

URL Semantic Analysis for Phishing Detection

Presentation ResCom 2014
May 13, 2014

Samuel Marchal
Ph.D. student
samuel.marchal@uni.lu

SnT – Interdisciplinary Centre for Security, Reliability and Trust - Luxembourg
TELECOM Nancy – University of Lorraine - France

What is Phishing ?

- Use of **social engineering** and **technical subterfuges** to steal consumers' data:
 - Identity information
 - Web-sites credentials
 - Credit card information
 - Etc.
- Cause **billions** of dollars of loss every year
- SoA detection techniques
 - Reactive blacklisting (e.g. PhishTank)
 - Web-site graphical analysis
 - Passive DNS analysis
 - URL lexical analysis



Observation (phishing)

Source: APWG

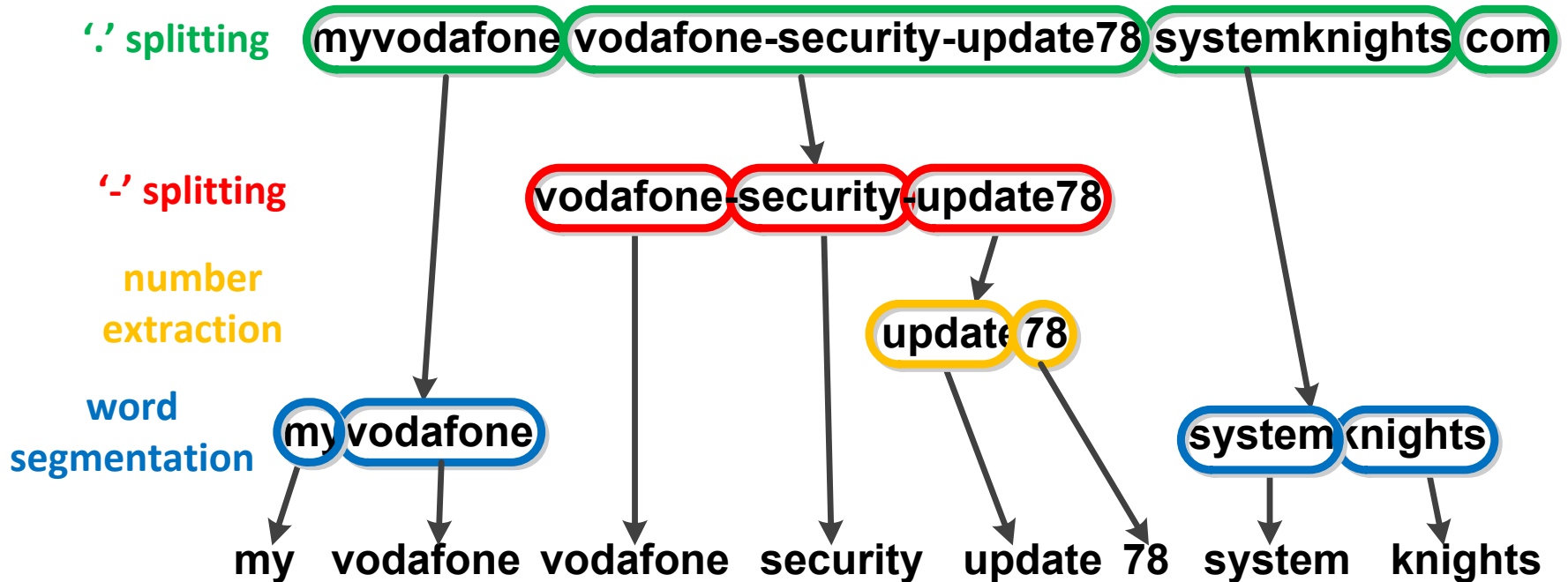
- Mainly relies on **social engineering** (easier)
- Use obfuscated URLs with similar patterns:
 - Targeted brands
 - Attractive keywords (secure, login, etc.)
- URLs are **meaningful**

