



# Click Trajectories: End-to-End Analysis of the Spam Value Chain

Presenter: Yu Shen

# Outline

- Background
- Model
- Data Collection
- Analysis
- Review Discussion

# Background

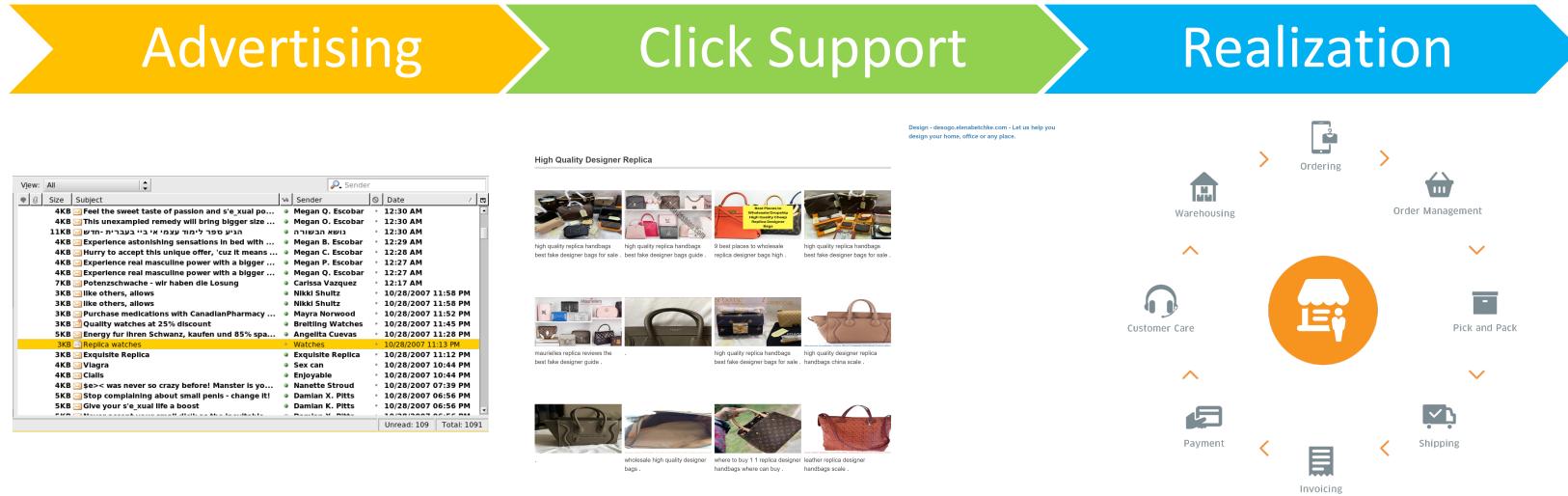


- Profitable industry – both spam and anti-spam
- Elements of the spam value chain are studied in isolation
- A decade ago (about 2000) spammers might have handled virtually all aspects of the business including email distribution, site design, hosting, payment processing, fulfillment, and customer service.
- Today's spam business involves a range of players and service providers.

# Modern Spam Working Model



## 3 stages – pipeline

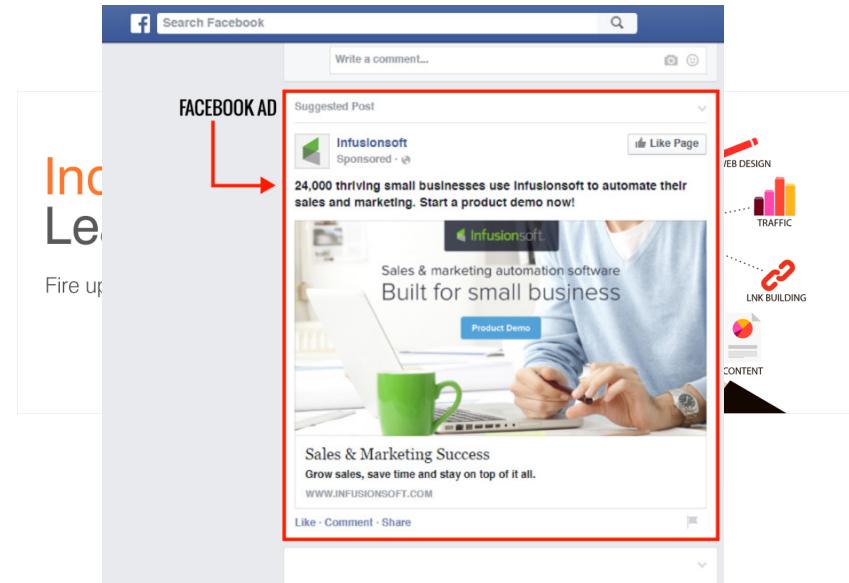
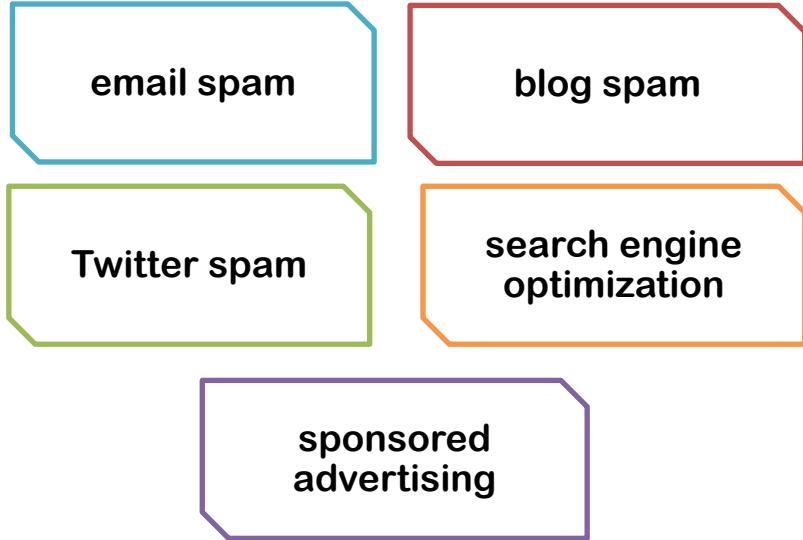


