# *Spam recognition by methods independent from text content*

## Virus Bulletin 2006 Montréal

**Ralf Iffert**
*ISS C-Force*

**Mark Usher**
*ISS C-Force*

INTERNET | SECURITY | SYSTEMS®

*Ahead of the threat.*™

# Conventional spam filters are ineffective

- **Circumvented by random text**
- **Outsmarted by spams without any text**
- **RBLs fooled by changing IPs with high frequency**

**INTERNET|SECURITY|SYSTEMS®**

# Introduction of two spam detection methods independent from text analysis

- **Structure Analysis**
  - Analysing the HTML structure of the email
- **Flow Analysis**
  - Analysing the flow of incoming emails

INTERNET|SECURITY|SYSTEMS®

# *Structure Analysis*

# Basic Idea

- **Remove all content from HTML part**

- **Calculate a hash on the remaining HTML structure**

- **Add hash to a database that is used for spam analysis**

INTERNET|SECURITY|SYSTEMS®

**Re: ParambYcy news**

Datei  Bearbeiten  Ansicht  Extras  Nachricht  ?

Antworten  Allen antw...  Weiterleiten  Drucken  Löschen  Zurück  Weiter  Adressen

**Von:** Crescencia Bottom
**Datum:** Mittwoch, 17. Mai 2006 19:46
**An:** de-info@iss.net
**Betreff:** Re: ParambYcy news

Viagra **$ 69,95** (10 tablets)                                    r1R
Valium **$ 105,45** (30 tablets)                          hWUS6**w**jcqee
Cialis **$ 99,95** (10 tablets)                            dCP9y**j**xlfin
                                                          i8kN**p**xrpjc
And many other http://kew76.obosome.com                         VZ
                                                                 uo

Zurück  Weiter  Adressen

Viagra **$ 69,95** (10 tablets)                                    v80
Valium **$ 105,45** (30 tablets)                          e0rOR**x**jqjte
Cialis **$ 99,95** (10 tablets)                           psk3Q**x**mhvlf
                                                          eADj**p**esdjr
And manyother http://oyh42.obosome.com                           a5
                                                                 I6

INTERNET|SECURITY|SYSTEMS®
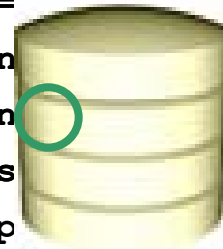
**Extract from source code of sample 1**

```
<DIV><FONT face=3DArial size=3D3><span style=3D" float : right ">r</span>V<=
span style=3D" float : right ">e</span>i<span style=3D" float : right ">e</=
span>a<span style=3D" float : right ">q</span>g<span style=3D" float : righ=
t ">c</span>r<span style=3D" float : right ">j</span>a <FONT color=3D#EC370=
7><STRONG>$ 69<span style=3D" float : right "> w </span>,95</STRONG></FONT>=
 (1<span style=3D" float : right "> S6 </span>0 t<span style=3D" float=
 : right "> WU </span>abIets)</FONT></DIV>
```