Idea



**Perform a semantic analysis of URLs
To detect phishing**

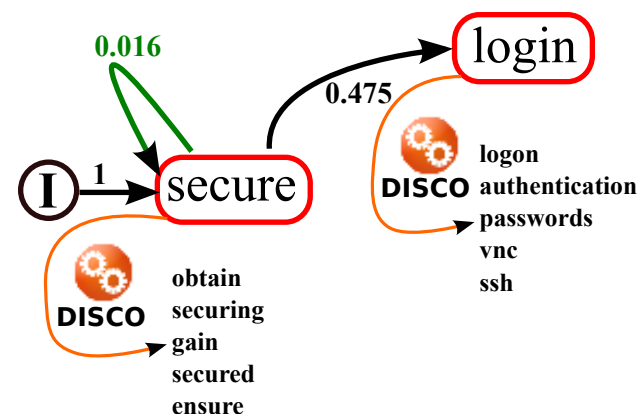
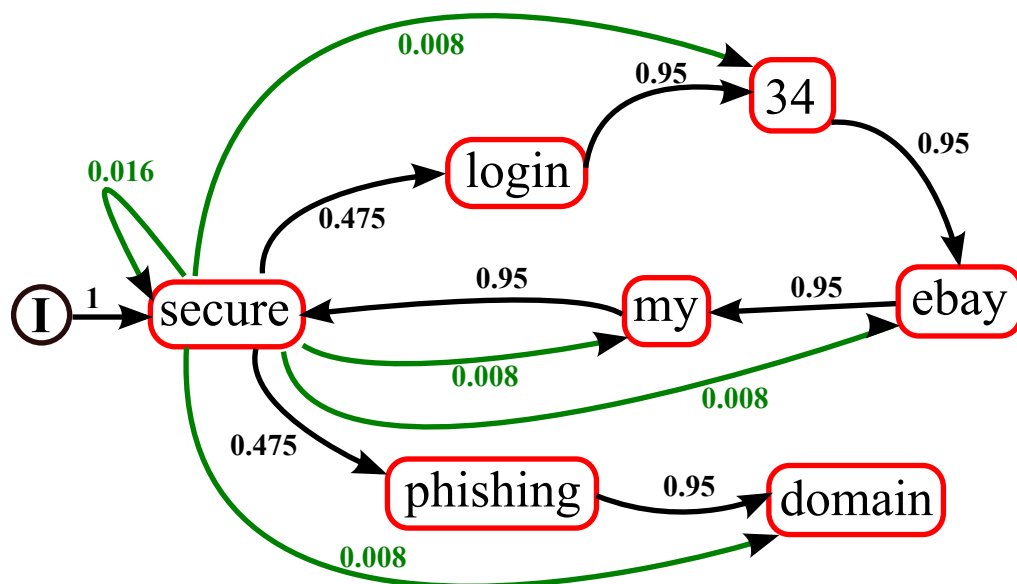
URL pre-processing



➔ Set of words composing the URL

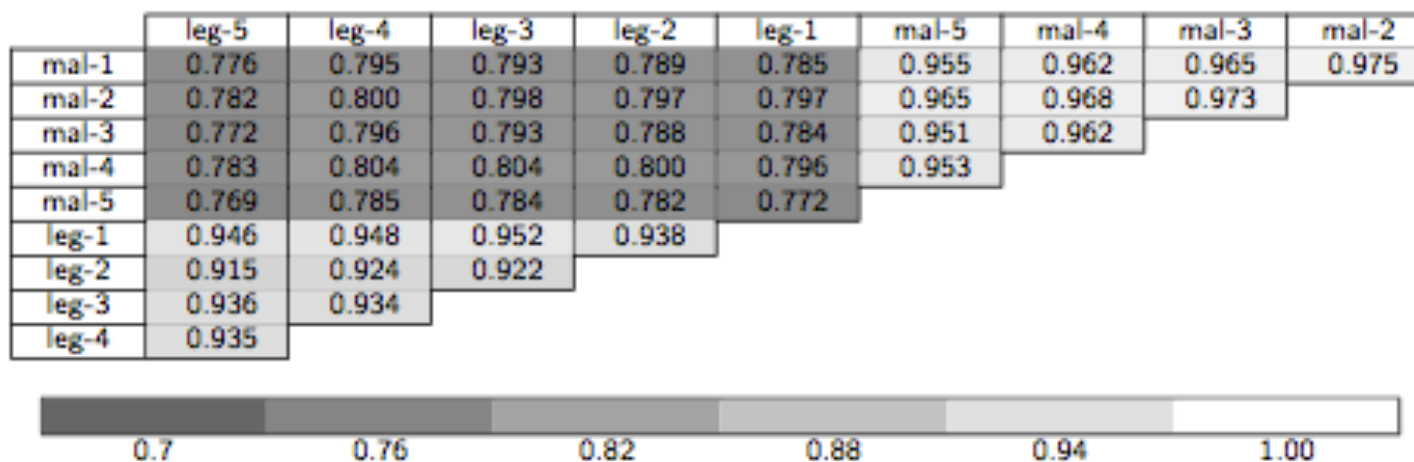
1. Phishing domains prediction (i.e. automatically generate domains likely to be registered by phishers)

- Study the composition of several known phishing domains
- Build a Markov chain representing their composition pattern
- Extend the generation model with semantic tools



2. Identify URL set maliciousness

- Group set of URL according to common features (e.g. pointing to a common @IP)
- Semantically compare with labeled sets of URL (legitimate / phishing)
- New metrics proposed to quantify semantic similarity between two sets of words



URL Semantic Analysis for Phishing Detection

Presentation ResCom 2014
May 13, 2014

Samuel Marchal
Ph.D. student
samuel.marchal@uni.lu

SnT – Interdisciplinary Centre for Security, Reliability and Trust - Luxembourg
TELECOM Nancy – University of Lorraine – France

Thank you