



## Spam Campaign Detection, Analysis, and Investigation

*By*

**Son Dinh, Taher Azeb, Francis Fortin, Djedjiga Mouheb and Mourad Debbabi**

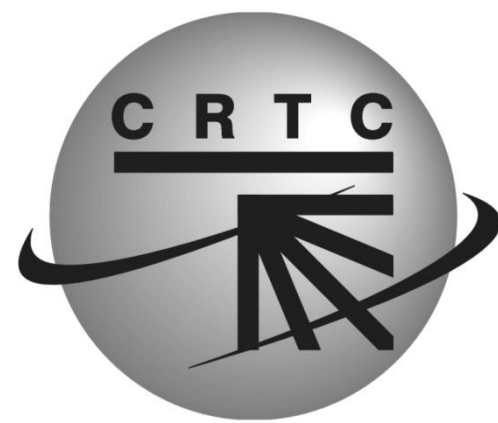
*Presented At*

The Digital Forensic Research Conference

**DFRWS 2015 EU** Dublin, Ireland (Mar 23<sup>rd</sup>- 26<sup>th</sup>)

DFRWS is dedicated to the sharing of knowledge and ideas about digital forensics research. Ever since it organized the first open workshop devoted to digital forensics in 2001, DFRWS continues to bring academics and practitioners together in an informal environment. As a non-profit, volunteer organization, DFRWS sponsors technical working groups, annual conferences and challenges to help drive the direction of research and development.

**<http://dfrws.org>**



# **SPAM CAMPAIGN**

## **DETECTION, ANALYSIS AND INVESTIGATION**

S. Dinh, T. Azab, F. Fortin, D. Mouheb, M. Debbabi

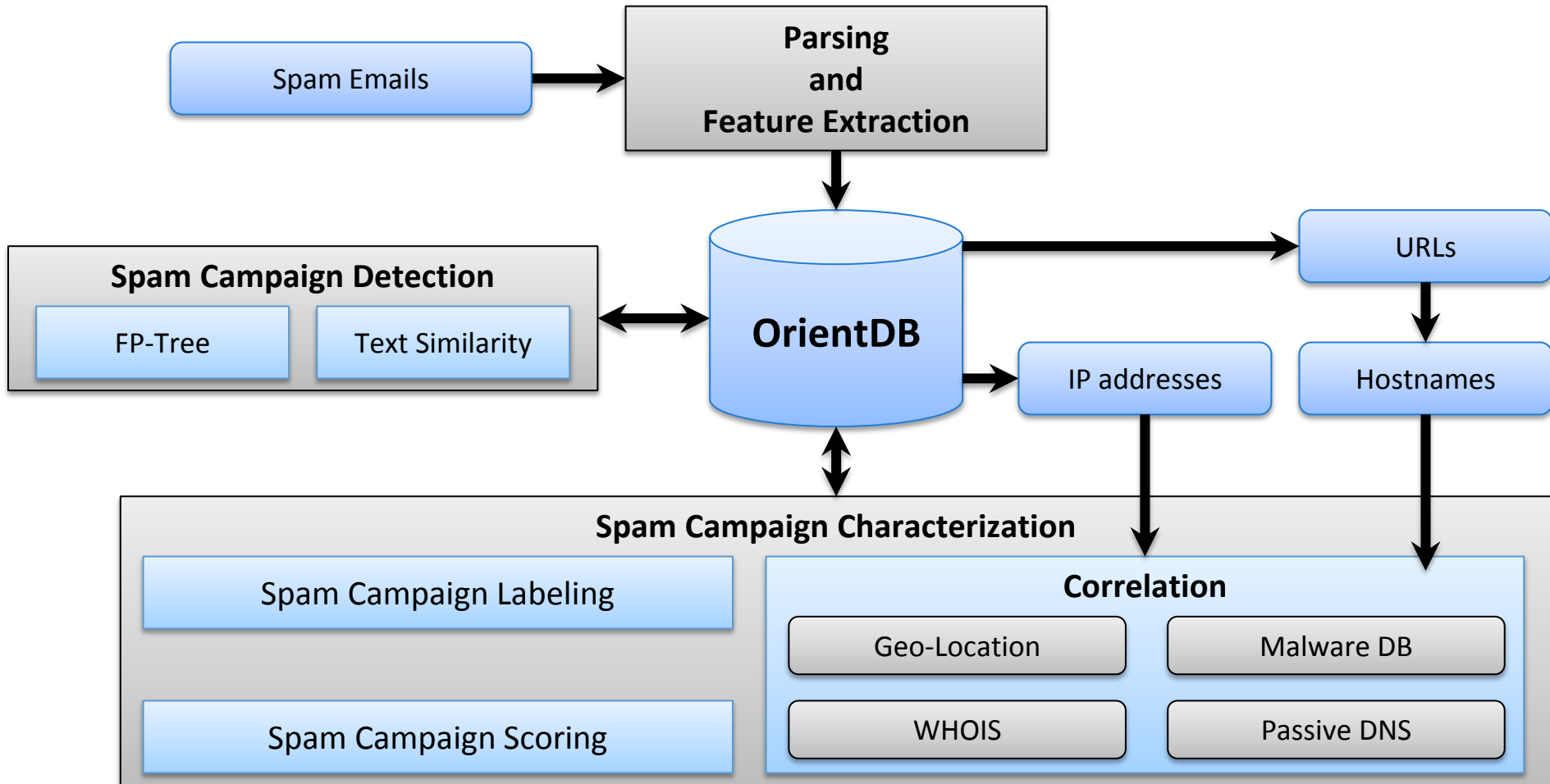
# Outline

2

- ❑ Architecture
- ❑ Central Database
- ❑ Parsing and Feature Extraction
- ❑ Spam Campaign Detection
- ❑ Spam Campaign Characterization
- ❑ Results

# Architecture

3



# Central Database

4

- Relational database management system (RDBMS)
- The number of spam emails exceeds 6 *trillion* in Q4 2014
  - McAfee Labs Threats Report (February 2015)
- OrientDB:
  - Flexibility (document database)
  - Interconnectivity (graph database)
  - Scalability
- Documents: spam campaigns, emails, IP addresses, domain names, attachments
- Connections and interconnections between spam campaigns, emails, IP addresses, domain names and attachments

Campaign Class
List of messages
Number of messages
Duration
Summary
Score
Label
...

Message Class	
Headers	The Origination Date Field
	...
	Trace Fields
	...
	X-Headers
Body	...
	Text
	URL(s)
	Attachment(s)

IP Class	
	IP
Geolocation	City
	Country
	...
	ISP
Malware	Hash(es)
	Timestamp(s)

Attachment Class
File Name
File Type
File Size
Hash
Base64
...

Hostname Class	
	Hostname
PDNS	First Seen
	Last Seen
	IP(s)
	Count
	WHOIS_ID

WHOIS Class
2 <sup>nd</sup> Level Domain
Creation Date
Expiration Date
Registrar
Registrant
Name Servers
...

6

# Preprocessing

# Parsing and Feature Extraction

7

## □ Parser

- ▣ Parse and store standard header fields (RFC5322)
- ▣ Unpack email bodies (single-part or multi-part)
- ▣ Extract embedded URLs

## □ Feature Extractor

### ▣ Content Type

- text/\*
- multipart/\*
- image/\*
- application/\*
- etc.