**Extract from source code of sample 2**

```
<DIV><FONT face=3DArial size=3D3><span style=3D" float : right ">v</span>V<=
sig://4A3DDB22F7C25943 (structhash)
span style=3D" float : right ">e</span>i<span style=
span>a<span style=3D" float : right ">j</span>g<span
t ">q</span>r<span style=3D" float : right ">j</span
1><STRONG>$ 69<span style=3D" float : right "> x </s
 (10<span style=3D" float : right "> OR </span>&nbsp
float : right "> Or </span>s)</FONT></DIV>
```
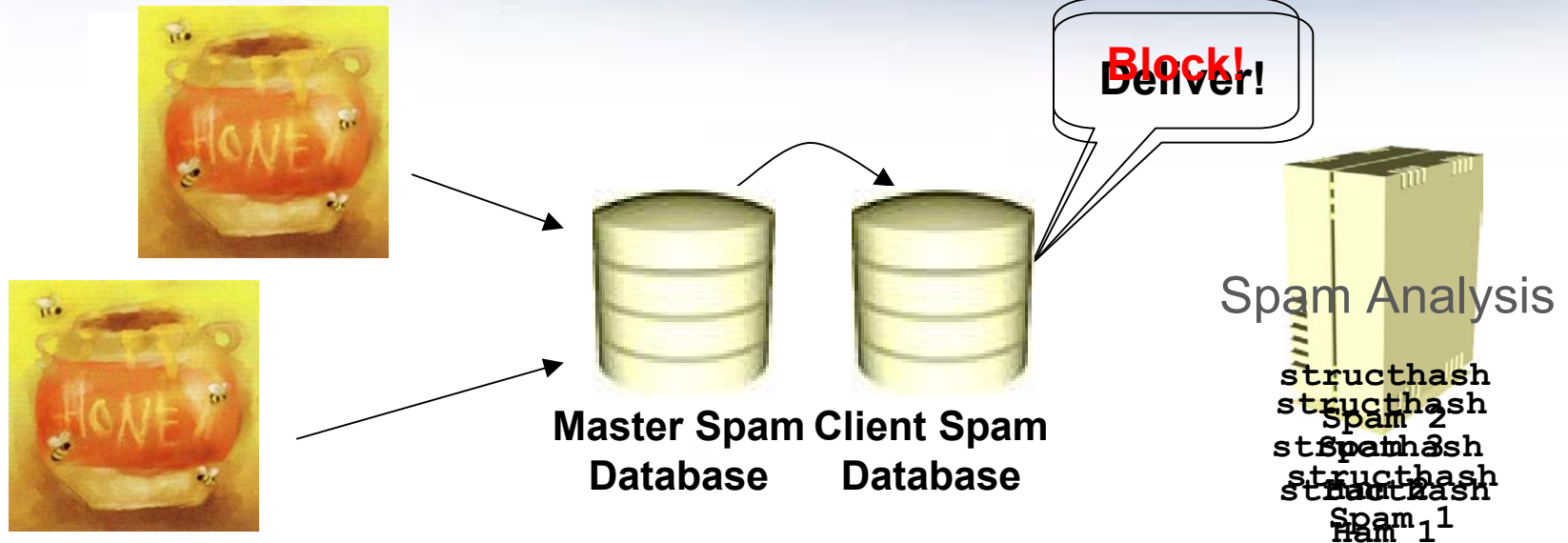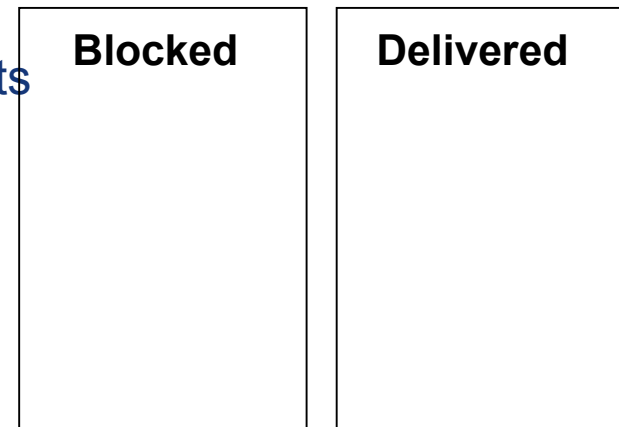
**Master Spam Database**

| Manufacturer's side | Client's side |
|---|---|



**Block!**

**Deliver!**

Spam Analysis

structhash
structhash
Spam 2
structhash
Spam 3
structhash
structhash
Ham 2
Spam 1
Ham 1

**Master Spam** **Client Spam**
**Database** **Database**

# Experimental Results

- **Voluminous inflow of spam**
- **Master Spam Database**
  - Directly delivered from the spam honey pots
- **Detection ratio: 45.1%**
- **Client Spam Database** **1.6%**
  - **Exclusively:**
    - Used for the local spam analysis
- **False positive ratio: 0.010%**
    - Updated by the Master Spam Database

