

SPOTSPAM

Tackling spam at new frontiers

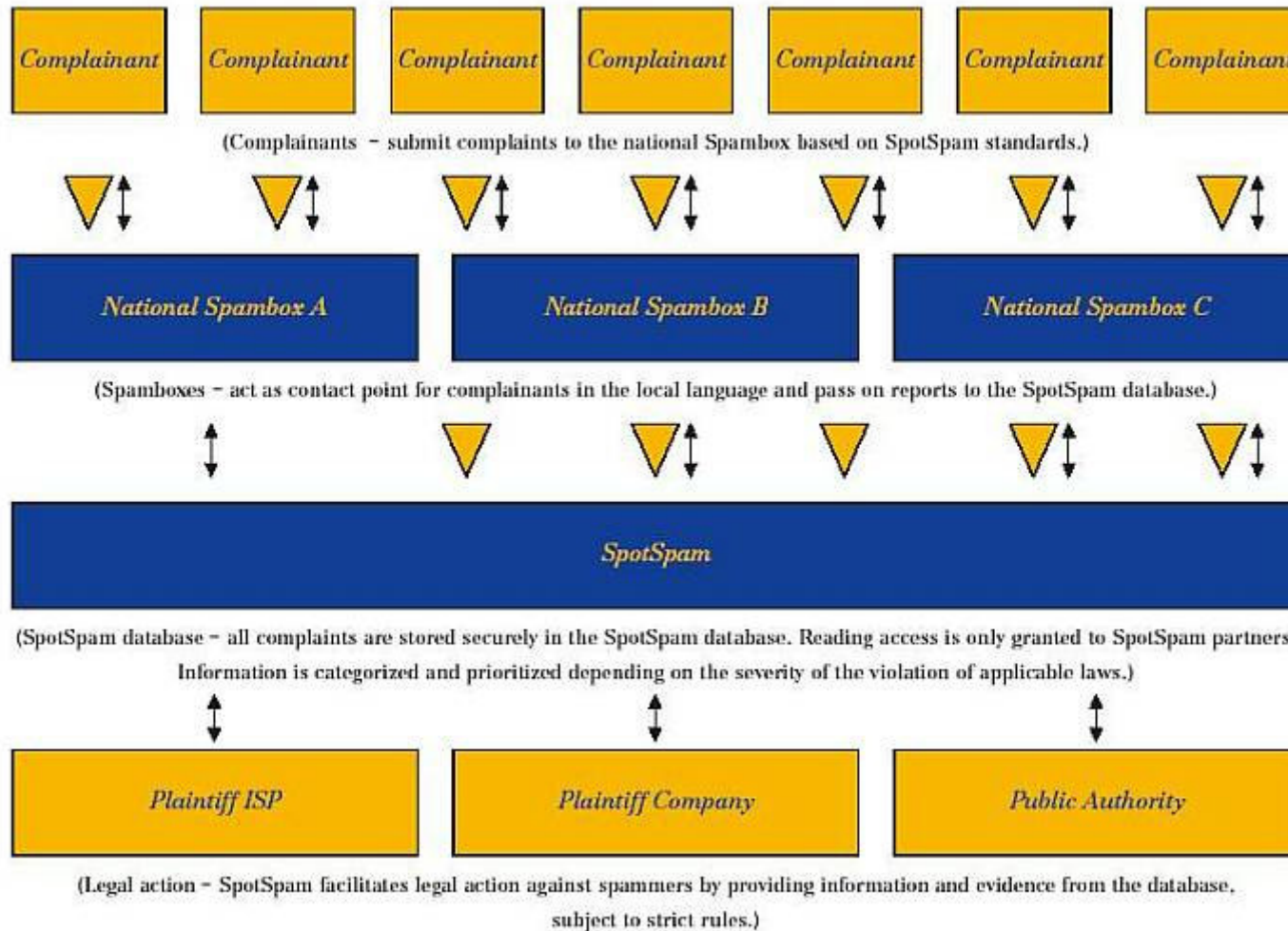
Przemek Jaroszewski
CERT Polska / NASK
przemek@cert.pl

