

**Jeremiah O'Connor
Thibault Reuille
Vinny LaRiza
Artsiom Holub
OpenDNS
March 2016**

The Security Wolf of Wall Street: Fighting Crime with High-Frequency Classification and Natural Language Processing

Introduction

When I first arrived at OpenDNS there was a plethora of awesome botnet, DGA, and time-series detection algorithms however there was not a lot of research around phishing detection, which made sense because when you think of phishing you think of E-mail, not necessarily DNS. This necessity to add a phishing detection algorithm to our arsenal is the primary motivation behind NLPRank, the model we will discuss in this paper. Further inspiration and training sets for the algorithm came from the awesome, homegrown, community-based phishing verification system Phishtank. Typically when an analyst performs hit review they view the site in question in a Tor Browser, their mind gains summary of page by looking at it, and they make decision whether it's malicious or not. The goal of NLPRank is to automate this process and detect dedicated and compromised phishing domains in our authoritative DNS log stream using unsupervised learning/topic modeling techniques in real-time fashion. Essentially we are trying to create an automated analyst. In this paper we will first give some background on phishing as an overall attack (Artsiom and Vinny), next we will give an overview of Avalanche, our streaming data processing pipeline (Thibault), and finally we will discuss the detection model NLPRank (Jeremiah).

Overview of Phishing Attacks

Phishing is the attempt to acquire sensitive information such as usernames, passwords, and credit card details (and sometimes, indirectly, money), often for malicious reasons, by masquerading as a

trustworthy entity in an electronic communication. While a phish can be approached from many different platforms (i.e. text message, phone call or email), we will review most common campaigns.

1. Phishing e-mail campaigns

E-mail can be a powerful persuasion device for attackers and con artists alike. It has become a basic mode of communication for many people and is considered crucial for many companies to run a successful business. People have grown so accustomed to e-mail that they rarely question the integrity of an e-mail's source or content.

Phishing scams are a form of cybercrime that involves defrauding users by acting as legitimate companies or organizations in order to obtain sensitive information such as passwords and login credentials.

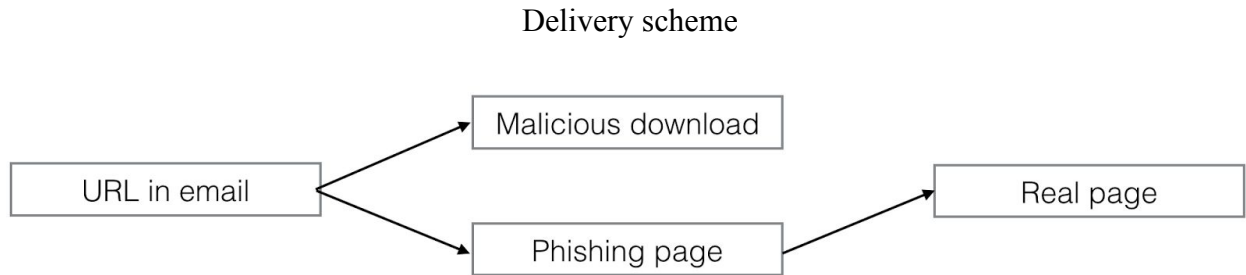
This brings to light the concept of Email Spoofing. E-mail spoofing is the forgery of an e-mail header so that the message appears to have originated from someone or somewhere other than the actual source. Because the SMTP provides little in the way of authentication or integrity checking, anyone with the requisite knowledge can connect to the server and use it to send messages. To send spoofed e-mail, senders insert commands in headers that will alter message information. It is possible to send a message that appears to be from anyone, anywhere, saying whatever the sender wants it to say. Thus, someone could send spoofed e-mail that appears to be from you with a message that you didn't write.

Another common attribute of phishing emails is the use of Link Manipulation. Link Manipulation is a fairly easy thing to accomplish. Essentially every link that's posted in an email can be changed to redirect to wherever the attacker wishes. An attacker can make it look like the link is going to the homepage of a legitimate company and redirect the victim to a malicious site they have set up.

While a lot of phishes tend to be general and not personal, there are some that are. Spear phishing, for instance, is an email or electronic communications scam targeted towards a specific individual, organization or business. These emails will target the victim personally, often by posing as a trusted individual or company.

The Cloning method is another attribute that's commonly found in the more high-quality phishes. Cloning is a type of phishing attack whereby a legitimate, and previously delivered, email containing an attachment or link has had its content and recipient address (or addresses) taken and used

to create an almost identical or cloned email. The attachment or Link within the email is replaced with a malicious version and then sent from an email address spoofed to appear to come from the original sender.



Example of phishing emails:

BankofAmerica campaign

From: "Bank of America" customerservice@bankofamerica.com
To: "Jane Smith" jane-smith12@gmail.com
Date: Wed, May 26, 2010
Subject: Fraud Alert – Action Required



Dear Customer,


At Bank of America, your satisfaction is our number one priority. We have recently added an Advanced Online Security option for our customers with online accounts. It is urgent that you go to our website and add Advanced Online Security to your account. Click on the following and update your information www.bankofamerica.com.

If you do not take these steps, in order to protect you, we will put a hold on your account, and you will be required to visit your local branch to verify your identity.

Thank you for helping us to make Bank of America the safest bank on the internet.

If you are receiving this message and you are not enrolled in online banking, [sign up now](#). New online members will automatically be enrolled in the Advanced Online Security program.

Sincerely,

Bank of America Online Security Department 

IRS campaign

From: IRS Online <ahr@irxt.com>
Reply-To: "noreply@irxt.com" <noreply@irxt.com>
Date: Thursday, April 11, 2013 12:15 PM
Subject: Final reminder: Notice of Tax Return. ID: I3H583326/13



Department of the Treasury
Internal Revenue Service

04/11/2013

Reference: I3H583326/13

Claim Your Tax Refund Online

Dear Taxpayer,

We identified an error in the calculation of your tax from the last payment, amounting to \$ 319.95.

In order for us to return the excess payment, you need to create a e-Refund account after which the funds will be credited to your specified bank account.

Please click "Get Started" below to claim your refund:

[Get Started](#)

Below are the top ten recent phishing alerts reported by Fraudwatch International.

[Bank of America – Bank of America – Important Notice](#)

[Westpac Bank – Your Account Has Been Blocked](#)

[PayPal – Resolve remote access](#)

[Chase Bank – INFORMATION ABOUT YOUR ACCOUNT](#)

[Outlook – FW: \(WARNING\) Microsoft Account Suspension](#)

[Chase Bank – Fwd: Online Message From Chase Online](#)

[Apple Store – About your last Transaction](#)

[PayPal – We were unable to process your request](#)

[Bank of America – Fwd: We Have Temporaly disabled Your Account!!!!!!](#)

[PayPal – \[PayPal\] Alert from our security seystem !](#)

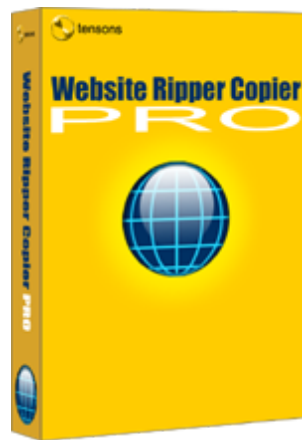
2. Phishing web pages

Phishing web pages can generally be separated into 2 categories, 'compromised' and 'dedicated.' A compromised website is a site that wasn't created for the purposes of phishing, but was later "hacked," or compromised, with the aim of installing a phish page on the site without the owner's knowledge. Dedicated phishing sites are created specifically for the purpose of being used as a phish, and can often be a play off of the name of the company itself (i.e. a Facebook phish url might be something like facebook.cn). The PhishTank.com website, which is a website that allows anybody to submit phishing urls to the site which then allows the community there to vote and decide if the site is a phish or not, handles both compromised and dedicated phishes, however sees many more compromised sites on average.

2.1. Dedicated phishing web pages (Typosquatting)


Typosquatting is a well-known security problem. In a typosquatting campaign, a malicious actor will target one or more well-known websites or brands and register domains very similar to the legitimate domain. Typosquatting easily solves one of the biggest hurdles for these bad actors: delivery of the malicious content. In typosquatting, users just show up. In some cases, effects can be relatively mild, such as: the user is redirected to objectionable material; the user is presented items for purchase from storefronts of questionable repute; or the user sees content that unfavorably portrays the intended brand or site. Effects can also be much worse. APWG reports that in the fourth quarter of 2015, 20,320 phishing websites were detected. Compared with the 67,765 unique phishing email campaigns reported during that period, phishing websites present a lower, while still significant, risk. The malicious actor can spoof a real site to harvest login credentials, place backdoors on a system, install ransomware, or really anything else of his choosing.