# Modern Spam Working Model - Advertising



TEXAS A&M UNIVERSITY  
Engineering

- ❖ All activities focused on reaching potential customers and enticing them into clicking on a particular URL



# Modern Spam Working Model - Advertising



TEXAS A&M UNIVERSITY  
Engineering

- ❖ Evolution of delivery and defense

open SMTP  
proxies

well-distributed  
IP blacklisting

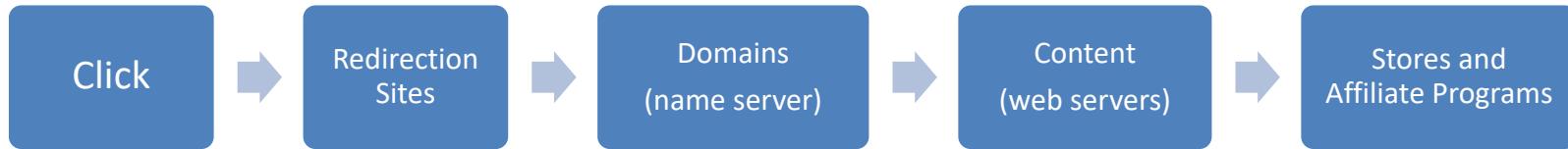
Botnets & Webmail spam  
& IP prefix hijacking

# Modern Spam Working Model – Click Support



TEXAS A&M UNIVERSITY  
Engineering

- Direct user's browser to a Web site of interest, after clicking on the embedded URL



- Redirection sites ( counter defensive measures )
  - Two ways: a legitimate third party | Manage DNS name resources by themselves