- eco (Association of German ISPs) runs a hotline for reporting illegal content, including spam
- The hotline was approached by Microsoft looking for specific reports concerning Hotmail services
- Reports from eco's hotline users helped in firing a successful case against a spammer
- The idea came... why not make it on broader scale – with more reporting points and more parties interested in chasing spammers?
- eco and NASK put up a succesful EU project called SpotSpam

SpotSpam as an EU project

- The project was run under EC's Safer Internet Programme
- Consortium: eco (Association of German Internet Providers) & NASK
- Support: Microsoft
- Goals of the project: prepare legal and technical basis for gathering and sharing evidence against spammers
- Timeline: September 05 – September 07

What is SpotSpam?



- Multiple spam reporting points (spamboxes) are established across Europe
- Spambox operators sign an agreement with SpotSpam
- Complainants register with their local Spambox. They must agree to submit signed evidence in case a court case is launched. They must also certify that all reports submitted by them will actually be spam.
- The reports are stored in central database, which can be queried against IP ranges, email addresses, message subjects etc.
- Basic queries can be run by any registered party but only some indicative data is returned (eg. how many reports are found to match given criteria)
- Interested parties can request full (personal) data upon indetification if it is required for launching a court case

- Parties potentially interested in database querying or reports
 - **Internet providers (abuse of infrastructure)**
 - **Email providers (email forgery, quality of service)**
 - **Trademark owners (illegal „replicas“)**
 - **Law enforcement / government agencies**
 - **CSIRTs?**

Currently input can be accepted from:

- The prototype spambox application
- Unix mailbox files (bulk submission, mainly for testing)
- HTTP POST request (preferred for external cooperation)
- Email forwarding

The prototype database

- From each reported message the following set of information is extracted:
 - individual attachments
 - IP addresses and associated whois data
 - e-mail addresses
 - spamvertized URLs, associated domain, IP(s) and whois data
- Messages are clustered into spam campaigns to help identify waves of related spam and focus on broader picture rather than individual reports
- The campaigns are automatically classified according to their content (helps with priorities)
- Lots of information about IP addresses, domain addresses and their relations is collected, including whois information and geolocalisation

- Messages are clustered according to similarity of attachments, calculated as percentage of common Rabin fingerprints

Wikipedia quote:

Given an n -bit message m_0, \dots, m_{n-1} , we view it as a polynomial of degree $n-1$ over the finite field $\text{GF}(2)$.

$$f(x) = m_0 + m_1x + \dots + m_{n-1}x^{n-1}$$

We then pick a random irreducible polynomial $p(x)$ of degree k over $\text{GF}(2)$, and we define the fingerprint of m to be

$$f(x) \bmod p(x)$$

which can be viewed as a polynomial of degree $k-1$ or as a k -bit number.

- In short words, Rabin fingerprints in string algorithms:
 - **$F(m_0m_1\dots m_n) = m_0 * t^{n-1} + m_1 * t^{n-2} + \dots + m_n \bmod P$ is a fingerprint of a given substring**
 - **Time complexity of computing $F(m_1m_2\dots m_{n+1})$ from $F(m_0m_1\dots m_n)$ is $O(1)$**

- Other possibilities to explore
 - **Similarities of structure**
 - **Similarities of certain header fields**
 - **Common IP addresses**
 - **and plenty more..**

- Classification happens on two levels: individual reports and campaigns
- A reporter can assign type of spam to his report
- When enough reports in a campaign are of the same type X (in terms of both percentage and absolute number), the whole campaign is assumed to be of type X
- An operator can manually assign spam type to a campaign, in which case all messages are assumed to be of this type
- A Naive-Bayes classifier is used with number of classes equal to arbitrary number of spam types we want to recognise
 - **training happens when types are manually assigned (on any level)**
 - **NB attempts to classify new reports without pre-determined type**

The database can be queried against several fields:

- Subject contents
- IP addresses
- Email addresses
- URLs

Full text search is not really a good option for a very large database. But... we already have some Rabin fingerprints (of most popular substrings) and we can calculate more.