Usually created as web clones utilizing software : Website Ripper Copier, NCollector Studio, WebSiteSniffer, WebCopier, SiteSucker, WebCopy, Wget, etc.



SiteSucker

and services like SiteCloner, clonezone, copyanywebsite, etc. :



SiteCloner is a PHP script which lets you make static **clones of any website**. Wether the website was built using plain HTML or a CMS like Wordpress, Joomla, Drupal or anything else.

Leave a message

http:// Enter the URL of the Website you would like to clone

[Main](#) [My Clones](#) [Gallery](#) [Log in](#)

clone zone

TWEAK THE WORLD'S MOST POPULAR WEBSITES

GET STARTED NOW

Depending on what software or service were used different breadcrumbs can be found, but the main goal is create identical page.

Most popular methods attackers use to fake domain names:

1. Registering domain in different domain zones, for example:

wellsomebank.com

wellsomebank.us

wellsomebank.net

wellsomebank.om

wellsomebank.org

1. Registering domain using homoglyphs - graphically identical or similar to each other symbols, so the address bar that contains homoglyph will look like legitimate page, for example:

weIIIs0mebank.com – l to capital I, o to 0

we11sornebank.com – l to 1, m to rn

wellsomedank.com – b to d

The homoglyph table can look like this:

l	1
o	0
i	j
m	rn
q	g
d	b

1. Registering subdomain with similar name or domain where [.] (between original domain and subdomain) replaced with [-], for example:

some.bank.com – original domain

some-bank.com – phishing domain

1. Registering domain using duplicated symbols, for example:

ssomebank.com, soomebank.com,

sommebank.com, someebank.com

somebbank.com, somebaank.com,

somebannk.com, somebankk.com

1. Registering domain with changed arrangement of the letters, for example:

smoebank.com, osmebank.com, sombeank.com,

soembank.com, sombaenk.com, someabnk.com

Dedicated phishing sites usually hosted using free, rogue or bulletproof hosting providers:

000webhost is the biggest free hosting provider on web. Free web hosting with PHP and MySQL. Domain netne.net is used by our clients to host their websites. Each subdomain on *.netne.net is managed by a different customer. And with data from investigate we detect over 500 domains associated with malware and phishing threats.

? INVESTIGATE

Visualize

Constrain Reg

Showing 500 results for **.*\netne\.net**

Search maxed out at 500 results. For complete results please narrow your search.

Domain Name	Security Categories
voktakte.netne.net	Malware
netease126-163login-logout.netne.net	Malware
daynasour.netne.net	Malware
facebookmobile.netne.net	Malware
sgfsecure.netne.net	Malware
paypal-account.netne.net	Malware
safranet.netne.net	Phishing, Malware

Another highly abused registrant let you register .om domains.

MARCARIA.com
SM NETWORK

Login to your account:

OK

View Cart

New User

Recover Password

Home Page

About Us

Domain Services

Other Services

Search

Contact Us

English

MARCARIA.com

.om

Domain Tools

Home > Domain Name Registration > Search Results

Search Results

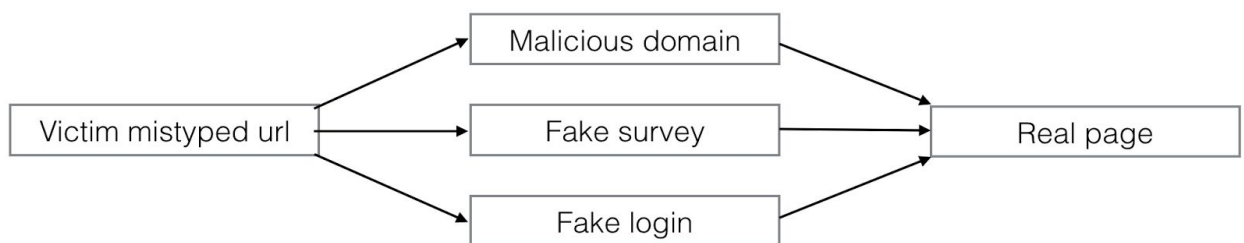
✓ DOMAIN	STATUS	YEARS	PRICE *
✓ wellssfargoc.om	AVAILABLE	1	\$ 268
Total Amount:			\$ 268

Research from [endgame.com](#) indicates how malicious authors uses this types of TLDs:

Registrant	Number of domains registered	Examples
Ahmed Al Amri	96	walgreens[.]om hotelsc[.]om bankofamerica[.]om youtube[.]om reddit[.]om
MuscatNet LLC	103	linkedin[.]om facebookc[.]om targetc[.]om live[.]om
Hassan Jaafar	80	yahoo[.]om linkedinc[.]om googlec[.]om baidu[.]om youtubec[.]om gmail[.]om
Shammy Shloma Sommet	16	xbox[.]om woot[.]om adidas[.]om hilton[.]om

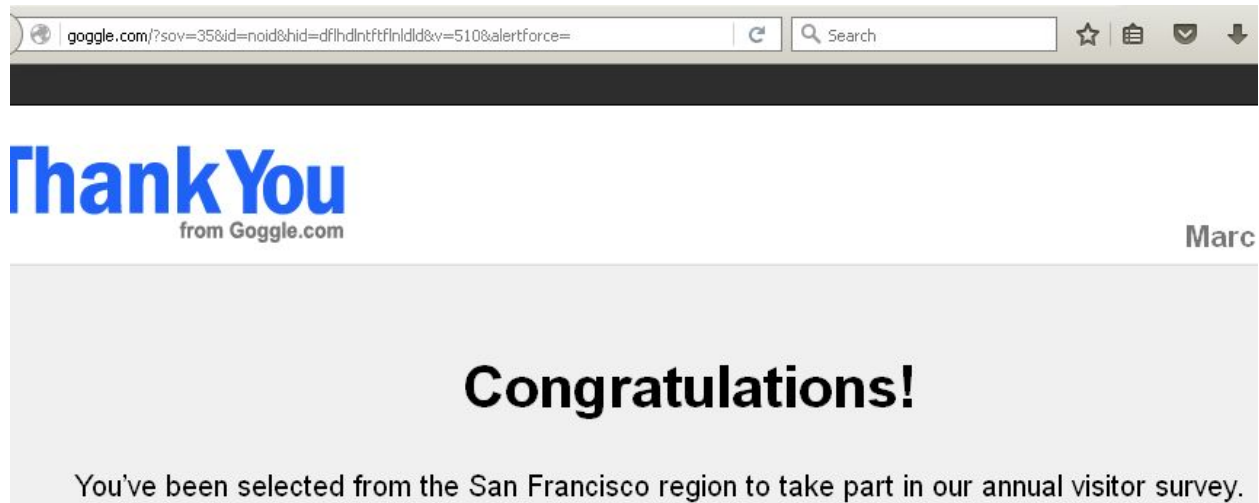
Many researchers also prove that most of the machines serving up these domains have severe unpatched vulnerabilities, including some which could provide arbitrary remote access. That is, these hosts could easily be exploited by other actors to serve up alternate (possibly worse) malicious content than what's currently being served.

Delivery scheme



The redirect to real web-page doesn't have place sometimes, but most sophisticated and intelligent phishing campaigns do it to hide malicious activity, that happens in between the redirects. Most users think they mistyped credentials and login to the legitimate site.

Example of fake survey on domain goggle[.]com:



2.2. Compromised phishing web pages

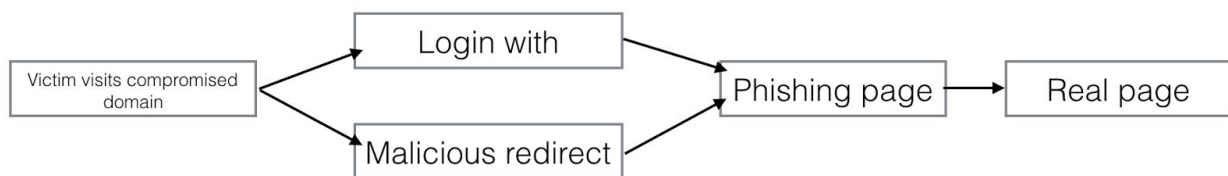
Compromised pages usually serves phishing pages within the domain. Implementation of them can vary from “login with” buttons to adding phishing web page as subdomain. This type of phishing often targets social media, blogs, news portals, e-commerce pages. In this case malicious author who got access to compromised domain embeds redirect to phishing page.

Accessing one of these pages tends to lead the user's browser to a few different web pages in a very short period of time, with the ultimate destination having content that may not even be relevant to the URI accessed in the first place. The redirections are in place for a few different reasons:

- The original URI can be made to appear somewhat legitimate, obscuring the path users will be forced to go down upon access.