# Modern Spam Working Model – Realization



TEXAS A&M UNIVERSITY  
Engineering

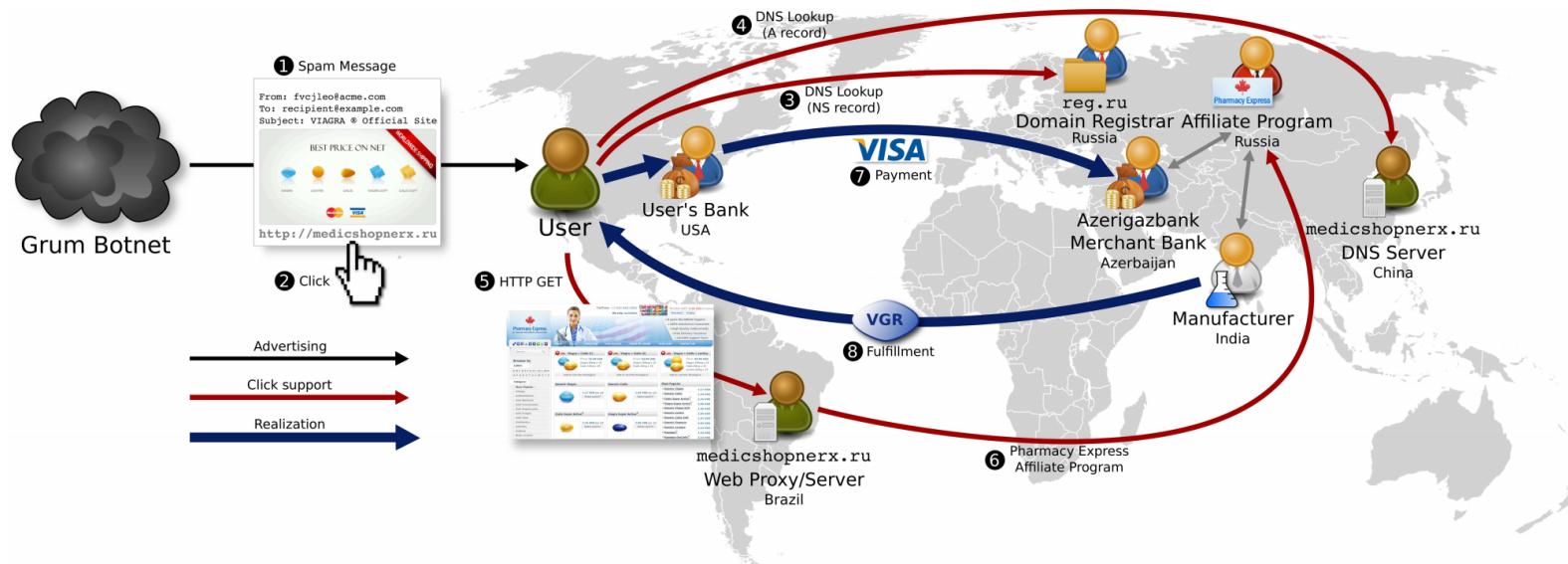
- Payment services



- Fulfillment



# Modern Spam Working Model – Example



# Data Collection - Spam-Advertised URLs



TEXAS A&M UNIVERSITY  
Engineering

- Duration: August 1, 2010 to October 31, 2010

Two sources:

**Third party**  
including multiple commercial  
anti-spam providers

**By botnet**  
from own botfarm environment

Feed Name	Feed Description	Received URLs	Distinct Domains
Feed A	MX honeypot	32,548,304	100,631
Feed B	Seeded honey accounts	73,614,895	35,506
Feed C	MX honeypot	451,603,575	1,315,292
Feed D	Seeded honey accounts	30,991,248	79,040
Feed X	MX honeypot	198,871,030	2,127,164
Feed Y	Human identified	10,733,231	1,051,211
Feed Z	MX honeypot	12,517,244	67,856
Cutwail	Bot	3,267,575	65
Grum	Bot	11,920,449	348
MegaD	Bot	1,221,253	4
Rustock	Bot	141,621,731	13,612,815
Other bots	Bot	7,768	4
<b>Total</b>		968,918,303	17,813,952

# Data Collection - Crawler data

Two parts: DNS & Web

DNS Crawler:

- Difficulty: fast flux - minimize central points of weakness
  - Solution: X. Hu, M. Knysz, and K. G. Shin. *RB-Seeker: Auto-detection of Redirection Botnets*. In Proc. of 16th NDSS, 2009.
- query servers repeatedly to enumerate the set of domains collectively used for click support

# Data Collection - Crawler data



TEXAS A&M UNIVERSITY  
Engineering

Web Crawler:

- replicates the experience of a user clicking on the URLs derived from the spam feeds
- captures any application-level redirects (HTML, JavaScript, Flash)
- Screenshots by Screengrab! Extension



- Firefox working parallelly fetching the web contents

# Data Collection - Crawler data

Result: nearly complete coverage: crawls over 98% of the URLs received

<i>Stage</i>	<i>Count</i>
Received URLs	968,918,303
Distinct URLs	93,185,779 (9.6%)
Distinct domains	17,813,952
Distinct domains crawled	3,495,627
URLs covered	950,716,776 (98.1%)

Notes: crawl URLs that account for only 20% of the whole

Explanation: 80% of domains, and the corresponding 2% URLs are from the domain-poisoning spam (Rustock) - no not reflect real sites

# Data Collection - Content Clustering & Tagging



TEXAS A&M UNIVERSITY  
Engineering

Focus on 3 types of products:

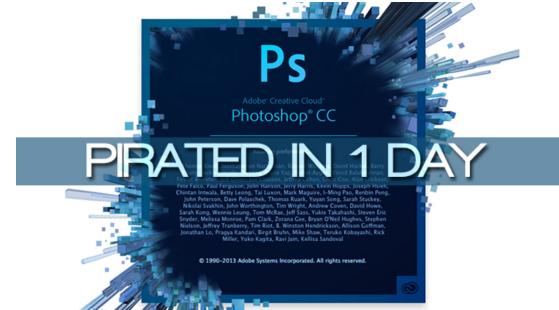
pharmaceuticals



replicas



software



Methodology: match sites with lexically similar content structure

# Data Collection - Content Clustering & Tagging



TEXAS A&M UNIVERSITY  
Engineering

## Step 1: Content clustering

- Basis: HTML text
- Approach: q-gram similarity
- Fingerprint: a set of multiple independent hash values over all 4-byte tokens of the HTML text

Q-gram language models:

**This is Big Data AI Book**

<i>Uni-Gram</i>	This	Is	Big	Data	AI	Book
<i>Bi-Gram</i>	This is	Is Big	Big Data	Data AI	AI Book	
<i>Tri-Gram</i>	This is Big		Is Big Data	Big Data AI	Data AI Book	

# Data Collection - Content Clustering & Tagging



TEXAS A&M UNIVERSITY  
Engineering

## Step 2: Category tagging

- Separates clusters that sells goods we want
- Avoid false negative: be intentionally conservative
  - potentially including clusters that turn out to be false positives to ensure that we include all clusters that fall into one of our categories
- identify interesting clusters by generic keywords
- Manual work: check for false negatives
  - Randomly select 1,000 pages from 675 untagged clusters, but find none missed

# Data Collection - Content Clustering & Tagging



TEXAS A&M UNIVERSITY  
Engineering

## Step 3: Program tagging

- Label clusters with specific program tags
  - E.g. distinct storefront brands of EvaPharmacy
- Methodology – vague...