### ▣ Character Set

- utf-8
- iso-8859-1
- windows 1250
- shift-jis
- koi8-r
- etc.

### ▣ Subject

- Decoded to Unicode
- ▣ URL Tokens
- ▣ Attachment



# Parsing and Feature Extraction

8

## □ Feature Extractor

### ▣ Email Layout

#### ▣ text layout

- T for a paragraph with only text
- U for a URL
- N for an empty line
- For example:  
TNTNUNN

#### ▣ HTML layout

- Top levels of the DOM tree
- For example:

```
<html>
-<head>
-</head>
-<body>
--<p></p>
--<br />
--<div></div>
-</body>
</html>
```

#### ▣ multipart layout

- The structure of the multipart email
- For example:

```
multipart/mixed
-multipart/alternative
--text/plain
--multipart/related
---text/html
---img/jpg
-application/pdf
```

# 9 Detection & Characterization

# Spam Campaign Detection

10

- Distance between two spam emails:
  - ▣ w-shingling and the Jaccard coefficient
  - ▣ Context Triggered Piecewise Hashing (CTPH)
  - ▣ Locality-Sensitive Hashing (LSH)
- Hierarchical, partitional, neural network-based, kernel-based clustering techniques
- Problems:
  - ▣ Scalability
  - ▣ Obfuscation techniques
    - Spam email templates
    - Randomly picked paragraphs from books or Wikipedia articles
    - Randomly generated subdomains and fast-flux service networks

# Spam Campaign Detection

11

- Frequent-Pattern Tree (FP-Tree)
  - ▣ The more frequent a feature is, the more it is shared among spam emails
  - ▣ Less frequent features correspond to the obfuscated parts
  - ▣ Two scan of the dataset:
    - First scan to compute the number of occurrences for each feature
    - Second scan to insert feature vectors into the tree
  - ▣ The cost of inserting a feature vector  $f_v$  into the FP-Tree is  $O(|f_v|)$ , where  $|f_v|$  is the number of features in  $f_v$ .
  - ▣ The **Content Type** feature is put at the beginning of each feature vector
  - ▣ The unique ID of each email is kept at the end of each feature vector
  - ▣ Embedded URLs are split into tokens

# Spam Campaign Detection

12

## □ Spam Campaign Identification

### ▣ Traverse the tree

### ▣ Conditions:

- The number of children  $\geq \text{min\_num\_children}$
- The average count of children  $\geq \text{freq\_threshold}$
- The path to root must contain **one** feature type not in `n_obf_features`
- The number of leaves of the sub-tree  $\geq \text{min\_num\_messages}$

### ▣ If a node satisfies the conditions:

- The leaves of the sub-tree are spam emails from the same campaign
- The path from the root to this node contains the common features
- The sub-tree is then removed from the tree



text/plain

root

iso-8859-1

NTNTTNTTTTNN



NTNTTNTTTTTTNUNNN

wrend.ru

Spam Email

...REPLICA WATCHES...

Pens

Bracelet

Lighters

Bags

Louis Vuitton Bags & Wallets

...Replica watches...

...replica watch

facap.ru

sj -split-here- ALL MAJOR DESIGNER REPLICA WATCHES 5

hn -split-here- wrend.ru 13  
sj -split-here- I can't believe you helped save over \$200 on this Cufflinks 3

sj -split-here- Pens 1

sj -split-here- Bracelet 4

sj -split-here- Lighters 12

sj -split-here- Bags 2

sj -split-here- Louis Vuitton Bags & Wallets 1

sj -split-here- View Our Wholesale Rolex Replica watches Today 4

sj -split-here- What to look for when purchasing a replica watch 1

1357345410.P25162M720949Q5997399.accidentally.kelly.st 1

1357330171.P25162M601968Q5945802.accidentally.kelly.st 1

1357230179.P25162M258010Q5453140.accidentally.kelly.st 1

1357256634.P25162M333393Q5537306.accidentally.kelly.st 1

1357223678.P25162M123283Q5430727.accidentally.kelly.st 1

1357223678.P25162M121824Q5430726.accidentally.kelly.st 1

1357219382.P25162M9230Q5403316.accidentally.kelly.st 1

1357254111.P25162M264054Q5525460.accidentally.kelly.st 1

1357240068.P25162M698145Q5482776.accidentally.kelly.st 1

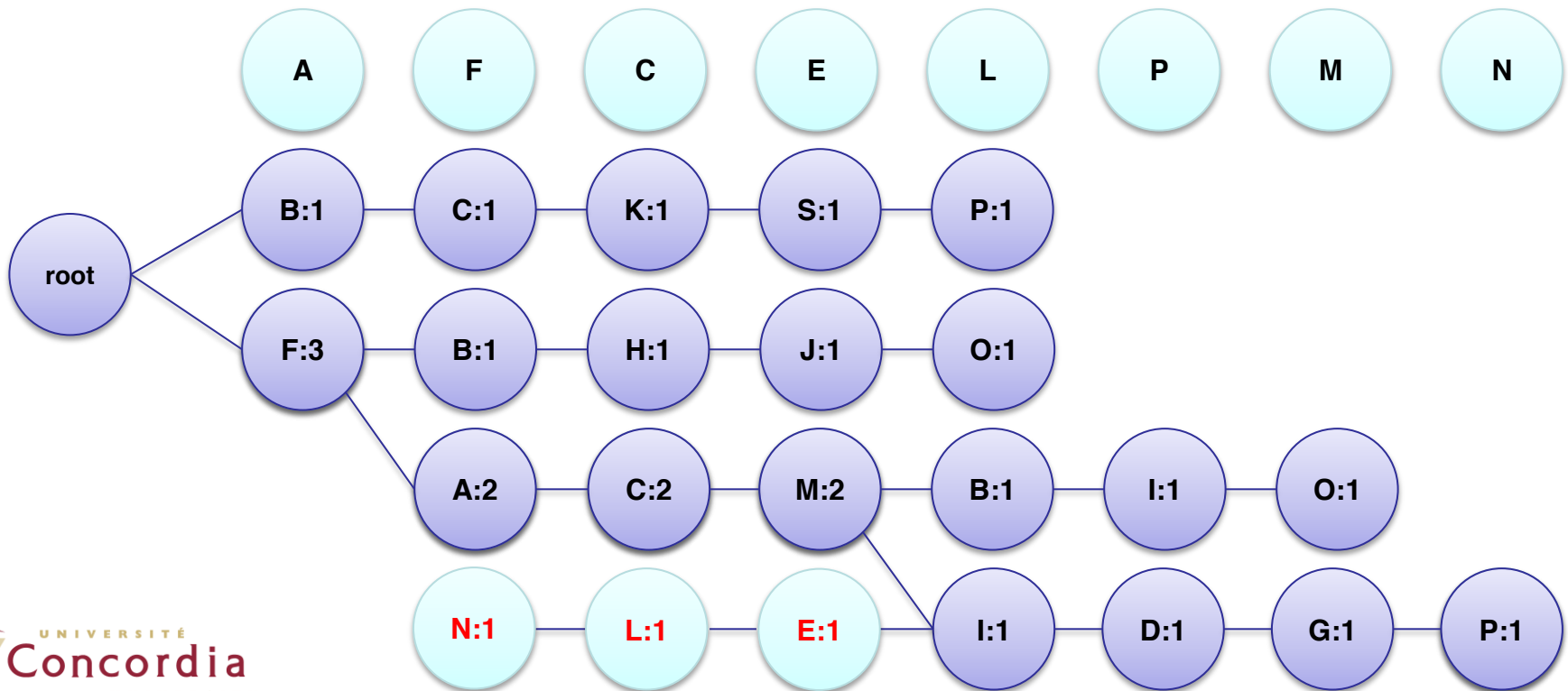


# Spam Campaign Detection

16

## □ Incremental Frequent-Pattern Tree

- Spam campaigns may last for a long period of time and therefore should be identified in their early stage
- Feature vectors are extracted as soon as spam emails arrive and are inserted into the FP-Tree from the root level