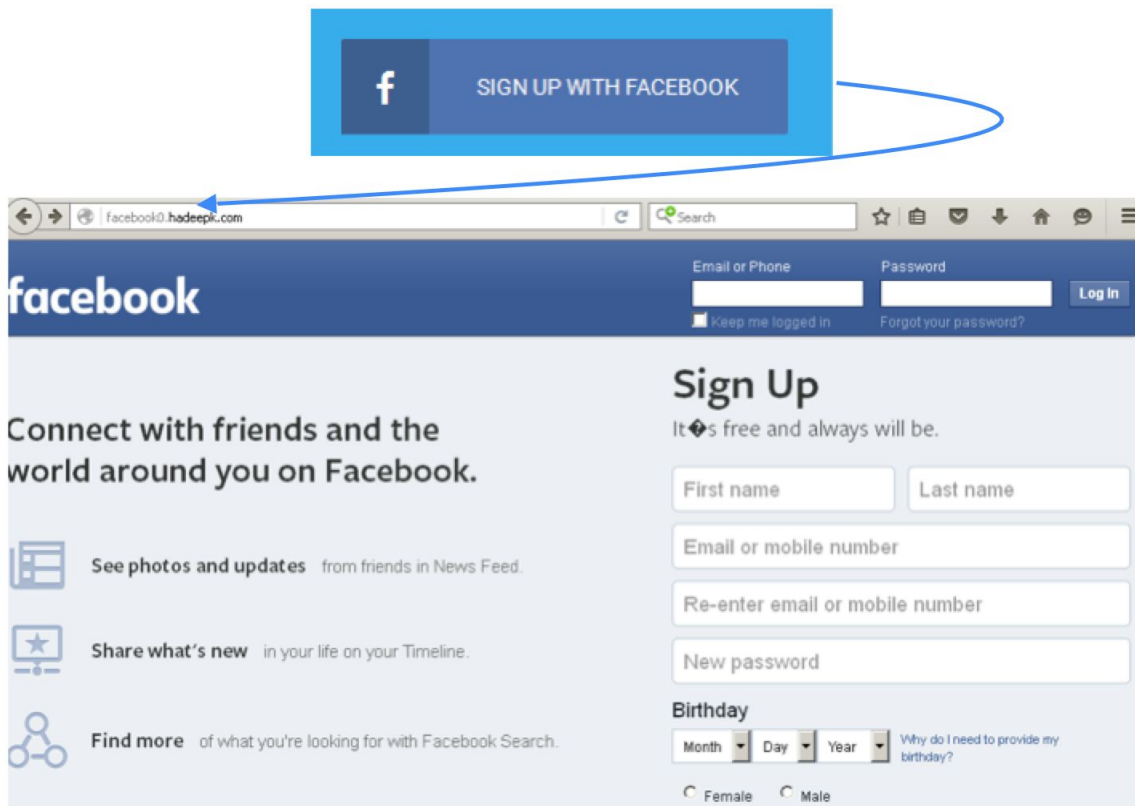
- The malicious actors can redirect the users to targeted platform-specific and / or location-specific content that may entice a naïve user to continue their journey further down the rabbit hole.
- The actors can change the destination web pages in an instant by modifying one or more of the redirect pages, thus allowing for easy pivoting to new pages or servers much like an incredibly frustrating game of Whack-A-Mole.
- Tracking cookies can be generated along the way to the ultimate destination and placed within the user's browser cache to surreptitiously monitor their behavior and provide further means for the actors to monetize a user's unfortunate trip to their site.

Delivery scheme



Examples:

Sign in with Facebook redirect to subdomain facebook0




Google drive phish served as index.php file




Welcome to Google Docs. Upload and Share Your Documents Se

Sign in with your email address to view or download attachme



Select your email provider

 Gmail

Sign in with Gmail

Email ***

Password ***

Sign in to view attachment

☐ Stay signed in [Need help?](#)

Access your documents securely, no matter your location



PAYPAL PHISHING FLOW

by OpenDNS Security Labs



a phishing email
is sent with a link
to redirectly-
paypal.com



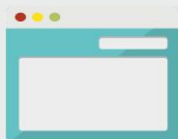
user clicks
on the link and is
redirected to a
fraudulent PayPal
site: security-
paypal-
center.com



the user is asked
to log into the
phishing site with
a PayPal
username and
password



user is prompted
to input a credit
card number, PIN
and CCV – input is
validated to ensure
that the card
number is real



user is prompted
to verify
information such
as name, address,
telephone number
and social security
number



user is redirected
to a page stating
that suspicious
activity has been
detected on the
account



user prompted
to upload a scan of
valid
government
issued ID – e.g., a
passport or
driver's license



user directed to
legitimate
PayPal.com
website as if
nothing happened



attackers use/sell
your personally-
identifiable
information (PII) –
which is far more
valuable than a
simple credit card
number



brought to you by :

OpenDNS

Big Data Engineering/Data Pipeline

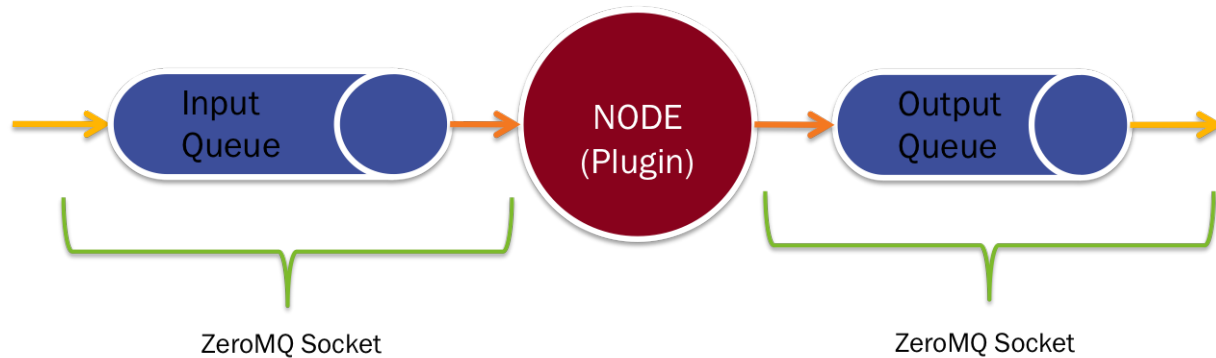
In a world where threat actors move fast and the Internet evolves in a non-deterministic fashion, turning threat intelligence into automated protection has proven to be a challenge for the information security industry. While traditional threat research methods will never go away, there is an increasing need for powerful decision models that can process data in a real-time fashion and scale to incorporate increasingly-rich sources of threat intel. This talk will focus on one way to build a scalable machine learning infrastructure in real-time on a massive amount of DNS data. In this talk, we will offer a sneak peek into how OpenDNS does scalable data science on ~80B DNS queries per day. We will touch on two core components, Big Data engineering and Big Data science, and specifically how they are used to implement a real-time threat detection systems for large-scale network traffic.

To begin, we will detail Avalanche, a stream processing framework that helps OpenDNS data scientists create their own data processing pipelines using a modular graph-oriented representation. Each node acts as a data stream processor running as a process, thread or EC2 instance. In this graph database, the edges represent streaming channels connecting the different inputs and outputs of the nodes. The whole data pipeline can then easily be scaled and deployed to hundreds of instances in an AWS cloud.

The Avalanche project's paradigm is to translate the approach that the finance world has been using for decades in high frequency or quantitative trading and apply it to traffic analysis. Applying intelligent detection models as close as possible to the data source holds the key to build a truly predictive security system, one where requests are classified and filtered on the fly. In our particular case at OpenDNS, we see a strong interest in integrating such a detection pipeline at the resolver level.

The Avalanche Project is where concepts from high frequency trading combine with high frequency classification. Avalanche follows the master/slave grid architecture. We manage the cluster with Boto & Fabric. Currently Avalanche is benchmarking around 30000 messages per

second per process. Avalanche uses ZeroMQ which is a very fast messaging system. The system works with a master-slave infrastructure. Here is an example of node (plugin) and edge (output stream) structure.



Detection

We will next discuss how we integrate our statistical model NLP-Rank into Avalanche, and show some benchmarks. NLP-Rank at its core is a fraud detection system that applies machine learning to the HTML content of a domain's web page to extract relevant terms and identify whether the content is potentially malicious or not. In this sense we are automating the security analyst's decision-making process in judging whether a website is legitimate or not. Typically when an analyst performs a review for a domain or URL in question, the analyst visits the site in a TOR browser, analyzes the content, and identifies the themes/summarize the page before deciding whether it's a fake or a false positive. In this talk, we will describe how we have automated this process at OpenDNS. The goal of our model is to basically automating this process for brand name associated phishes, and give back a probability on how close they think it is to a phish.