# Data Collection - Content Clustering & Tagging

Affiliate Program		Distinct Domains	Received URLs	Feed Volume	RxRev	RX Rev Share	299	9,696	0.04%	
RxPrm	RX-Promotion	10,585	160,521,810	24.92%	Medi	MediTrust	24	6,156	0.01%	
Mailn	Mailien	14,444	69,961,207	23.49%	ClFr	Club-first	1,270	3,310	0.07%	
PhEx	Pharmacy Express	14,381	69,959,629	23.48%	CanPh	Canadian Pharmacy	133	1,392	0.03%	
EDEX	ED Express	63	1,578	0.01%	RxCsh	RXCash	22	287	<0.01%	
ZCashPh	ZedCash (Pharma)	6,976	42,282,943	14.54%	Staln	Stallion	2	80	<0.01%	
DrMax	Dr. Maxman	5,641	32,184,860	10.95%	<b>Total</b>	<b>Royal</b>	<b>Royal Software</b>	<b>572</b>	<b>2,291,571</b>	<b>0.79%</b>
Grow	Viagrow	382	5,210,668	1.68%		<b>EuSft</b>	<b>EuroSoft</b>	<b>1,161</b>	<b>694,810</b>	<b>0.48%</b>
USHC	US HealthCare	167	3,196,538	1.31%		<b>ASR</b>	<b>Auth. Soft. Resellers</b>	<b>4,117</b>	<b>65,918</b>	<b>0.61%</b>
MaxGm	MaxGentleman	672	1,144,703	0.41%		<b>OEM</b>	<b>OEM Soft Store</b>	<b>1,367</b>	<b>19,436</b>	<b>0.24%</b>
VgREX	VigREX	39	426,873	0.14%		<b>SftSl</b>	<b>Soft Sales</b>	<b>35</b>	<b>93</b>	<b>&lt;0.01%</b>
Stud	Stud Extreme	42	68,907	0.03%		<b>Total</b>	<b>7,252</b>	<b>3,071,828</b>	<b>2.12%</b>	
ManXt	ManXtenz	33	50,394	0.02%	<b>ZCashR</b>	<b>ZedCash (Replica)</b>	<b>6,984</b>	<b>13,243,513</b>	<b>4.56%</b>	
GlvMd	GlavMed	2,933	28,313,136	10.32%	<b>UltRp</b>	<b>Ultimate Replica</b>	<b>5,017</b>	<b>10,451,198</b>	<b>3.55%</b>	
OLPh	Online Pharmacy	2,894	17,226,271	5.16%	<b>Dstn</b>	<b>Distinction Replica</b>	<b>127</b>	<b>1,249,886</b>	<b>0.37%</b>	
Eva	EvaPharmacy	11,281	12,795,646	8.7%	<b>Exqst</b>	<b>Exquisite Replicas</b>	<b>128</b>	<b>620,642</b>	<b>0.22%</b>	
WldPh	World Pharmacy	691	10,412,850	3.55%	<b>DmdRp</b>	<b>Diamond Replicas</b>	<b>1,307</b>	<b>506,486</b>	<b>0.27%</b>	
PHOL	PH Online	101	2,971,368	0.96%	<b>Prge</b>	<b>Prestige Replicas</b>	<b>101</b>	<b>382,964</b>	<b>0.1%</b>	
Aptke	Swiss Apotheke	117	1,586,456	0.55%	<b>OneRp</b>	<b>One Replica</b>	<b>77</b>	<b>20,313</b>	<b>0.02%</b>	
HrbGr	HerbalGrowth	17	265,131	0.09%	<b>Luxry</b>	<b>Luxury Replica</b>	<b>25</b>	<b>8,279</b>	<b>0.01%</b>	
RxPnr	RX Partners	449	229,257	0.21%	<b>AffAc</b>	<b>Aff. Accessories</b>	<b>187</b>	<b>3,669</b>	<b>0.02%</b>	
Stmul	Stimul-cash	50	157,537	0.07%	<b>SwsRp</b>	<b>Swiss Rep. &amp; Co.</b>	<b>15</b>	<b>76</b>	<b>&lt;0.01%</b>	
Maxx	MAXX Extend	23	104,201	0.04%	<b>WchSh</b>	<b>WatchShop</b>	<b>546</b>	<b>2,086,891</b>	<b>0.17%</b>	
DrgRev	DrugRevenue	122	51,637	0.04%	<b>Total</b>	<b>7,530</b>	<b>15,330,404</b>	<b>4.73%</b>		
UltPh	Ultimate Pharmacy	12	44,126	0.02%	<b>Grand Total</b>	<b>69,002</b>	<b>365,395,278</b>	<b>100%</b>		
Green	Greenline	1,766	25,021	0.36%						
Vrlty	Virility	9	23,528	0.01%						