# Spam Campaign Characterization

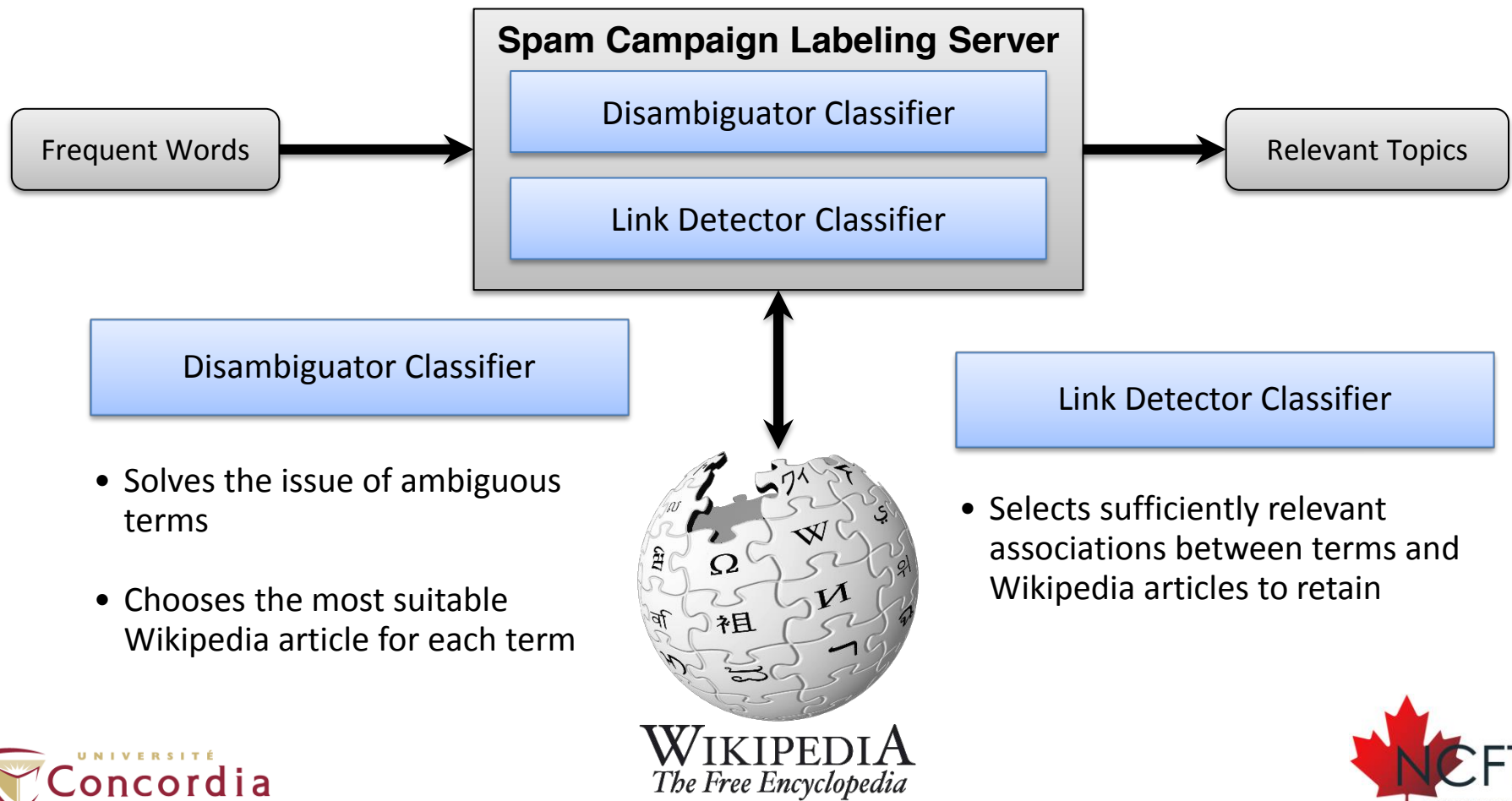
17

- 
- ```
graph LR; IP[IP addresses] --> Geo[Geo-Location]; IP --> MD[Malware Database]; IP --> PD[Passive DNS]; Host[Hostnames] --> MD; Host --> PD; Dom[2nd-level domains] --> WHOIS[WHOIS];
```
- IP addresses**
- From header fields
    - Received
    - X-\*
  - From URLs
    - using passive DNS
- Hostnames**
- From URLs
- 2<sup>nd</sup>-level domains**
- Geo-Location**
- City, Country
  - ISP
  - Organization
- Malware Database**
- Passive DNS**
- Timeslot (First/Last seen)
  - IP address(es)
  - Count
- WHOIS**
- Name servers
  - Domain status
  - Creation date
  - Expiration date
  - Registrar & registrant
- UNIVERSITÉ  
Concordia