We will also discuss the unique heuristics of NLP-Rank, and the natural language processing techniques it uses. Additionally, we will discuss the design and implementation of our phishing classification system. We will provide an overview of data preprocessing techniques and the information retrieval/natural language processing techniques used by our classifier. We

will then discuss how Avalanche manages the results of NLP-Rank, how we add those results to our block lists and periodically retrain our corpus, and Avalanche's overall performance.

Heuristic 1:

Text-processing is resource intensive so when working with 80B+ DNS Queries a day we want to throw out as much data as possible. NLP-Rank filters traffic by ASNs, discarding traffic from all the ASNs that are associated with legitimate brands. For example you would expect a domain advertising a gmail security update to come from a domain associated with Google. In this way an ASN basically acts as one's zip code on the internet. We have built an ASN map of all legit companies of top spoofed brand names (ex. Google, Wells Fargo, etc.). Checks if ASN of IP of FQDN is in this map and not associated with top brand name ASNs.

Heuristic #2:

For the next step we have defined a 'malicious language' among FQDNs to be applied to the FQDNs in our authoritative logs. For this task we ran statistics on common n-gram collocations (co-occurrences) of dictionary words and brand names across a large set of malicious/phishing domains. Additionally we applied a common NLP technique called 'stemming', which gets the root of the word. Common stemmed dictionary words were: secur, updat, install, serv, etc. For example: paypal-security-update.info. These are characteristics exhibited by both APT and phishing domains.

Here are examples of APT domains:

Dark Hotel (Kaspersky):

- **adobeupdates[.]com**
- **adobeplugins[.]net**
- **adoberegister[.]flashserv[.]net**
- **microsoft-xpupdate[.]com**

Carbanak (Kaspersky):

- **update-java[.]net**
- **adobe-update[.]net**

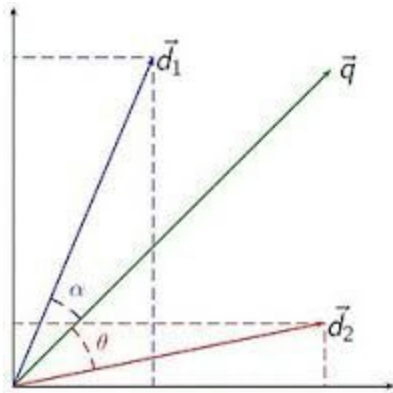
APT 1 Domains (Mandiant):

- **gmailboxes[.]com**
- **microsoft-update-info[.]com**
- **firefoxupdate[.]com**

We use a combination of edit-distance on substrings algorithm for brand names and regex patterns/automata theory flag based on those.

Heuristic #3:

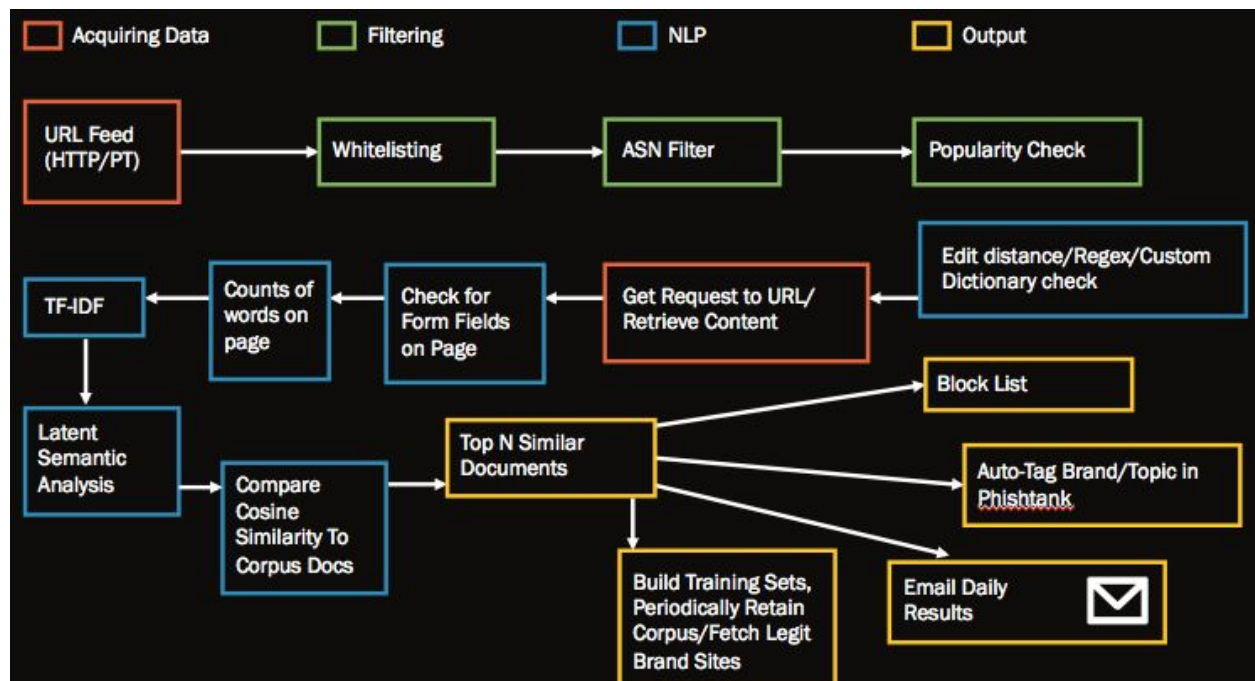
For domains that are flagged by the edit-distance and regex layer, the next step in the process is to do a GET request to retrieve the HTML content from the page. Since crawling HTML content is a heavyweight process we have created a distributed crawler to perform asynchronous requests to fetch the content when given an input list of domains. After we fetch the content we apply some minor preprocessing/cleaning and tokenization of the content on the page. Next step is converting the documents to vector space by tallying up the word counts, and then we apply TF-IDF (Term-Frequency-Inverse Document Frequency) as a weighting schema to put more weight on important words in the document. Term frequency is defined as # of times term t appears in document d . It is important to note term relevance does not increase proportional with term frequency. Inverse Document Frequency is defined as the # of documents that contain term t . TFIDF would then be defined as TF weight * IDF weight. Once we have the TFIDF weighted vectors, we apply an Unsupervised Algorithm, Latent Semantic Indexing (LSI), which applies a linear algebra technique called Singular Value Decomposition, which performs dimensionality reduction/ clustering on the term-document vectors. Now our model is ready and upon each input we convert our input document to LSI (vector) space and then compare against each of the documents in the corpus-matrix using cosine similarity (normalized dot product).



This returns a score between -1 and 1 on how similar the query is to the documents in the matrix. -1 meaning exactly opposite, 0 meaning independent, 1 meaning they are the same. One of the innovative things about this type of NLP is that we are using an unsupervised algorithm to generate features (topics) to make a supervised decision (comparing it to labeled corpus dataset). To quote Michal Sofka, the Cisco CTA Cognitive Threat Analytics Lead, *“Unsupervised Learning is the future. It’s all about the features.”*

Overview of the System:

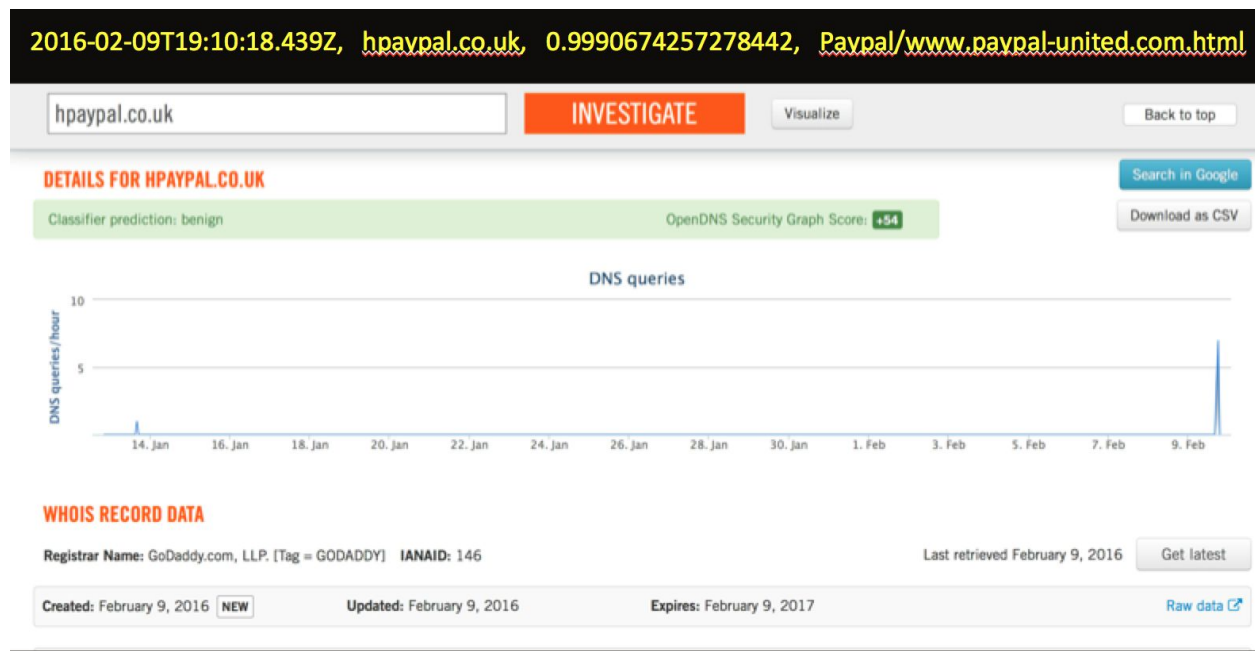
Figure 2 gives a visual overview of the whole system:

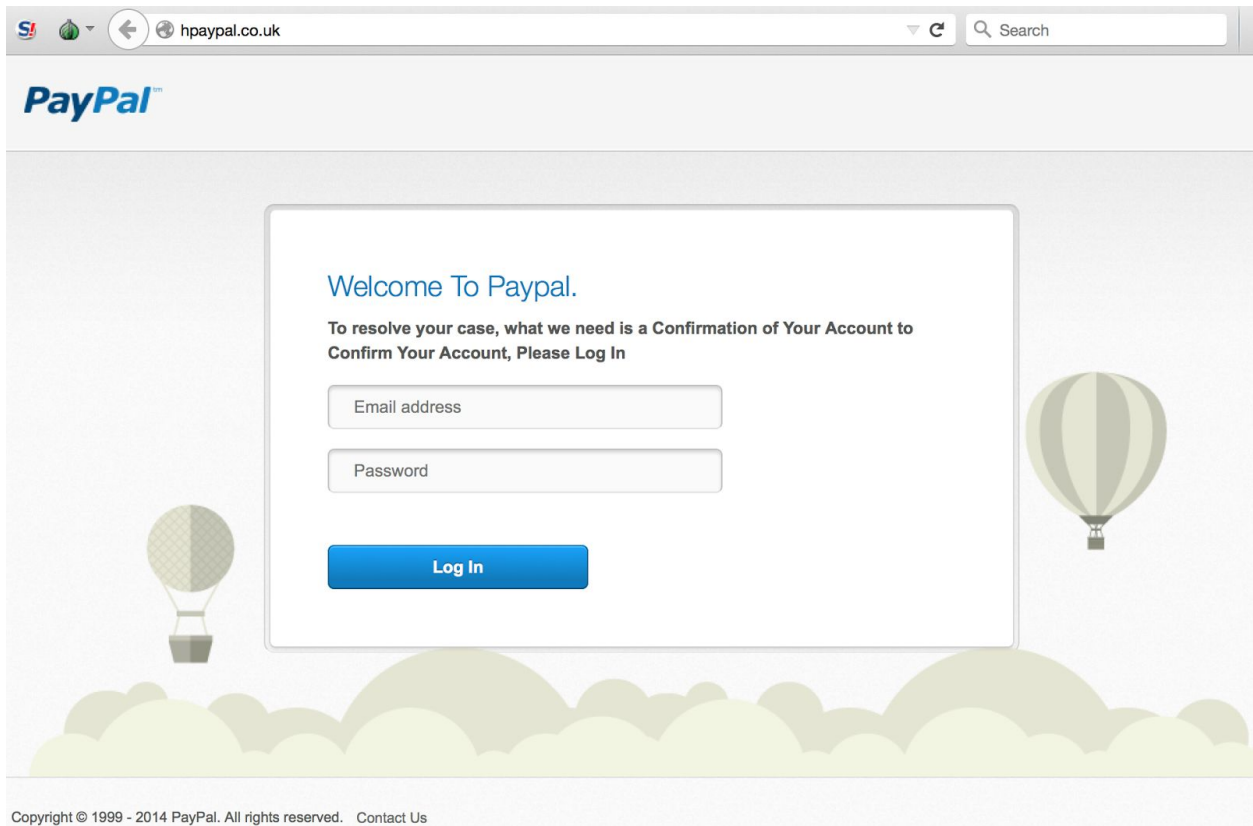


Results:

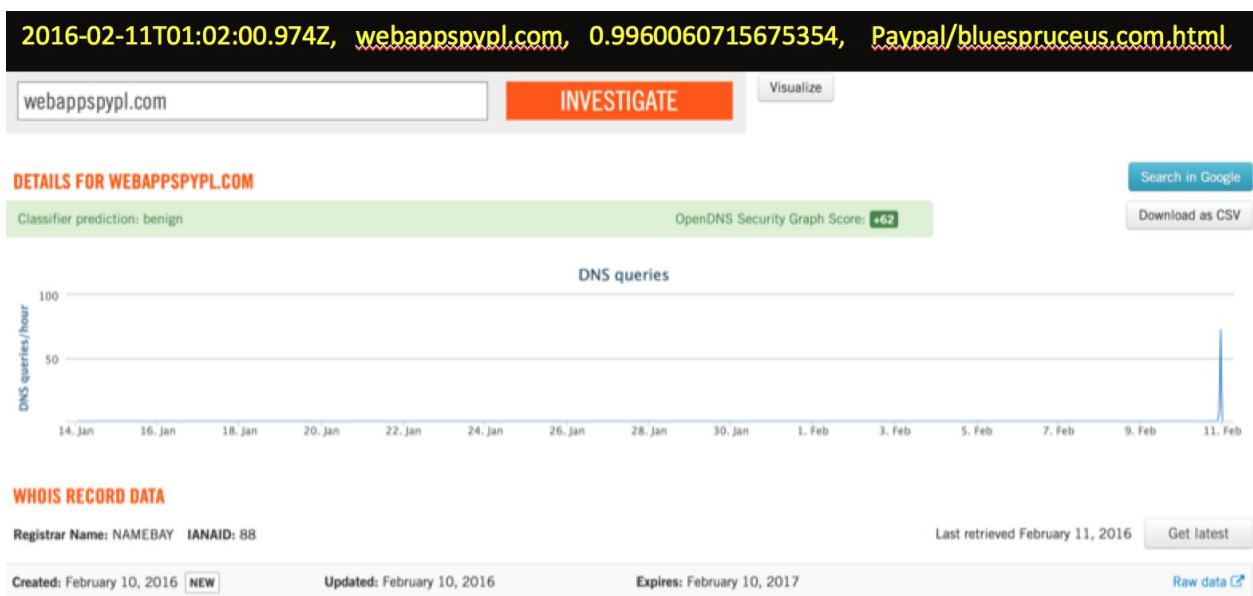
For our results we are catching a combination of compromised and dedicated phishing domains. The compromised domains we are detecting exhibit domain shadowing features. Most of the results among DNS traffic are dedicated hosting. Here are a few examples:

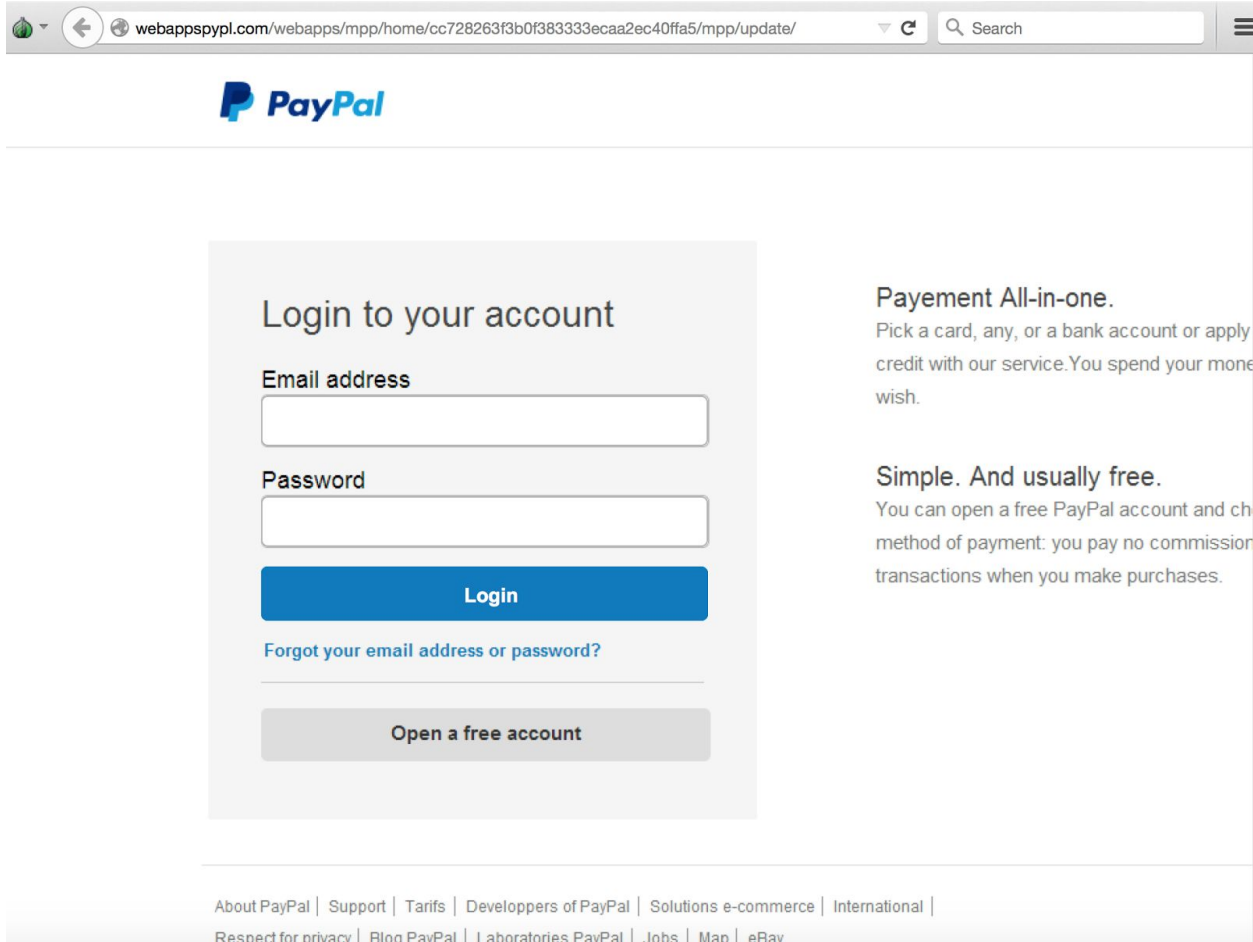
hpaypal[.]co[.]uk detected with 99.9% probability, February 9th 2016, created February 9th, 2016





webappspypl.com detected with 99.6% probability, February 11, created February 10, 2016.





Here is a spoof on Gmail, google-draive[.]com detected with high confidence score of 95.5%

googledraive.com

0.9555937051773071,

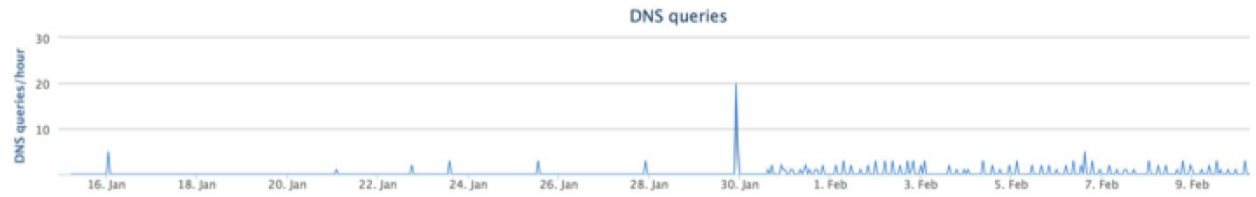
Google/www.rotisseriebuongusto.com.br

DETAILS FOR GOOGLEDRAIVE.COM

This domain is currently in the OpenDNS Security Labs block list

Classifier prediction: suspicious

OpenDNS Security Graph Score: **-94**



WHOIS RECORD DATA

Registrar Name: PDR Ltd. d/b/a PublicDomainRegistry.com IANAID: 303

Last retrieved February 6, 2016

Created: September 4, 2015

Updated: September 4, 2015

Expires: September 4, 2016

Email Address	Associated Domains	Email Type	Last Observed
s.miloshevich@yandex.ru	8 Total - 3 malicious	Administrative, Registrant, Technical	Current

Domain Name	Security Categories	Content Categories	Last Observed
alert-login-gmail.com	Malware, Phishing		Current
suporteng.com	Malware		Current
docsautentication.com			Current
g000glemail.com			Current
googledraive.com			Current
googlsupport.com			Current
membrana52.com			Current
pwdrecover.com			Current

Here was a site we picked up with high confidence:

security-appleinc.com 0.9436904191970825,Apple/www.lcloudid-ds.top.html



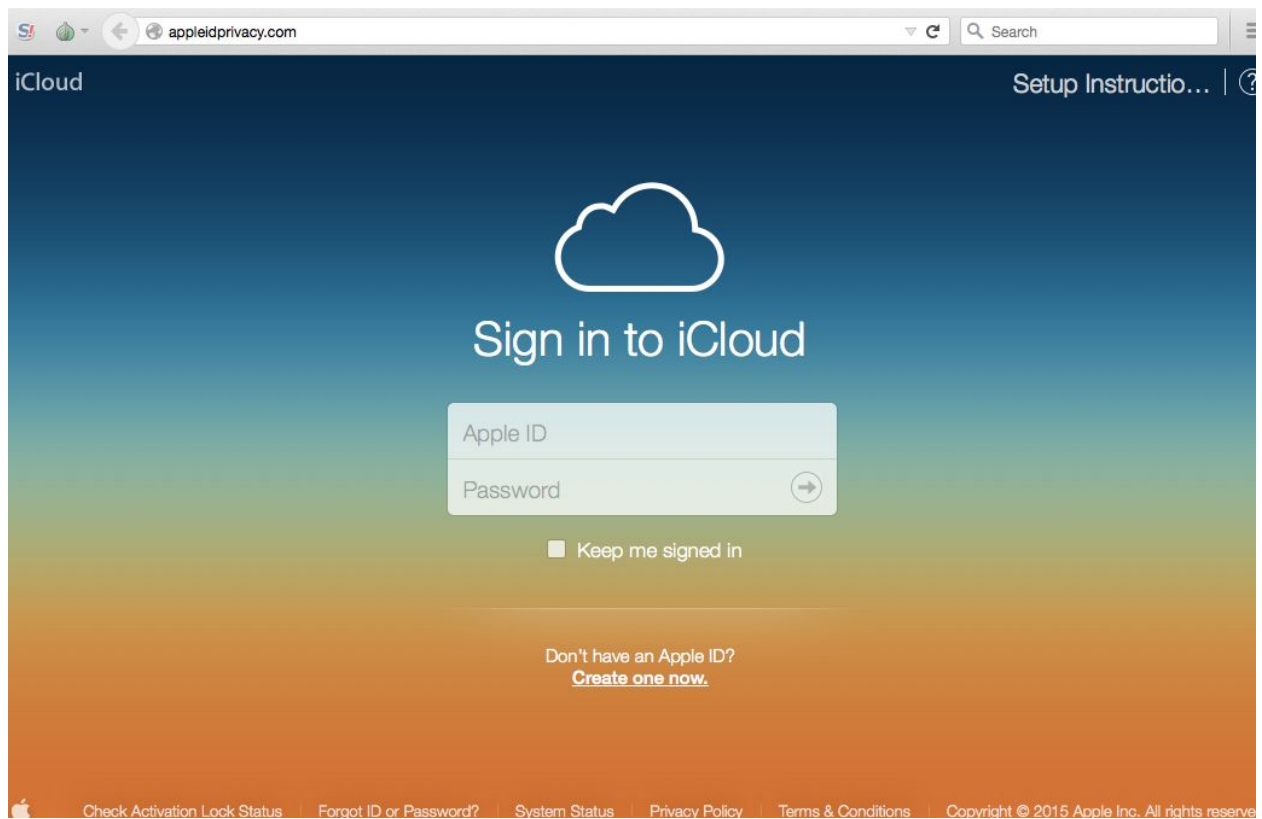
Pivoting on the email address associated with that domains:

DOMAINS ASSOCIATED WITH TRUSTEDMON@GMAIL.COM

Domain Name	Security Categories	Content Categories	Last Observed
appleidprivacy.com			Current
appleinc-security.com			Current
appleinc-support.com			Current
icloudprivacy.com			Current
inc-appleid.info			Current
secure-appleinc.com			Current
security-apple.com			Current
security-appleinc.com			Current

Showing 8 of 8 results

All very similar pages:



Here is an example where you can also pivot around IP address:

2016-02-26T16:15:01.418Z www[.]gmail-remind[.]tk 0.9633020758628845

Google/www.gmail-edit.pw.html



只需一個帳戶，便可通行所有 Google 產品與服務。

登入以繼續前往 Gmail

登入

☒ 保持登入狀態

[忘記密碼?](#)