# Data Collection - Content Clustering & Tagging

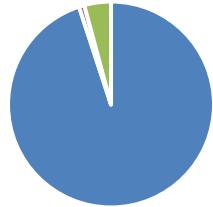


TEXAS A&M UNIVERSITY  
Engineering

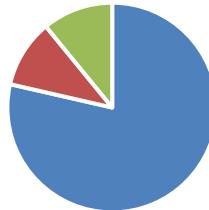
- Overall result

Stage	Pharmacy	Software	Replicas	Total
URLs	346,993,046	3,071,828	15,330,404	365,395,278
Domains	54,220	7,252	7,530	69,002
Web clusters	968	51	20	1,039
Programs	30	5	10	45

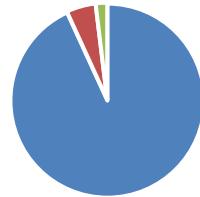
URLs



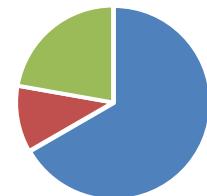
Domains



Web clusters



Programs



■ Pharmacy ■ Software ■ Replicas

# Data Collection - Purchasing

- Goal: identify differences or similarities in suppliers
- Basis: contents (e.g., lot numbers) and packaging (nominal sender, packaging type, etc.)
- Method: place multiple purchases from each major affiliate program or store “brand”, ordering the same “types” of product from different sites
- Process: attempt 120 purchases, 76 authorized, 56 settled

# Data Collection - Overview



TEXAS A&M UNIVERSITY  
Engineering

## 1 Feed Collection

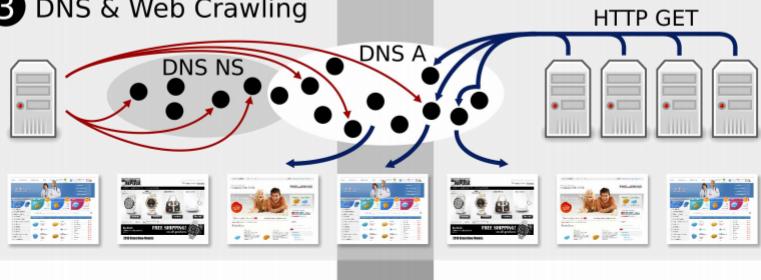


## 2 URL Extraction



<http://cheapdrugz.com>  
<http://pillsale.cn>

## 3 DNS & Web Crawling



## 4 Content Clustering



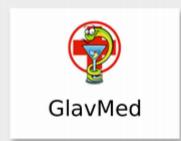
## 5 Content Tagging



Rx  
Promotion



Ultimate  
Replica



GlavMed

## 6 Selective Purchasing



# Analysis



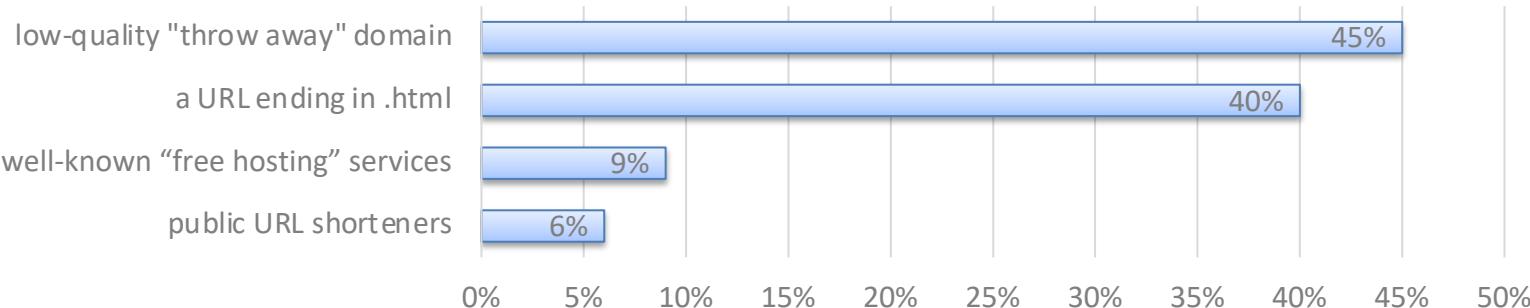
TEXAS A&M UNIVERSITY  
Engineering

- Major goal: identify any “bottlenecks” in the spam value chain



- [problem] Redirection in click support: use the final domain
  - URLs using redirection: 32%

