

What we now know about phishing websites

Richard Clayton

(joint work with Tyler Moore)



**UNIVERSITY OF
CAMBRIDGE**
Computer Laboratory

Luxembourg
25th March 2009

Academics & phishing

- Everyone can play! Display instant expertise!!
 - examine psychology, attempt to block spam, detection of websites, browser enhancements, password mangling, reputation systems etc
- Our approach : Security Economics
 - phishing will continue, because humans involved!
 - so we measure the impact, assess the effectiveness of countermeasures, work out how to change incentives so that problem tends to fix itself...

Flaws when researching reality

- Hard to report on an on-going understanding
 - papers have to be “novel research”
 - PhDs have to be “a contribution”
 - hence where the “real world” is tackled, tendency to pick the “low hanging fruit” and move on
- Errors in early papers often go uncorrected
 - “peer review” process needs knowledgeable peers
 - natural tendency not to want to report failures
 - natural tendency not to admit mistakes

Types of phishing website (Jan 2008)

- Misleading domain name (unusual at present)
 - `http://www.banckname.com/`
 - `http://www.bankname.xtrasecuresite.com/`
- Insecure end user or machine (76% of sites)
 - `http://www.example.com/~user/www.bankname.com/`
 - `http://www.example.com/bankname/login/`
- Free web hosting (17% of sites)
 - `http://www.bank.com.freespacesitename.com/`
- Random domains (after canonicalisation)
 - rock-phish 4%, fast-flux 1.4%, “ark” 1.4%

Rock-phish is different!

- Compromised machines run a proxy
- Domains do not infringe trademarks
 - name servers usually done in similar style
- Distinctive URL style

`http://session9999.bank.com.lof80.info/signon/`

- We track domains & IP addresses generically
- Usage of “fast-flux” from Feb’07 onwards
 - viz: resolving to 5 (or 10...) IP addresses at once

Phishing website lifetimes (hrs) 2007	# sites (8 weeks)	Mean lifetime	Median lifetime
Non-rock	1707	58.4	20
Rock-phish domains	419	94.3	55
Rock-phish IP addresses	122	124.9	25
Fast-flux rock-phish domains	67	454.4	202
Fast-flux rock-phish IP addresses	2995	124.6	20

How many visitors?

- Some (non rock-phish) sites had world readable “webalizer” statistics pages
 - could determine number of visitors on each day
 - 22 on day first reported, 24 next day and then tails off a bit (but NOT to zero)
- Some sites had world readable files of compromised credentials
 - about 50% were “die spammer die” responses

What's the co\$t of phishing?

- 56 days, 1448 banking websites (exclude eBay)
- Average lifetime was 57 hours
- Hence 33 real victims per site
- Gartner loss estimate of \$572/victim
- Hence \$178 million per year
- Rock-phish is half the spam... so \$350 million
 - NB: complete hand-waving !!!
 - and cf. Gartner total estimate of \$2 billion

Data Sources

- Originally mining PhishTank dataset
 - free and apparently accurate and substantial
- Now getting data from a brand owner and two brand protection companies (plus PhishTank and “Artists Against 419”)
 - PhishTank only has 48% of sites we know of
- Even the commercial “feeds” have common components, but turn out to be different...

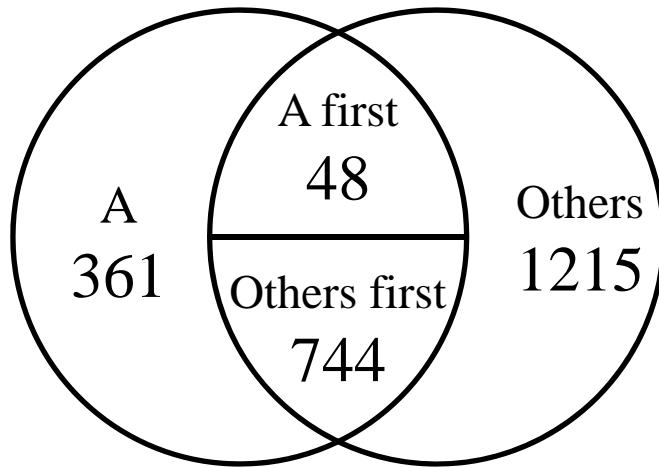
	PhishTank	BrandProtectA
URLs	10924	13318
Non-duplicate URLs	8296	8730
Unique URLs	3019	2585
Rock-phish domains	586	1003
Unique rock-phish domains	127	544