An external partner can only retrieve indicative numerical values while the operators are presented with full set of messages/campaigns that fit the criteria.

- The operator has access to all information about messages related to given URLs, emails or IP addresses.
- Data about misused URLs and IP addresses is periodically extracted from the database, mapped, and can be distributed to external partners (note: it does not include any information about reporters, message contents etc.)
- Complete report covering all data about a given campaign can be generated in pdf format. Such a report can be provided upon verified request for data.



Problems and lessons learned

- The main problem: the project delivered a pilot and EU support has ended
- While many parties have shown interest, the critical mass was never reached
- In many countries laws or practises are still inadequate

Further development?

- Automated screenshots
 - Better and more effective header analysis (emails and IP addresses)
 - Nicer presentation (eg. graphical geolocalisation)
 - Standalone application
 - Standalone applications with possibility of data exchange
-
- Any opinion will be valued and any help appreciated!

Campaign status

Name  

Created 03-09-2007 17:00:44 CEST

Modified 03-09-2007 17:00:44 CEST


Status

Priority

Type

Messages

x Name	Submit Date	Type	Attachments	Urls
<input type="checkbox"/> What IS OEM Software And Why DO You Care?	03-09-2007 16:57:20 CEST	Unknown	1	0
<input type="checkbox"/> Acrobat 8 PRO & Office 2007 \$79 NOW @ Joseph's WebSoft	03-09-2007 16:57:20 CEST	Unknown	1	0
<input type="checkbox"/> Buy OEM Software	03-09-2007 16:57:22 CEST	Unknown	1	0
<input type="checkbox"/> VISTA, ACR0BAT 8 PRO & OFFICE 2007 \$79 NOW at Rick [...]	03-09-2007 16:57:25 CEST	Unknown	1	0
<input type="checkbox"/> ACR0BAT 8 PRO & OFFICE 2007 \$79 NOW at Kari's WebShop	03-09-2007 16:57:23 CEST	Unknown	1	0
<input type="checkbox"/> VISTA, ACR0BAT 8 PRO & OFFICE 2007 \$79 NOW at Sami [...]	03-09-2007 16:57:25 CEST	Unknown	1	0
<input type="checkbox"/> Microsoft Office 2007, Acrobat 8 Pro 79\$ @ Walter' [...]	03-09-2007 17:00:04 CEST	Unknown	1	0

Subject Acrobat 8 PRO & Office 2007 \$79 NOW @ Joseph's WebSoft 
Campaign name [9716.03-09-2007 17:00:44](#)
Received 09-12-2006 20:41:54 CET
Reported 03-09-2007 16:57:20 CEST

Type 



Length 8720 bytes
SPF Status none
Number of URLs 0
Number of emails 6

Email	Role
alice@cert.pl	TO
alice@beregond.nask.waw.pl	OTHER
alice@cert.pl	LOCAL_RECIPIENT
alice@cert.pl	RECIPIENT
cofferwork@trannies-real.com	ENVELOPE_SENDER
cofferwork@trannies-real.com	FROM

Number of ips 3

IP	Role
195.187.245.33	LOCAL
195.187.245.33	RECIPIENTS_MX
66.184.2.193	SOURCE

Number of attachments 1

Filename	Content Type	Size
body  	text/plain; charset="us-ascii"	6610 bytes

```
Return-Path: <cofferwork@trannies-real.com>  
X-Original-To: alice@beregond.nask.waw.pl  
Delivered-To: alice@beregond.nask.waw.pl  
Received: from boromix.nask.net.pl (boromix.nask.net.pl [195.187.245.33])  
    by beregond (Postfix) with ESMTTP id B30E812C8132  
    for <alice@beregond.nask.waw.pl>; Sat, 9 Dec 2006 20:41:57 +0100 (CET)  
Received: from localhost (reverse.193.2.184.66.static.ldmi.com [66.184.2.193] (may be forged))  
    by boromix.nask.net.pl with SMTP id kB9JfqvH009819  
    for <alice@cert.pl>; Sat, 9 Dec 2006 20:41:54 +0100  
Message-ID: <000001c71bc95c9b7f08050100007f@localhost>
```