Redirect Pattern Percentage



# Analysis



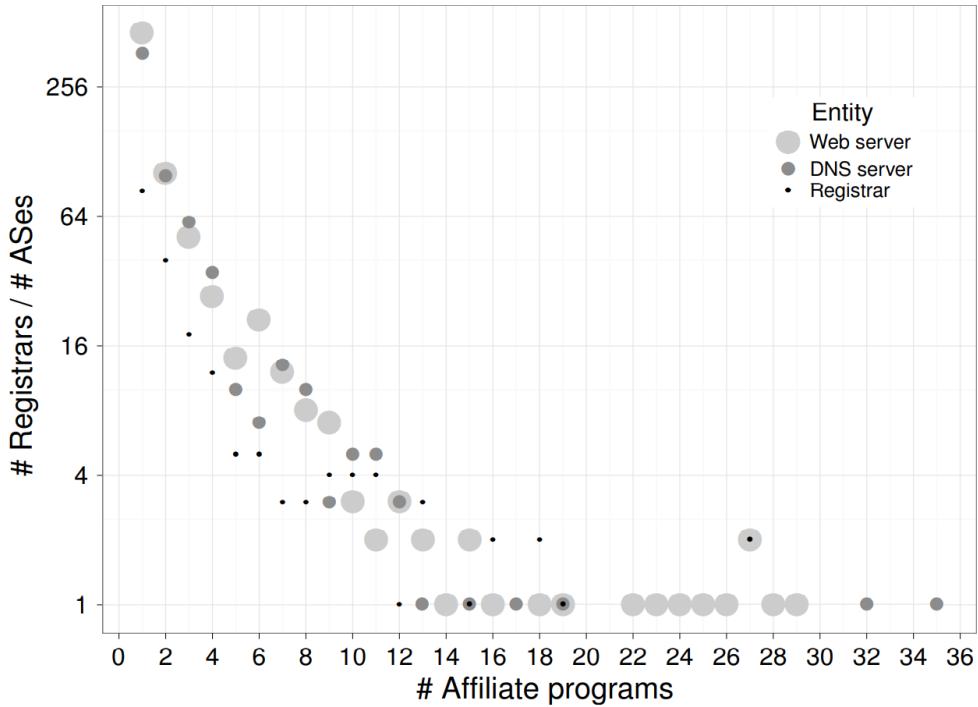
TEXAS A&M UNIVERSITY  
Engineering

- examine the set of resources used by sites for each affiliate program
- Network infrastructure sharing (at least 1):
  - Domain name
  - Registrar
  - Provision server

# Analysis - Click Support



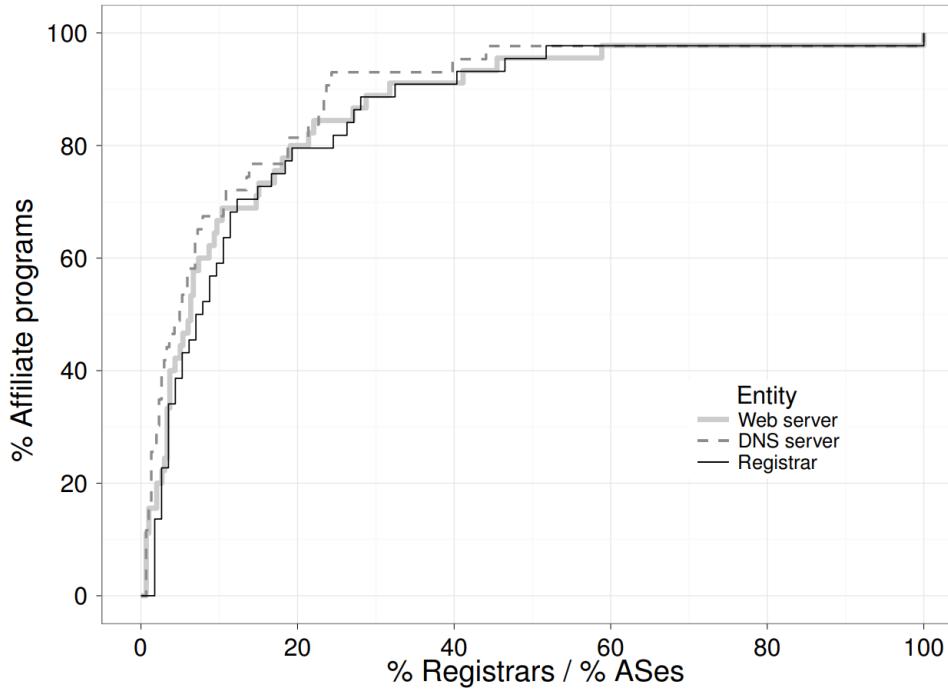
TEXAS A&M UNIVERSITY  
Engineering



Network infrastructure sharing among affiliate programs is concentrated in a small number of registrars and Autonomous Systems (ASes)

# Analysis - Click Support

Distribution of infrastructure among affiliate programs



# Analysis - Realization

## Hypothesis

- Realization infrastructure is the province of affiliate programs

## Fact

- Consistency in payment processing and fulfillment between different instances

## Support

- Email templates & Order number formats

# Analysis - Realization



TEXAS A&M UNIVERSITY  
Engineering

## Merchant banks

<i>Bank Name</i>	<i>BIN</i>	<i>Country</i>	<i>Affiliate Programs</i>
Azerigazbank	404610	Azerbaijan	GlvMd, RxPrm, PhEx, Stmul, RxPnr, WldPh
B&N	425175	Russia	ASR
B&S Card Service	490763	Germany	MaxGm
Borgun Hf	423262	Iceland	Trust
Canadian Imperial Bank of Commerce	452551	Canada	WldPh
Cartu Bank	478765	Georgia	DrgRev
DnB Nord (Pirma)	492175	Latvia	Eva, OLPh, USHC
Latvia Savings	490849	Latvia	EuSft, OEM, WchSh, Royal, SftSl
Latvijas Pasta Banka	489431	Latvia	SftSl
St. Kitts & Nevis Anguilla National Bank	427852	St. Kitts & Nevis	DmdRp, VgREX, Dstn, Luxry, SwsRp, OneRp
State Bank of Mauritius	474140	Mauritius	DrgRev
Visa Iceland	450744	Iceland	Staln
Wells Fargo	449215	USA	Green
Wirecard AG	424500	Germany	ClFr