# Spam Campaign Characterization

18

## □ Spam Campaign Labeling



# Spam Campaign Characterization

19

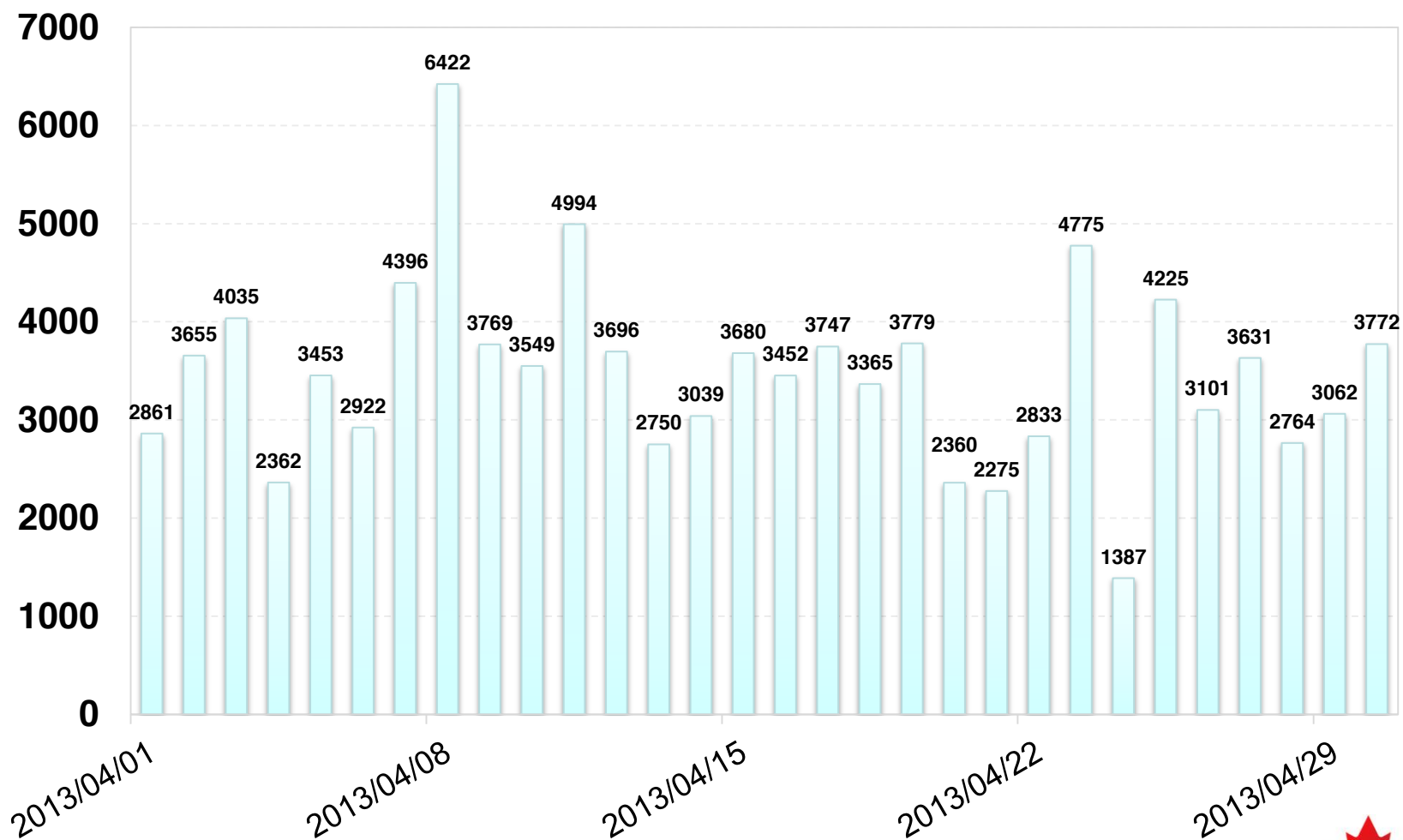
- Spam Campaign Scoring
  - An investigator may only need to pursue a particular objective
  - Each spam campaign is assigned a score computed using:
    - The number of spam emails inside a campaign
    - The number of IP addresses in Canada
    - The number of domain names that are resolved to IP addresses in Canada
    - The number of “.ca” TLDs that appear in the from field
    - The number of “.ca” TLDs that appear in the embedded URLs
    - The number of Canadian city names that appear in the content
    - The number of appearances of the string “Canada” in the content
    - The number of IP addresses that are associated with malware
    - The number of IP addresses that belong to a specific IP range
  - Each criterion has a customizable weight
  - The criteria have been verified by a law enforcement official