63% of total overlap (9380 URLs) from “PhishReporter”
remainder from 316 separate submitters

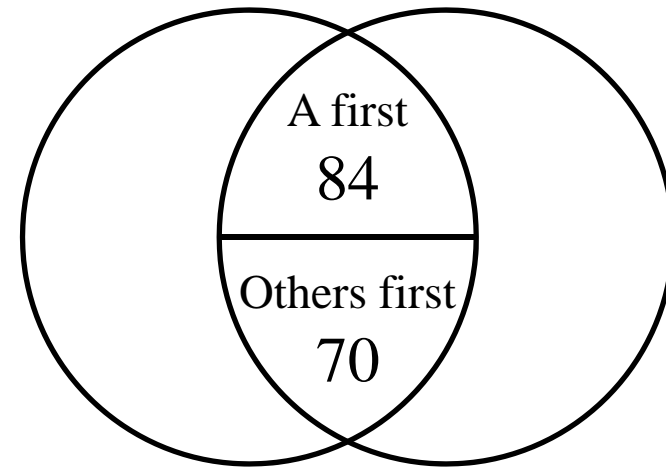
Verification time (average)	46 hours	8 seconds
Verification time (median)	15 hours	8 seconds

Feeds are not shared

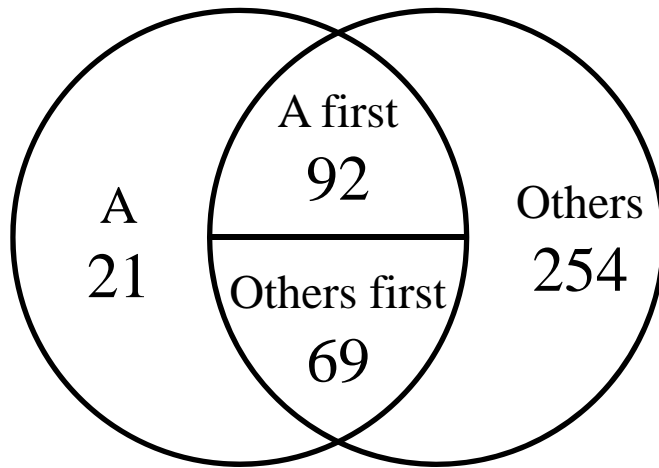
- Brand-protection companies obtain feeds from many places (including PhishTank)
- They run their own detectors
- They sell feeds, but don't share them
- Hence Company A, who sells services to Bank A1, can be unaware of sites detected by Company B – and doesn't take them down



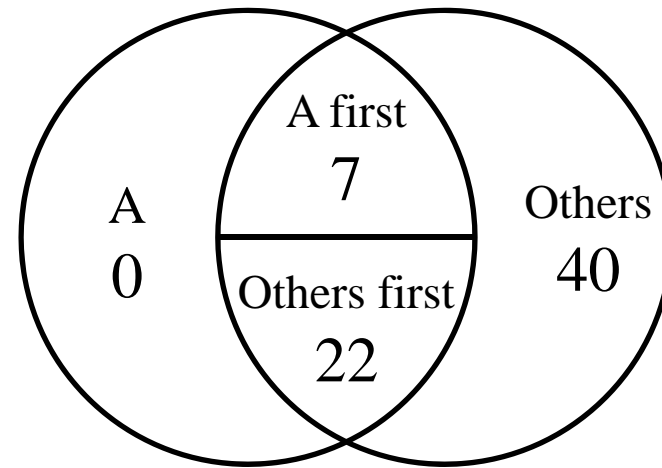
Ordinary phishing sites



Delay in detecting (hours)



Mean lifetime (hours)



Median lifetime (hours)

Bank A1's experience as a client of BrandProtection company A

Company A v Company B

- Same pattern continues for top 6 banks for Company A and B, and for all n clients
- However, less pronounced for B: which seems to have a better feed [or maybe just one that is much more aligned with ours!]
- But A's clients bigger and proportion missed goes up with size; so B's prowess may be more a structural issue than just extra effectiveness

This represents risk

- Longer lifetimes => more visitors (Webalizer logs)
- Hence we can assess impact of longer lifetimes:

Exposure figures (6 month totals)	A's banks		B's banks	
	Khour	\$m	Khour	\$m
Actual values	1005	276	78	32
Expected if sharing	418	113	61	28.5
Effect of no sharing	587	163	17	3.5

Hence...