# Analysis - Realization



TEXAS A&M UNIVERSITY  
Engineering

## Product Suppliers

<i>Supplier</i>	<i>Item</i>	<i>Origin</i>	<i>Affiliate Programs</i>
Aracoma Drug	Orange bottle of tablets (pharma)	WV, USA	CIFr
Combitic Global Caplet Pvt. Ltd.	Blister-packed tablets (pharma)	Delhi, India	GlvMd
M.K. Choudhary	Blister-packed tablets (pharma)	Thane, India	OLPh
PPW	Blister-packed tablets (pharma)	Chennai, India	PhEx, Stmul, Trust, CIFr
K. Sekar	Blister-packed tablets (pharma)	Villupuram, India	WldPh
Rhine Inc.	Blister-packed tablets (pharma)	Thane, India	RxPrm, DrgRev
Supreme Suppliers	Blister-packed tablets (pharma)	Mumbai, India	Eva
Chen Hua	Small white plastic bottles (herbal)	Jiangmen, China	Stud
Etech Media Ltd	Novelty-sized supplement (herbal)	Christchurch, NZ	Staln
Herbal Health Fulfillment Warehouse	White plastic bottle (herbal)	MA, USA	Eva
MK Sales	White plastic bottle (herbal)	WA, USA	GlvMd
Riverton, Utah shipper	White plastic bottle (herbal)	UT, USA	DrMax, Grow
Guo Zhonglei	Foam-wrapped replica watch	Baoding, China	Dstn, UltRp

# Analysis - Intervention analysis

From the point of defender: which intervention has the most impact?

- blocking its advertising
  - filtering spam
- disrupting its click support
  - takedowns for name servers of hosting sites
- interfering with the realization step
  - shutting down merchant accounts

# Analysis - Intervention analysis



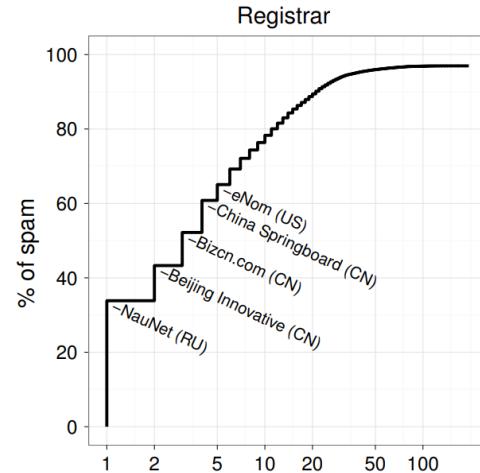
TEXAS A&M UNIVERSITY  
Engineering

- Strategy Evaluation
  1. overhead to implement
  2. Business impact (replacement cost & opportunity cost)