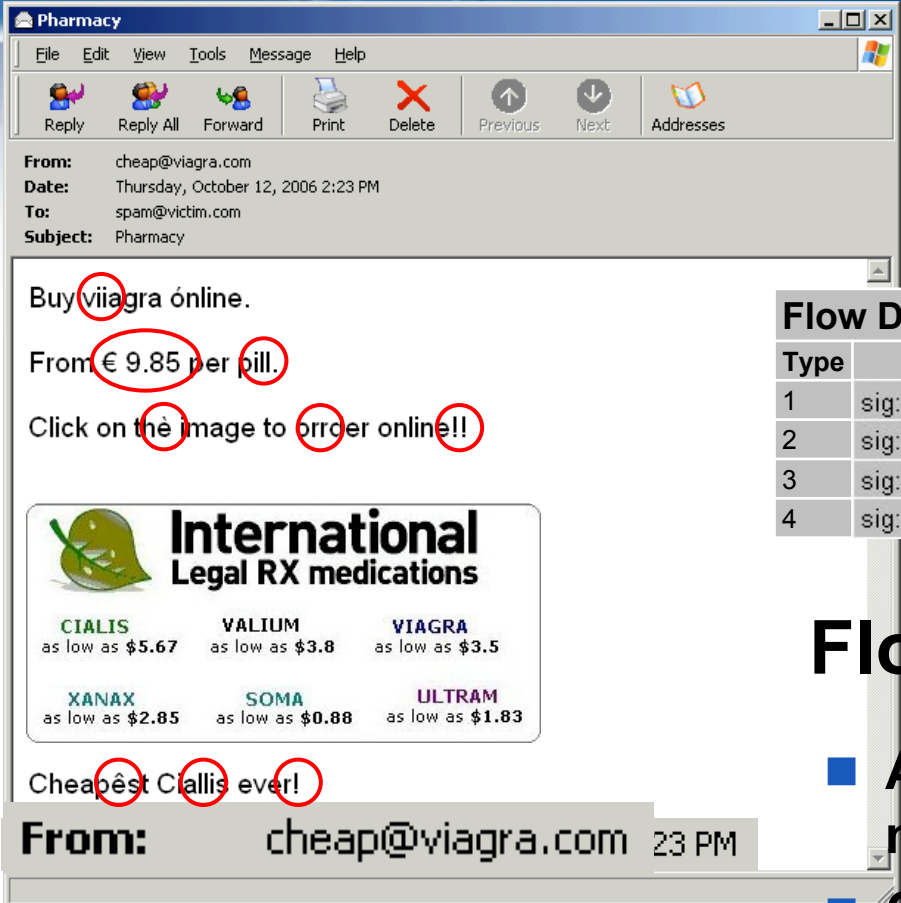| Blocked | Delivered |
|---|---|
|  |  |

# *Flow Analysis*

# Basic Idea

# Identifying "similar" emails arriving within a small time frame

# →Detection of whole spam threads

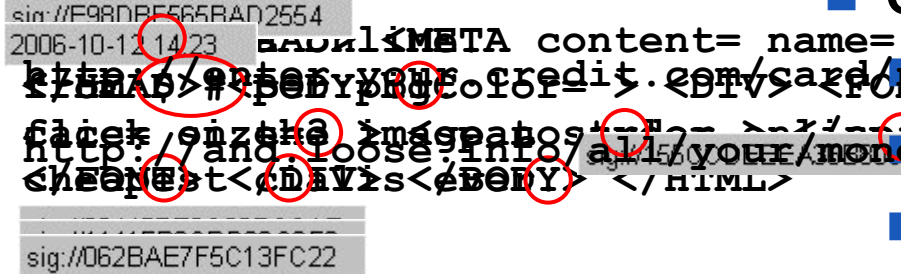# →Similarity of emails is determined by similarity measures

## Similarity Measures

**Flow Database**

| Type | SimilaritySignature | Sender | TimeStamp |
|---|---|---|---|
| 1 | sig://90445BE9C82FCCAE | sig://E98DBF565BAD2554 | 2006-10-12 14:23 |
| 2 | sig://1141FB0CBD68C0F8 | sig://E98DBF565BAD2554 | 2006-10-12 14:23 |
| 3 | sig://062BAE7F5C13FC22 | sig://E98DBF565BAD2554 | 2006-10-12 14:23 |
| 4 | sig://155CCDEEEA36B809 | sig://E98DBF565BAD2554 | 2006-10-12 14:23 |

## Flow Database

- Computed body text
- Administration of similarity measures
- Set of image URLs
- Columns of the Flow Database
  - Type of similarity measure
  - Set of attachments
  - Similarity signature
  - Time when entry was generated
  - Sender of the email