- Banks should force brand-protection companies to share feeds
 - cf the anti-virus community since 1993
- Brand-protection companies could form a “club” to prevent new entrants from free-riding
 - don’t have to make feeds free, just share them
- Side-note: free-riding by rock-phish attacked banks only works some of the time!

How are insecure machines found?

- Traditionally machines found by “scanning” hence interest in Intrusion Detection Systems, “slow scan” software etc etc
- We have been collecting Webalizer logs (wanted to count number of visitors to sites and hence calculate impact of prompt take-down)
- Webalizer parses referrer strings to determine search terms used to locate the sites....

Typical searches in weblogs

- Hand categorisation, but most were obvious
 - many searches for MP3s ! these were ignored
- Vulnerability
 - `phpizabi v0.848b c1 hfp1` (CVE-2008-0805)
- Compromise
 - `allintitle:welcome paypal`
- Shell
 - `c99shell drwxrwx`

Webalizer logs (June 07 – March 08)

- 2486 domains with world-readable logs
- 1320 (53%) had one or more search terms
- 25 cases where searches provably linked

	Domains	Phrases	Visits
Any evil search	204	456	1207
Vulnerability search	126	206	582
Compromise search	56	99	265
Shell search	47	151	360

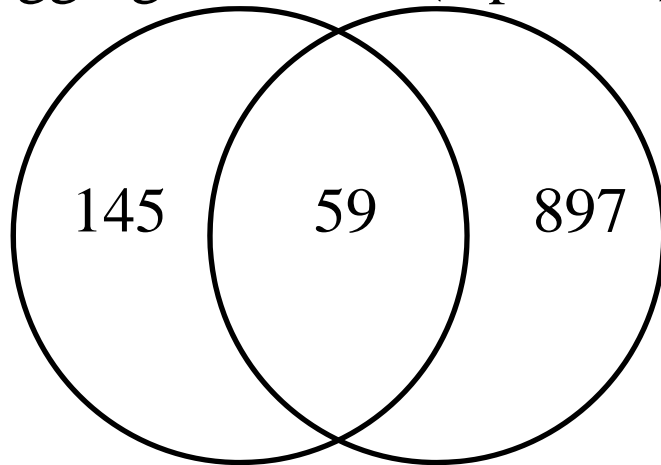
More statistics

- Assume Webalizer sites are a random sample of all sites (make up your own mind on that)
 - if so, then 95% confidence interval for incidence of “evil searching” (aka “dorks”) is 15.3% to 19.8%
- Did our own searches (thanks Yahoo!) on evil and non-evil terms and checked if phishing site
 - 1.9% sites found with evil terms used for phishing
 - 0.73% sites located by using non-evil terms (statistically significant difference)

Overlap of search results

Webalizer
logging data

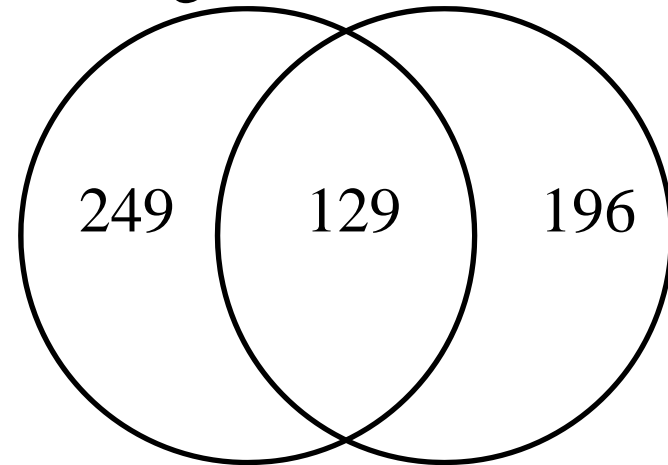
Yahoo!/Google
(April 08)



Many searches don't work any more, but lots more sites to attack!