IP details

IP 66.184.2.193
Whois server whois.arin.net
Status DONE
Check date Mon Sep 03 17:00:08 CEST 2007

Domain whois info

OrgName: Ideal Technology Solutions US Inc.
OrgID: ITEC
Address: 27777 Franklin Road
Address: Suite 500
City: Southfield
StateProv: MI
PostalCode: 48034
Country: US

NetRange: 66.184.0.0 - 66.184.127.255
CIDR: 66.184.0.0/17
NetName: ITS-USNET
NetHandle: NET-66-184-0-0-1
Parent: NET-66-0-0-0-0
NetType: Direct Allocation
NameServer: DNS1.IDEALAPPS.COM
NameServer: DNS2.IDEALAPPS.COM
Comment:
RegDate: 2005-09-12
Updated: 2006-11-21

Related messages [ikona](#)
Role SOURCE
Asn AS14359
Country code US
Asn name

Search Criteria

IP:

IP range: from: to:

Email:

Subject:

Url:

Strict:

Case sensitivity:


Return as:

Messages

Name	Submit Date	Type	Attachments	Urls
" Similarly, it says: "There's no convincing evide [...]	03-09-2007 17:00:03 CEST	Unknown	3	0
bootlegger corrode	03-09-2007 17:00:02 CEST	Unknown	2	0

Move all found messages to:

Message details

Subject " Similarly, it says: "There's no convincing evidence that chlorella benefits humans in any way. 

Campaign name [9168.03-09-2007 17:00:44](#)

Received 22-12-2006 17:04:21 CET

Reported 03-09-2007 17:00:03 CEST

Type 

Length 32765 bytes

SPF Status none

Number of URLs 0






Number of emails 6

Email	Role
rtiofz@jurilan.nl	FROM
rtiofz@jurilan.nl	ENVELOPE_SENDER
alice@cert.pl	TO
alice@beregond.nask.waw.pl	OTHER
alice@cert.pl	LOCAL_RECIPIENT
alice@cert.pl	RECIPIENT

Number of ips 4

IP	Role
195.187.245.33	LOCAL
195.187.245.33	RECIPIENTS_MX
85.106.130.58	UNTRUSTED
209.201.220.228	UNTRUSTED

Number of attachments 3

Filename	Content Type	Size
body  	text/plain; charset="iso-8859-1"	2940 bytes
body  	text/html; charset="iso-8859-1"	4086 bytes
innovator.gif  	image/gif; name="innovator.gif"	19978 bytes

Return-Path: <rtiofz@jurilan.nl>

X-Original-To: alice@beregond.nask.waw.pl

Delivered-To: alice@beregond.nask.waw.pl

Received: from boromix.nask.net.pl (boromix.nask.net.pl [195.187.245.33])
by beregond (Postfix) with ESMTP id 9E02812C805D
for <alice@beregond.nask.waw.pl>; Fri, 23 Dec 2006 17:04:40 +0100 (CET)

- **More information can be obtained from:**

- <http://www.spotspam.net/>
- mail@spotspam.net
- myself in person or by email: przemek@cert.pl

CERT POLSKA

zgłaszanie incydentów: cert@cert.pl

strona internetowa: www.cert.pl

tel. +48 (22) 523 12 74

fax. +48 (22) 523 13 99

adres pocztowy:

NASK - CERT Polska

ul. Wąwozowa 18

02-786 Warszawa

Polska

DZIĘKUJEMY ZA UWAGĘ

ZMYŚL TELEKOMUNIKACJI