20

# Results

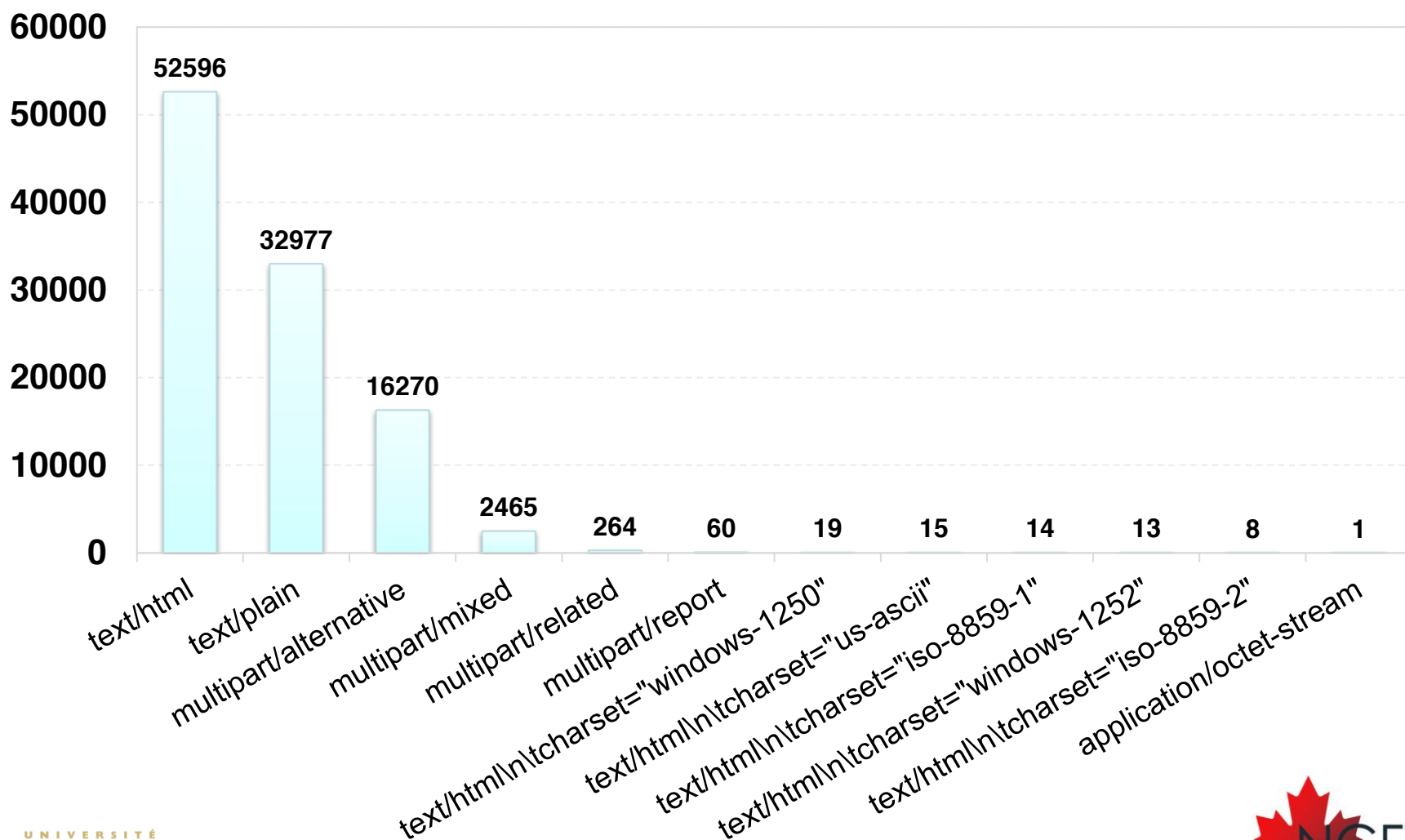
# Spam Data

21



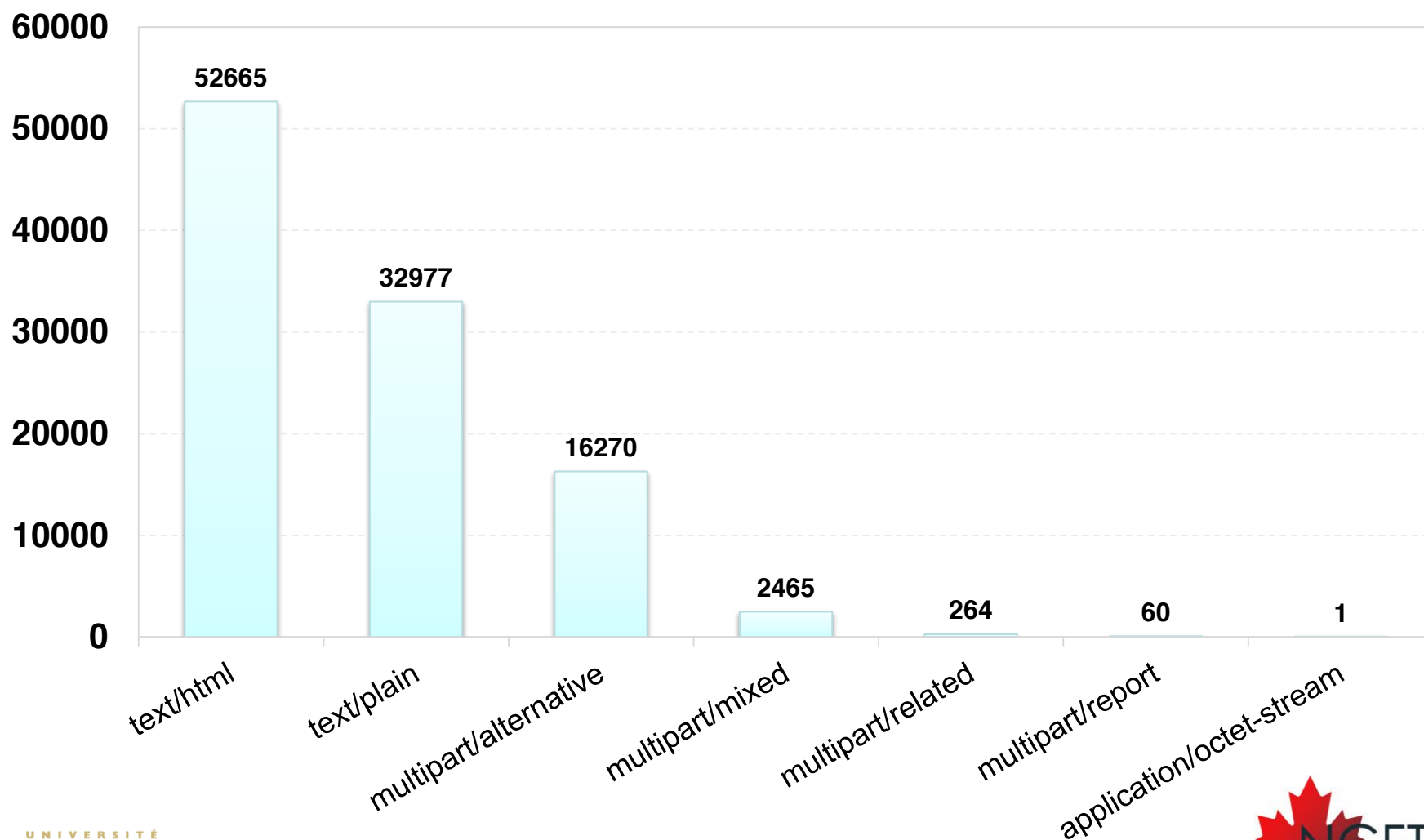
# Raw MIME Types

22

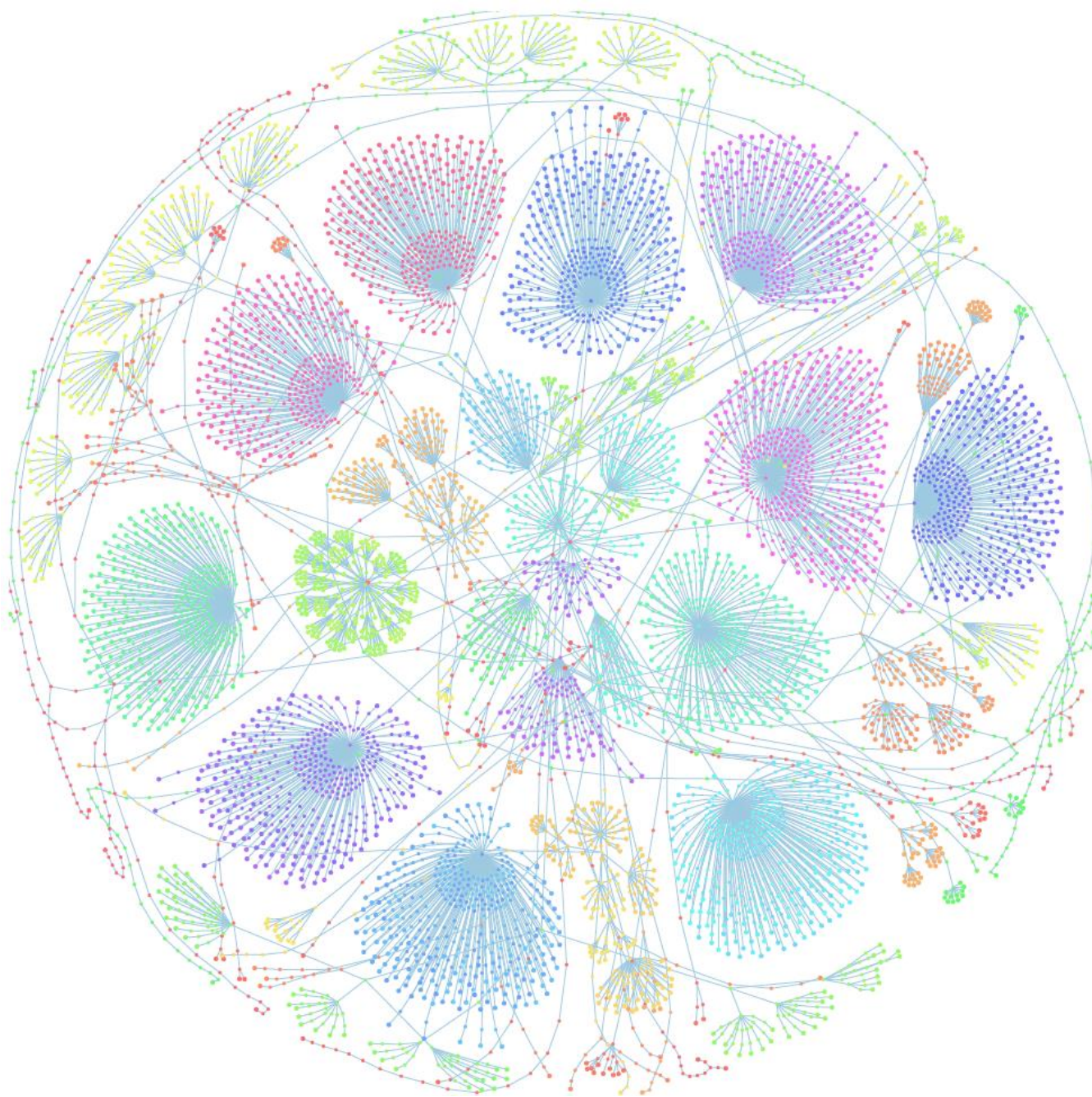


# Correct MIME Types

23