Google

Yahoo!



There's a surprising lack of overlap in the results

Recompromise

- Consider phishing pages on same site more than a week apart (likely a different attacker)
- 9% of all sites recompromised within 4 weeks, rising to 19% within 24 weeks
- For Webalizer sites this is 15% rising to 33%
- If evil search terms present then this becomes 19% rising to 48% (14% to 29% if no terms)
- This doubling is statistically significant!

Email spam

- Email drives visitors to phishing websites
- Combining data about when URLs are seen with website lifetimes gives better picture

	Phishing feeds		Spam feed	
	Total	Visited	Total	Visited
Ordinary	4250	3360	432	369
Fastflux	120	113	103	100

How to measure the harm ?

	Number of websites	Spam volume	Website lifetime (hours)
Ordinary	4250 (97%)	31.7%	25600 (72.6%)
Fast-flux	120 (3%)	68.3%	9674 (27.4%)

Losses to customers correlate to volume of spam (if it is convincing), but rapid removal of websites mitigates impact. Number of websites (and volume of spam) impacts public perceptions, possibly eroding trust.

Comparing take-down times

- Defamation – believed to be quick (days)
- Copyright violation – also prompt(ish)
 - experimentally “days” (albeit with prompting)
- Fake escrow agents
 - average 9 days, median 1 day
 - note that AA419 aware of around 25% of sites
- Mule recruitment sites (Sydney Car Center etc)
 - average 13 days, median 8 days

Phishing Lifetimes (hrs)	sites	mean	median
<i>Free-web hosting</i>			
all	395	47.6	0
brand-owner aware	240	4.3	0
brand-owner unaware	155	114.7	29
<i>Compromised machines</i>			
all	193	49.2	0
brand-owner aware	105	3.5	0
brand-owner unaware	155	103.8	10
<i>Rock-phish domains</i>	821	70.3	33
<i>Fast-flux domains</i>	315	96.1	25.5

Incentives

- Most of the take-down time variations are explainable in terms of incentives
 - the motivated complain again&again until removed
 - the banks are ignoring mule recruitment (not their problem) so just volunteers (vigilantes)
 - escrow faster than mule sites: attacking the innocent?
or maybe escrow.com is doing more than we think?
 - no-one's job to remove fake pharmacies (and no active volunteers) so their lifetime is ~2 months

Child Sexual Abuse Images (“CAI”)

- Provided with anonymised data by IWF
- Jan–Dec 2007 2585 domains
 - ignoring 8 (free-web?) domains with >100 reports
- Computed the initial take-down time (ignored recommitment): mean 21 days, median 11 days
- If we include sites with no removal at all then mean grows to 30 days (and counting)
 - median also grows by one day

Why so slow?

- In fact quick within the UK : IWF checks with police and then contacts the ISP
- But “not authorised” to act internationally
- Passes data via UK police to foreign forces
 - but may not reach local field office for a while
- Also pass to another INHOPE member
 - but (eg) NCMEC only act “when appropriate”
- Confusion of aims (removal/catch criminals)

Ongoing research agenda

- How many phishers are there ?
- How much phishing is phishing ?
- How do we fix the incentives to prevent phishing from being effective ?
- Phishing is now mechanised and uses standard kits – we'd like to disrupt them!
- Phishing attacks also involve spam: the timing of this is as relevant as site take-down times

What we are (currently) sure about

- The phishing site take-down industry is putting significant funds at risk by not co-operating
- The police are chasing the right gang!
- Search engines are widely used to find websites to compromise (and re-compromise)
- Takedown times are affected more by incentives than by formal structures
- Slowness of removal of CAI is a scandal

What we now know about phishing websites

BLOG: <http://www.lightbluetouchpaper.org/>

<http://www.cl.cam.ac.uk/~rnc1/>

<http://people.seas.harvard.edu/~tmoore/>

PAPERS: <http://www.cl.cam.ac.uk/~rnc1/publications.html>