*Further measures are conceivable…*

# Usage of the Flow Database

- **Conditioned by the two parameters:**
  - Threshold for similar emails
  - Time-frame that is monitored
- **Only contains information of emails received within the time-frame**

| Flow Database | | | Threshold: 3 similar emails | Time-frame: 10 min |
|---|---|---|---|---|
| Type | Similarity | Signature | Sender | TimeStamp |
| 1 | sig://90445BE9C82FCCAE | | sig://E98DBF565BAD2554 | 2006-10-12 14:23 |
| 2 | sig://1141FB0CBD68C0F8 | | sig://E98DBF565BAD2554 | 2006-10-12 14:23 |
| 3 | sig://062BA57FD0139C22 | | sig://E98DBF565BAD2554 | 2006-10-12 14:23 |
| 4 | sig://155CCDEEEA36R80? | | sig://E98DBF565BAD2554 | 2006-10-12 14:23 |
| | ?BD739DC | | sig://C350770FF2B567? | 2006-10-12 14:25 |
| | 77? | | sig://? | 2006-10-12 14:25 |
| 2006-10-12 14:29 | 44CF5D6CE | | sig://? | 2006-10-12 14:25 |
| 2006-10-12 14:29 | 3814C37C4 | | sig://? | 2006-10-12 14:25 |
| 2006-10-12 14:29 | | | | 2006-10-12 14:26 |
| 2006-10-12 14:29 | CBD68C0F8 | | sig://6937759F8A377BFE | 2006-10-12 14:26 |
| 3 | sig://77B27F4E2484457A | | sig://6937759F8A377BFE | 2006-10-12 14:26 |
| 4 | sig://BF844E7A1B8?...2 | | sig://6937759F8A377BFE | 2006-10-12 14:26 |
| 1 | sig://F51441B2B751?EE | | sig://0CF3A801615F4CE4 | 2006-10-12 14:27 |
| 2 | sig://1141FB0CBD68C0F8 | | sig://0CF3A801615F4CE4 | 2006-10-12 14:27 |
| 3 | sig://5201C2A2BE66FC2A | | sig://0CF3A801615F4CE4 | 2006-10-12 14:27 |
| 4 | sig://59203949BA35933C | | sig://0CF3A801615F4CE4 | 2006-10-12 14:27 |
| 1 | sig://7C061C27A?025EB | | sig://6D495F90A5CBC8DC | 2006-10-12 14:29 |
| 2 | sig://1141FB0CBD68C0F8 | | sig://6D495F90A5CBC8DC | 2006-10-12 14:29 |
| 3 | sig://B8CCF8F9?...? | | sig://6D495F90A5CBC8DC | 2006-10-12 14:29 |
| 4 | sig://C0561E4A370C83B2 | | sig://6D495F90A5CBC8DC | 2006-10-12 14:29 |

# Experimental Results

- **Detection ratio: 72.7%**
  - Exclusively: 0.4%
- **False positive ratio: ~0%**

**4 min**

**Precondition: About 50,000 spams per day or more!**

# Optimizations

- Global usage of Flow Database
- Delay email delivery according to time-frame

# Computing time

## Computing time depends on

- $P_i$ Averaged time to extract the similarity data of the $i^{th}$ similarity measure

- $M_i$ Averaged hash calculation time on the $i^{th}$ similarity data

- $R$ Time for one database request

- $N$ Number of similarity measures

→ **Computing time per mail:**
**of the $i^{th}$ similarity measure**

$$\sum_{i=1}^{N}\left(P_i + M_i + R\right)$$

INTERNET|SECURITY|SYSTEMS®

# Memory requirements Flow Database

**Required memory depends on**

- *M* Maximal number of entries in Flow Database
- *S* Size in bytes of each Flow Database entry

→ **Required memory:** *M\*S*

**Example**

**M depends on**

- *N* Number of similarity measures
- *E* Throughput in emails per minute
- *X* Time-frame in minutes used for the Flow Analysis

→ *M = N\*E\*X*

*= 4*
*= 600*
*= 10*

**S depends on**

- *C* Number of bytes consumed by a Flow Database entry
- *O* Memory overhead

→ *S = C+O*

*= 44*
*= 12*

→ **Required memory:** *N\*E\*X\*(C+O)*

*= 1.28 MB*

# Further approaches

- **Automated detection of random text**
- **Usage of visual features**
- **Image signatures invariant against random variations**



```
sig://C3B90474A349823E    sig://C3B90474A349823E    sig://C3B90474A349823E
```

INTERNET|SECURITY|SYSTEMS®

# *Many thanks for your attention!*

## *Q & A*

**Ralf Iffert
ISS C-Force
riffert@iss.net**

INTERNET | SECURITY | SYSTEMS®

*Ahead of the threat.™*