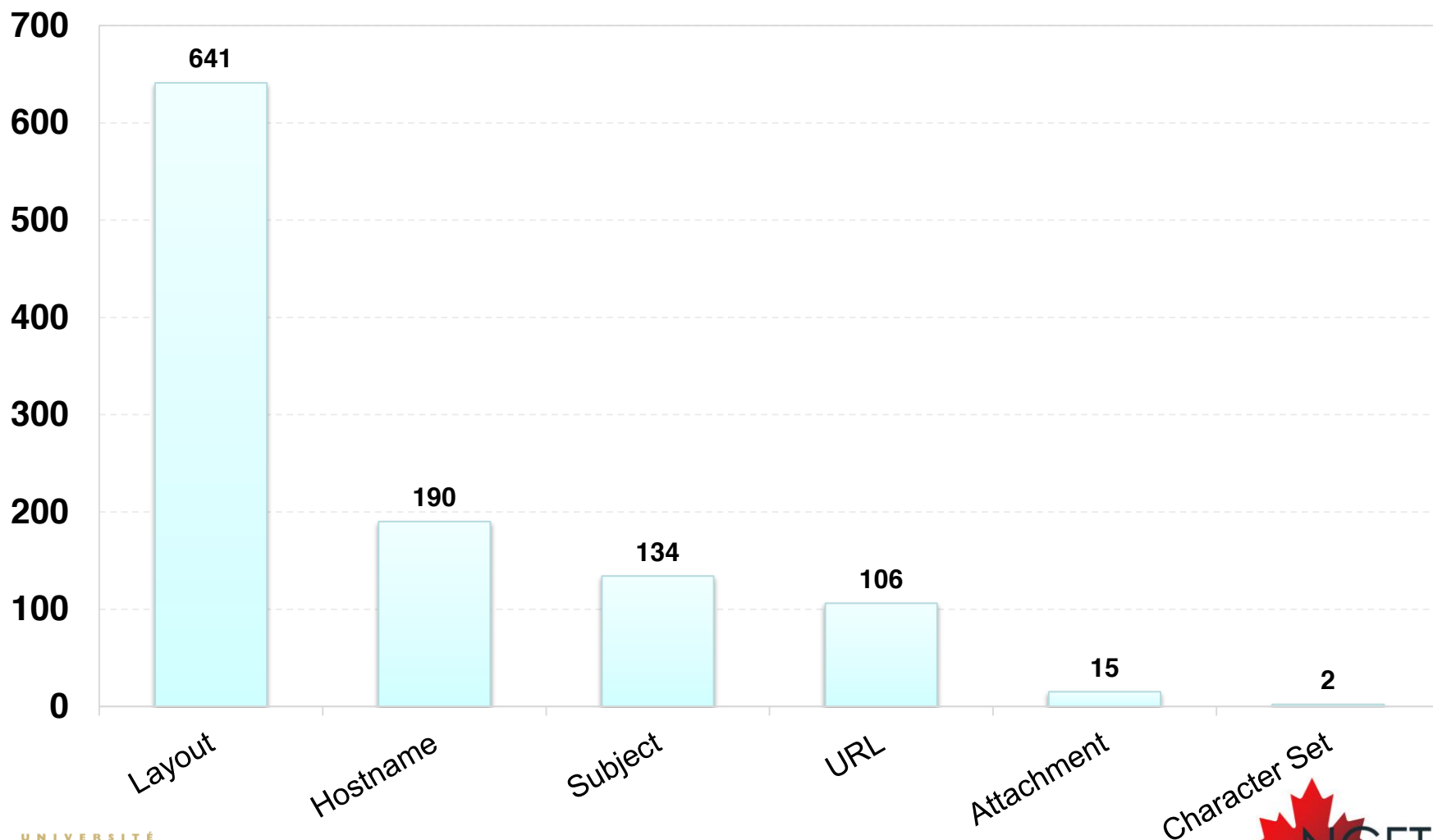






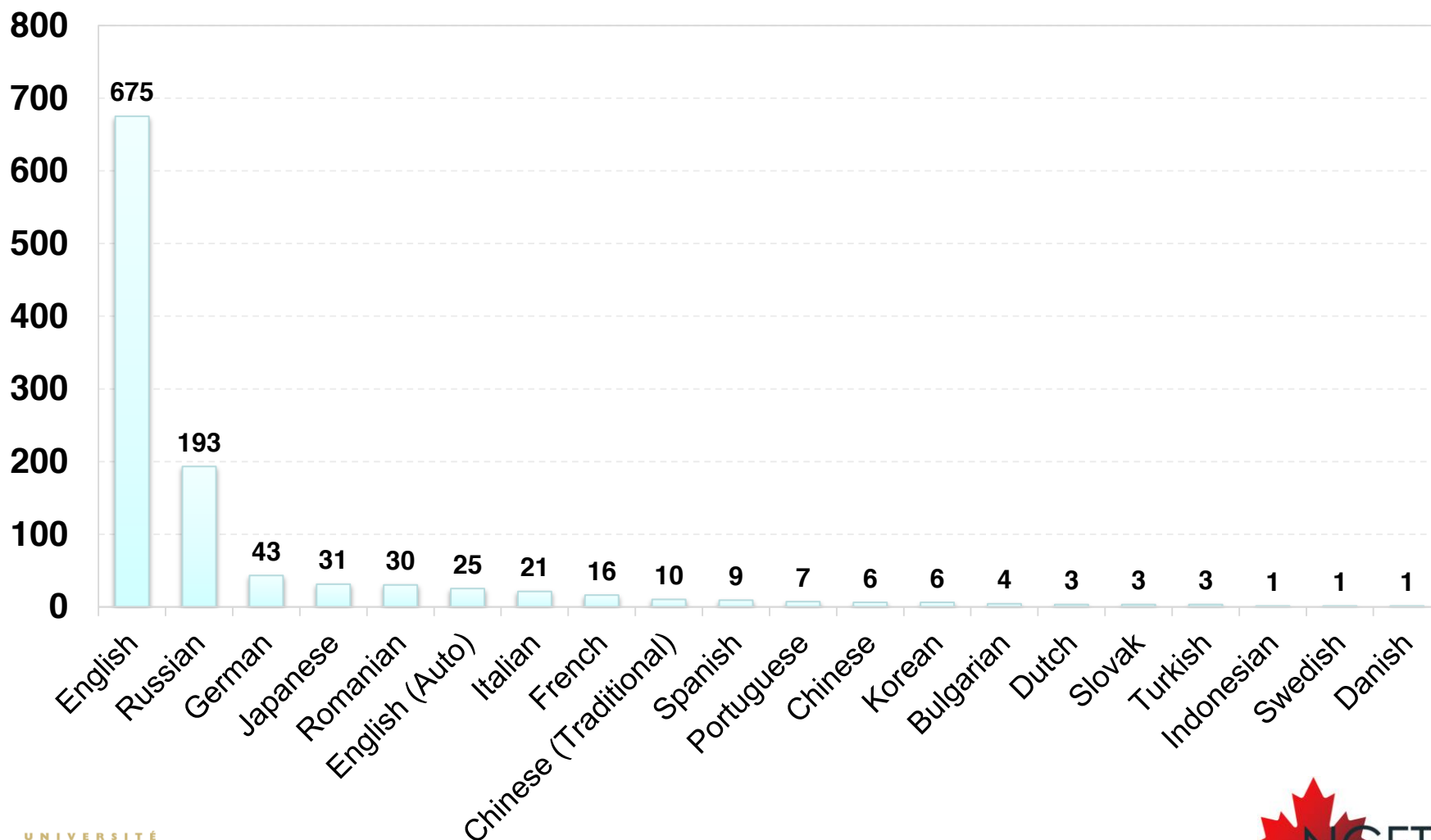
# Decisive Features

25



# Spam Campaign Languages

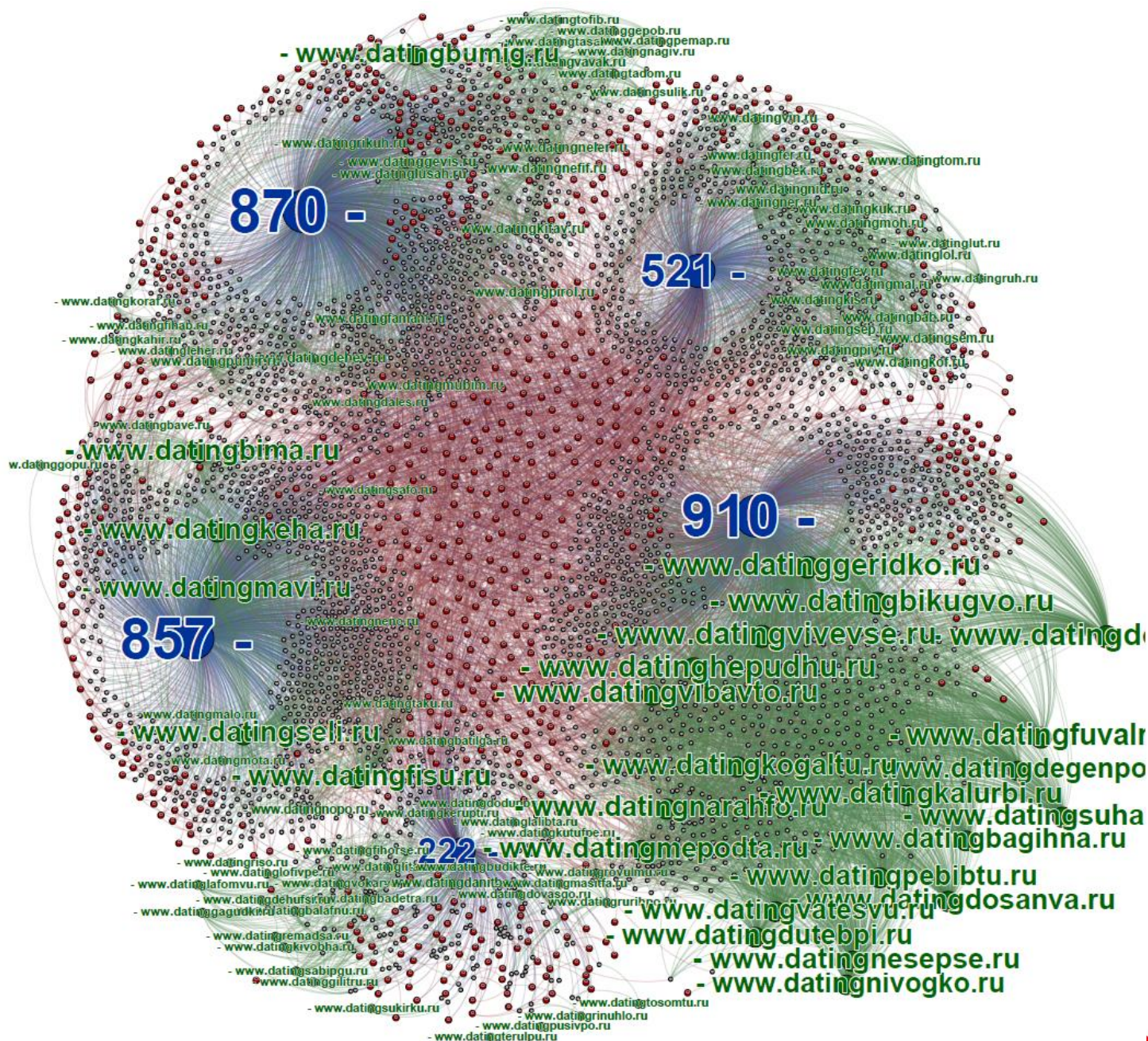
26



