### **Cluster 2**

Domains that, in their lifetime, have resolved to the very same IPs

#### **Cluster 3**

Domains that, in their lifetime, have resolved to the very same IPs

#### Singleton (removed)



## **Real output (example)**

Palsit



# **Classifying new domains**





### Roadmap

- 1. Modern cybercrime
- 2. Botnets
- 3. Communication channels
- 4. Domain generation algorithms (DGAs)
- 5. Detecting DGA-based botnets
- 6. Results



# Step 1 on real data





#### Step 2 on real data



Correct clusters found: Conficker, Bamital, SpyEye, Palevo



## **DEMO** (come talk to me offline)

← → C 🗋 /clustering

#### **DGA Clustering**









🖌 🏹 🔕 🔍 🎜

# **Ongoing research**

#### Non-english baseline

- Italian domain names? Swedish domain names?
- Non-ASCII domains?
  - п.com
  - 葉瑶ou.io
  - ♥★≈♥.tk

#### Word-based DGAs

- concatenate random, valid words instead of letters
  - also-is-dom-yesterday-a-new.com



# **Questions?**





Federico Maggi federico@maggi.cc Politecnico di Milano