## ❖ Registrar

significant concentrations among the top few providers

takedowns would seem to be an effective strategy



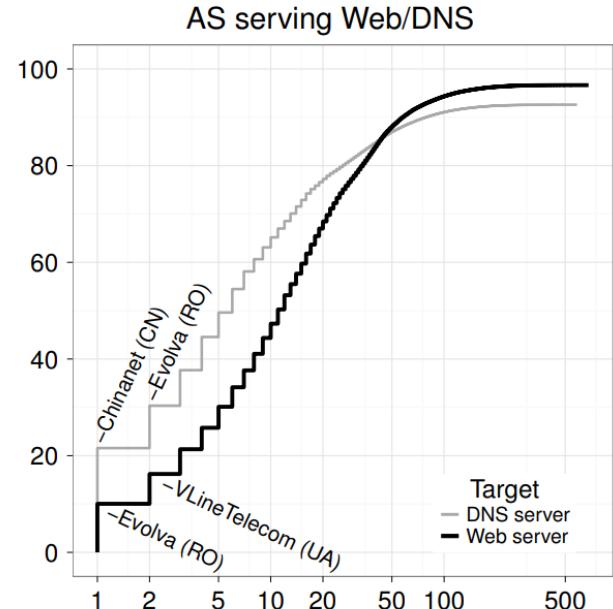
# Analysis - Intervention analysis

## ❖ Hosting

Analysis target: the number of distinct ASs that would need to be contacted

“worst case” model - that all such resources must be eliminated

Similar concentrations like registrar



# Analysis - Intervention analysis

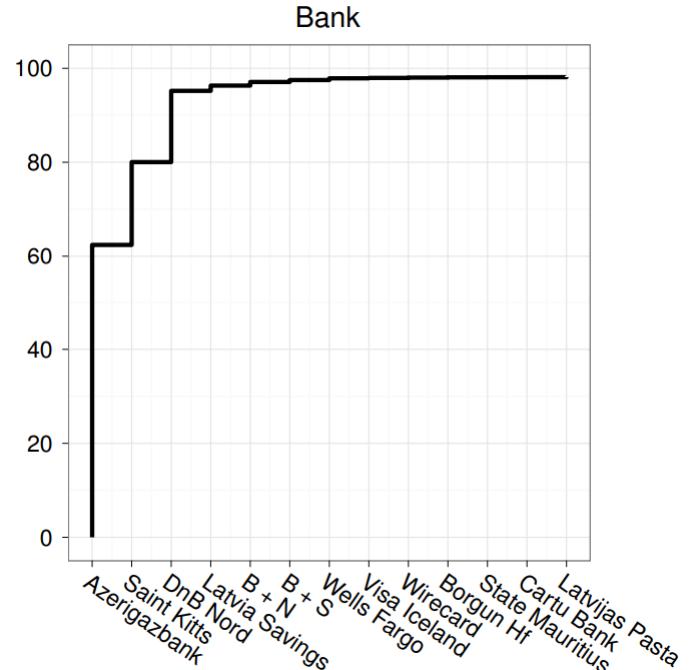
## ❖ Payments

Not buy from all domains

Assumption: bank use will be consistent across domains belonging to the same brand or affiliate program

95% transactions are with 3 banks:

- Azerigazbank
- Saint Kitts
- Latvia Savings



# Analysis - Intervention analysis



TEXAS A&M UNIVERSITY  
Engineering

- ❖ availability of alternatives & switching cost

	Alternatives	Switching Cost
Hosting	Vast (>1M)	Low (almost none)
Registrar	Few (roughly 900 gLTD registrars)	Low (spammers show great agility)
Bank	Few (a small part of the banks accept processing high-risk transactions)	High (setup fees, legitimate merchant account)

# Analysis – Policy Options



TEXAS A&M UNIVERSITY  
Engineering

- ❖ Directly pressure the merchant banks to stop doing business with such merchants
  - Slow
  - Incongruities in intellectual property protection
  
- ❖ refuse to settle certain transactions with the spam-supporting banks
  - Update quickly
  - Legal basis

# Review

## ❖ Merit - Good overview analysis

Reviewer1

- The paper provides a big picture of the spam ecosystem, which will clearly benefit future research in this direction.
- The authors divide the spam chain to advertising, click support and realization, providing a comprehensive understanding of the spam business.

Reviewer2

- Comprehensive thinkings. First full analysis on the entire value chain.

Reviewer3

- The paper brought me a full view of the spam ecosystem. Even if I'm not familiar with how spam works, I can still understand the paper to some degree.

# Review



TEXAS A&M UNIVERSITY  
Engineering

## ❖ Merit – Rich background info & graphics

Reviewer4

- This paper was very well written and thought out. One of the primary strengths is the detailed background information. ... Figure 2 really was set up very well to display step by step how the data collection and processing workflow occurred.

Reviewer5

- Specifically, the beginning sections discussing the background research and related work goes in depth about the step by step implementation of spam systems and how they try to target their customers.

Reviewer6

- The graphic well shows how much of a global industry this has become hardening attempts for one government to stop this industry.

## ❖ Weakness – Lack of detailed explanation on techniques

Reviewer3

- For some techniques mentioned in this paper like maximum likelihood methods and q-gram, I would like to see the more detailed examination of how they using it to process data.

Reviewer7

- The approach for content clustering is out-of-date. We have many more powerful tools today.
- In conclusion, they mentioned that they "described a framework for ...". But the framework is not clear in the paper.

Reviewer8

- When classifying clusters, if using image recognition to categorize three kinds of web sites, personally assumed, would be better than HTML contents analysis. At least untagged groups would be fewer.

## ❖ Different Views: Is payment the bottleneck? / Is the policy practical?

Reviewer1

- Some of the proposed intervening measures are hard to realize. For example, one approach for intervening at the payment tier of the value chain suggested in the paper is to directly engage the merchant banks and pressure them to stop doing business with such merchants. This measure needs support from the government. It is much more difficult to realize than designing the anti-spam system.

Reviewer8

- The conclusion, "payment" has few connections with prior research approaches, classifying, content analysis. This makes the whole paper a bit waste of efforts.

Reviewer6

- With the use of a VISA card issuer that is willing to help provide non traditional amount of information to the study I would have liked for the authors to ask the issuer directly how they feel about any of the policy recommendations.

# Review



TEXAS A&M UNIVERSITY  
Engineering

## ❖ Other questions:

Reviewer1

- The data collection extracts URL. However, some spammers may not use the URL. They may only use contact information like a phone number.

Reviewer7

- The paper is not well-organized. From their analysis, we can get gradually understand the spam value chain. They illustrated each stage very clearly. But they missed an overview of the chain. It would be better if they can spend some words on the overview of the chain.



TEXAS A&M UNIVERSITY  
**Engineering**