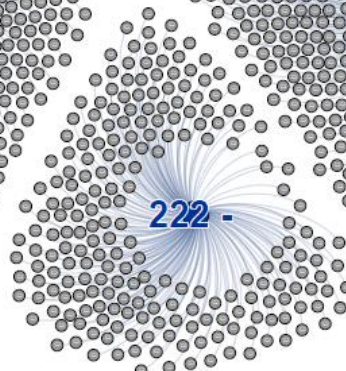
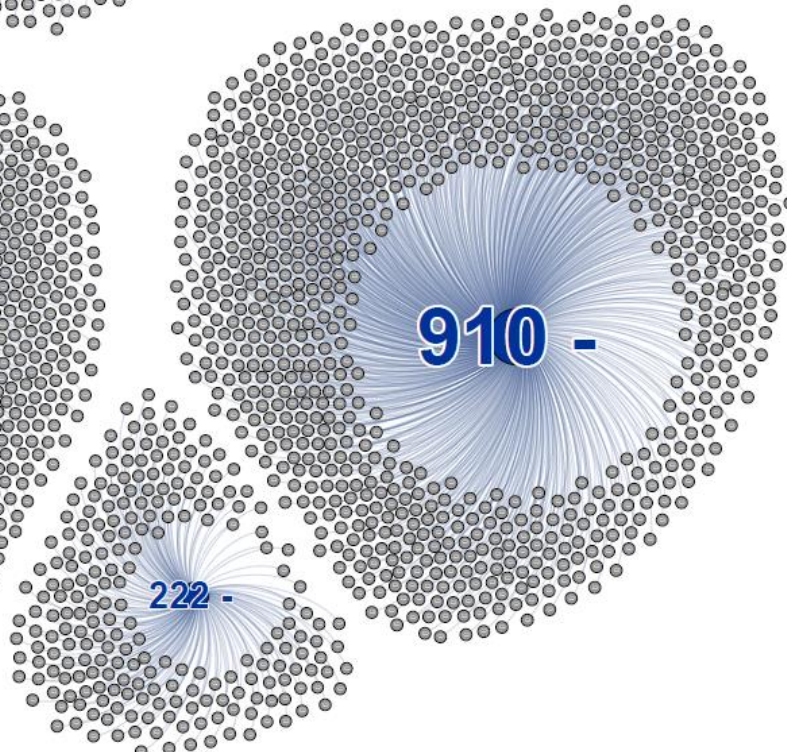
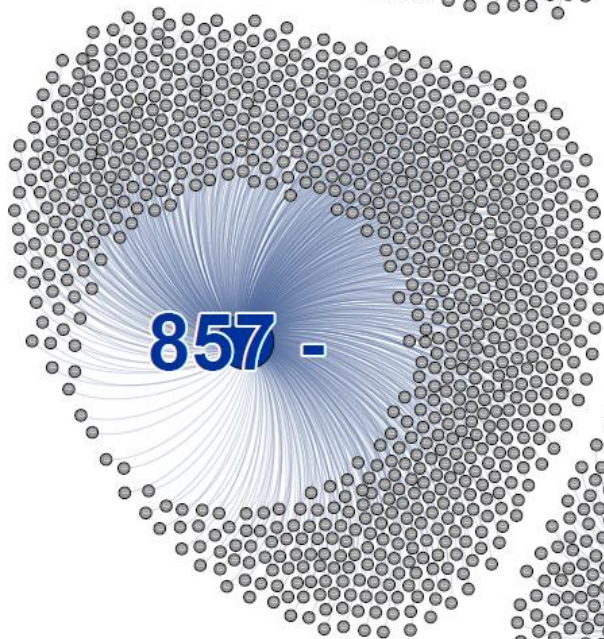
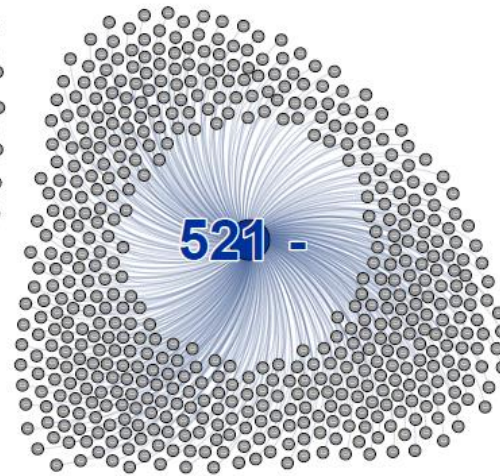
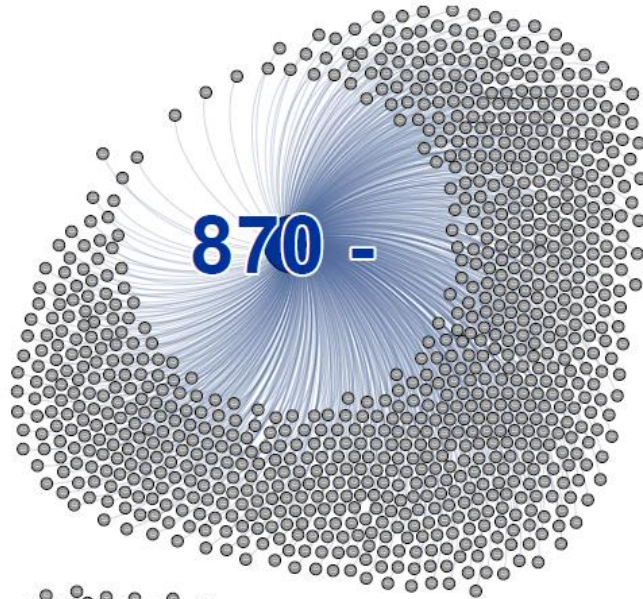
## 27



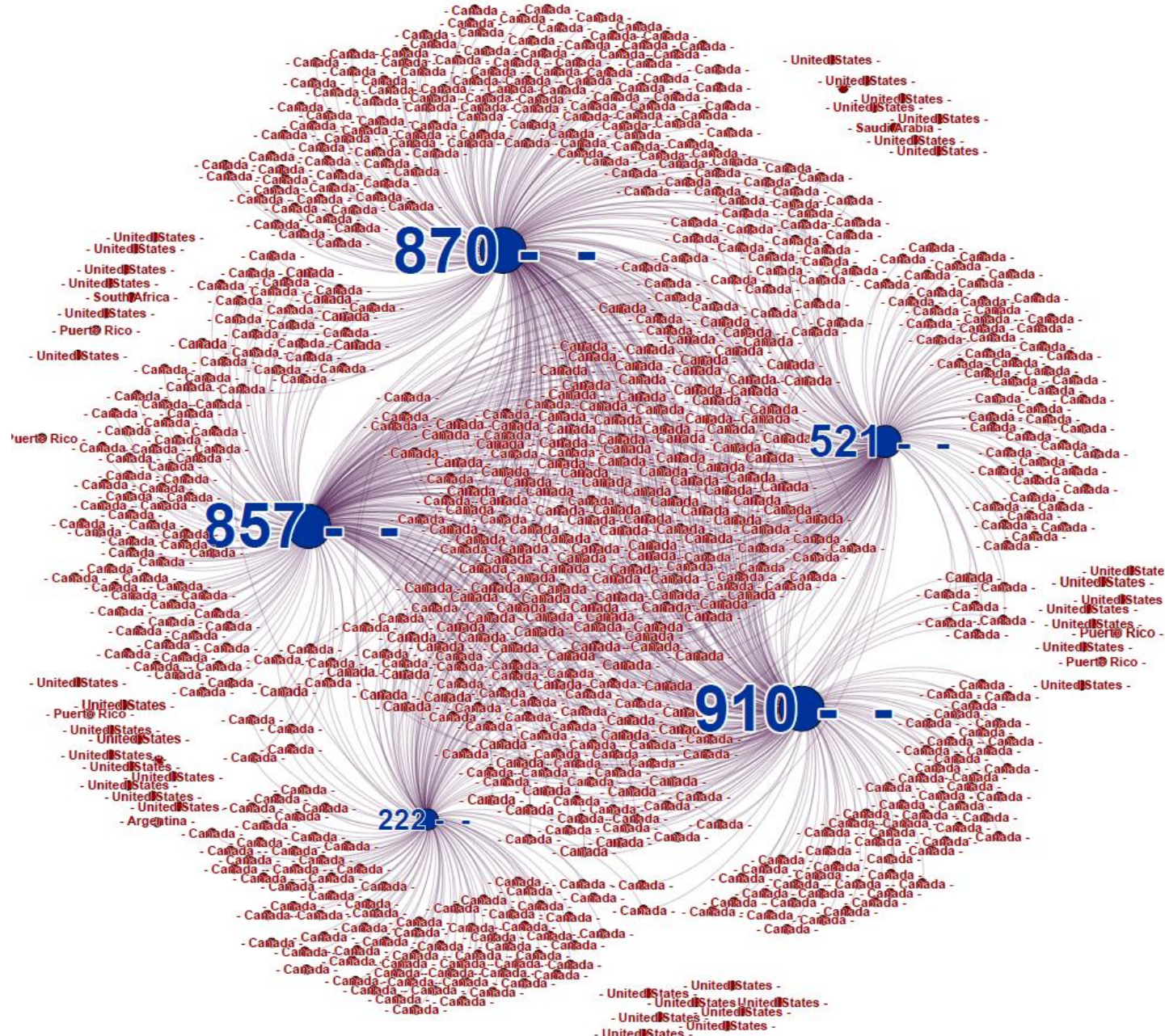












# Conclusions

31

- Spam:
  - ▣ Still a serious threat itself
  - ▣ An instrument for other cyber attacks
- A spam campaign detection, analysis and investigation system that:
  - ▣ Identifies spam campaigns (on-the-fly)
  - ▣ Aggregates different data sources to expose campaign characteristics
  - ▣ Labels spam campaigns
  - ▣ Scores spam campaigns
- Future work:
  - ▣ Improve spam campaign labeling (more suitable for spam contents)
  - ▣ Spam campaign categorization (threat-based)



# Acknowledgements

32

- Concordia University
  - ▣ Concordia Institute for Information Systems Engineering (CIISE)
- National Cyber-Forensics and Training Alliance (NCFTA)  
Canada
- Canadian Radio-television Telecommunications Commission (CRTC)