[使用其他帳戶登入](#)

只要一個 Google 帳戶，即可體驗 Google 的各項服務



https://investigate.opendns.com/ip-view/104.207.132.165

104.207.132.165 INVESTIGATE Visualize Back to top

MALICIOUS DOMAINS HOSTED BY 104.207.132.165
gmail-edit.pw yahoo-maintain.pw gmail-safety.pw yahoo-safety.com gmail-retry.tk

FEATURES

Known domains hosted at this IP	20
LD2 domains count	16
LD3 domains count	20
LD2-1 domains count	15
LD2-2 domains count	20
LD2 domains diversity	0.8
LD3 domains diversity	1
LD2-1 domains diversity	0.75
LD2-2 domains diversity	1

KNOWN DOMAINS HOSTED BY 104.207.132.165
sg-images.yahoo-images.com www.gmail-remind.tk www.gmail-secure.tk www.yahoo-noreply.tk yahoo-noreply.tk gmail-retry.tk www.gmail-retry.tk phpinfo.pw www.gmail-safety.pw www.yahoo-safety.com www.gmail-edit.pw www.yahoo-maintain.pw yahoo-maintain.pw yahoo-protect.com www.gmail-maintain.tk www.yahoo-protected.tk www.yahoo-operation.tk accounts-163.tk www.yahoo-protect.com yahoo-protect.tk

www.yahoo-noreply.tk Search

YAHOO!
奇摩

服務說明

Yahoo奇摩讓你左右逢源，盡如人意。

無與倫比的 Yahoo奇摩電子信箱、重大地方新聞和國內外新聞、財經、運動、音樂和影視等精采內容。探索網上大千世界，一覽人間五光十色。

YAHOO!
奇摩

登入您的帳號

電子信箱

密碼

☒ 保持我的登入狀態

登入

無法存取自己的帳號？

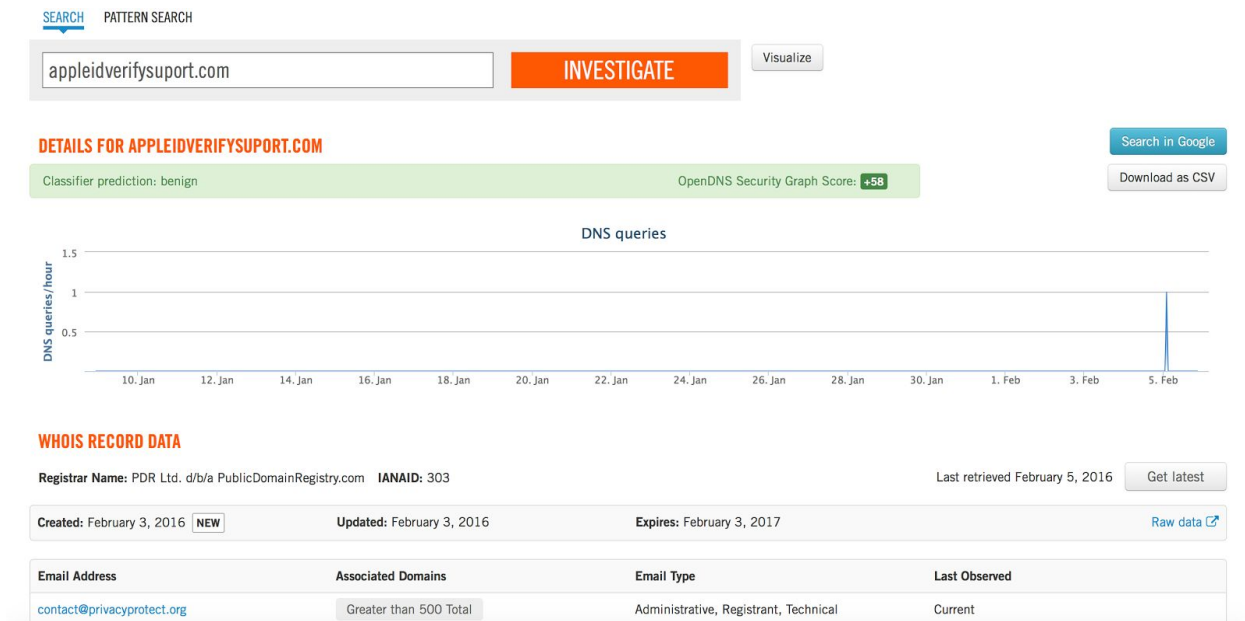
第一次使用 Yahoo奇摩？
註冊新帳號

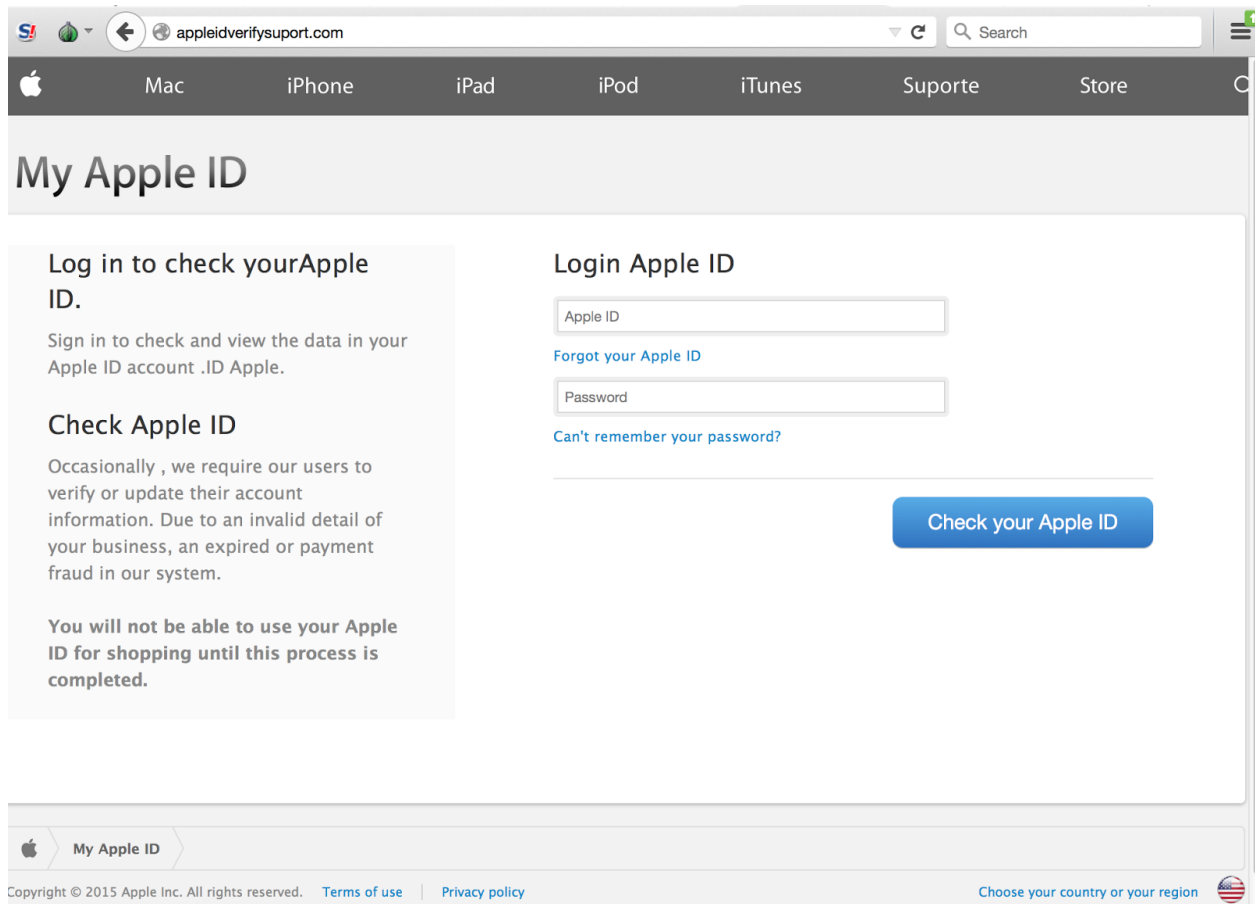
服務條款 | 隱私權

Here is an example of an Apple page created February 3rd, and then we see the first query on February 5th, detected with high probability after only 1 query, appleidverifysuport[.]com:

2016-02-05T22:06:28.109Z appleidverifysuport.com

0.8936428427696228,Apple/www.apple.uk.id-vrfy.email.html





Then after we added it to the corpus we were getting really high rate hits:

2016-03-02T04:21:00.552Z icloud[.]account-id[.]com 0.9970257878303528

Apple/appleidverifysuport.com.html

2016-03-17T10:10:37.374Z applesecurelogin[.]com 0.9799830317497253

Apple/appleidverifysuport.com.html

Here are some examples of domain shadowing/compromised sites:

The regular site ruralmaquinas[.]com[.]br:



Here is the compromised page located multiple subdomains deep:

`paypallimeited[.]com[.]cgi[.]bin[.]merchantpaymentweb[.]cmd[.]flowsession[.]05swp[.]ruralmaquinas[.]com[.]br`

paypalmeited.com.cgi.bin.merchantpaymentweb.cmd.flowsession.05swp.ruralmaquinas.com.br/743t

Search

PayPal

Personal Identification

Email

Password

Personal Informations

First Name

Last Name

Date of Birth

Address

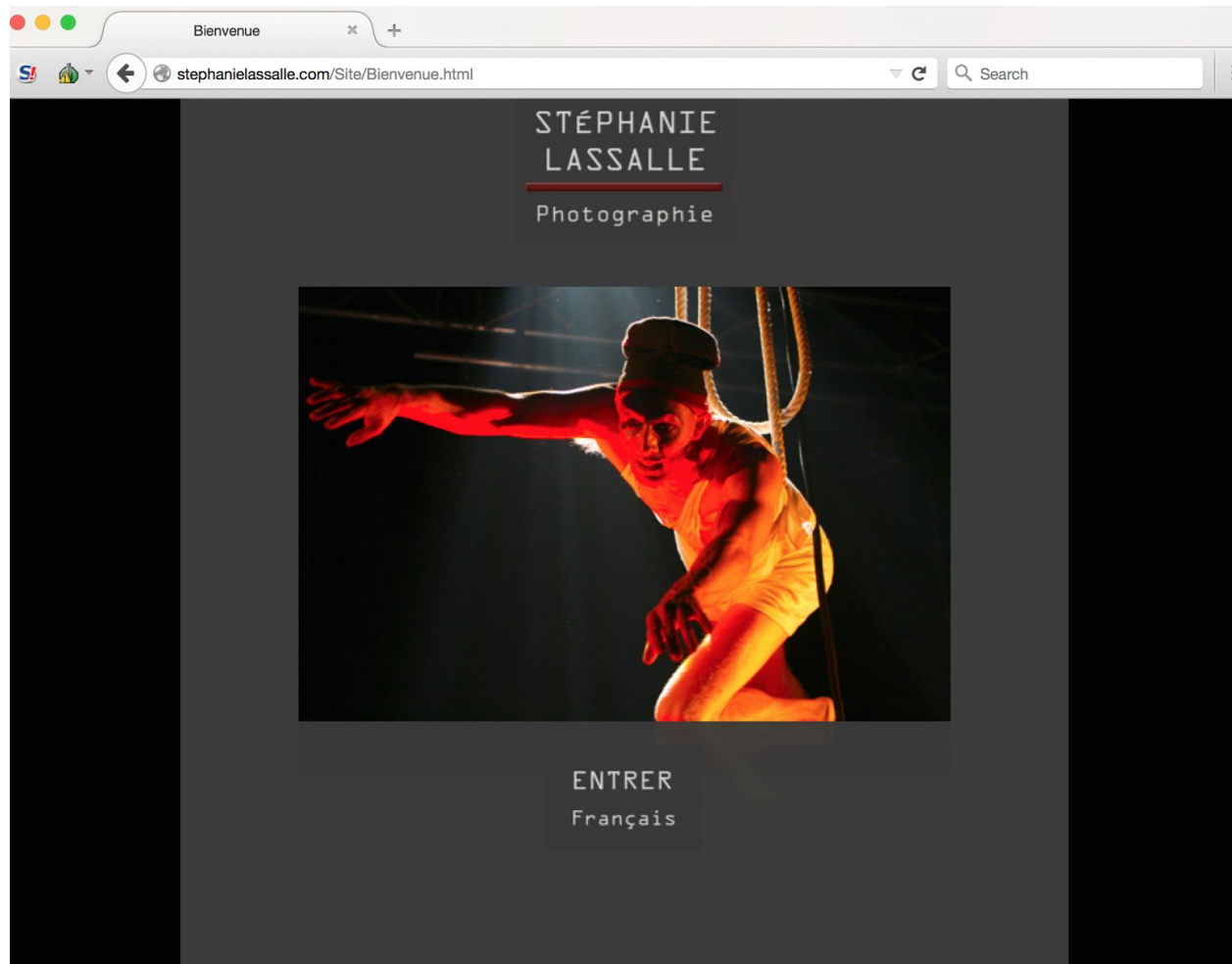
All in one pay.
Pick a card, any card... or a bank account. It's your money, you choose how to spend it.

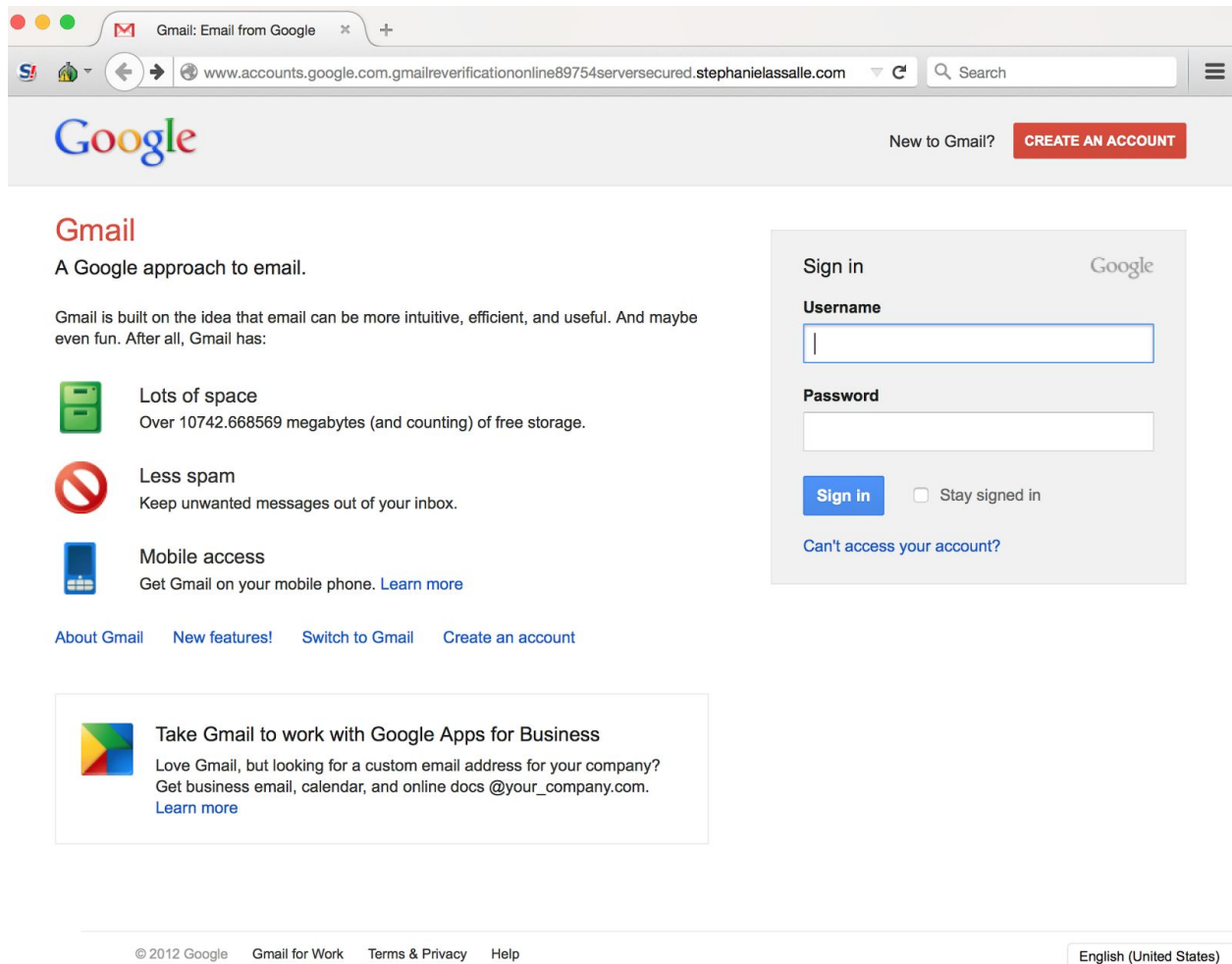
Simple. And usually free.
It's free to sign up for a PayPal account, and we don't charge you a transaction fee when you buy something, no matter how you choose to pay.

Here is another one:

2016-03-26T10:47:26.379Z

**www[.]accounts[.]google[.]com[.]gmailverificationonline89754serversecured[.]stephaniel
assalle[.]com 0.9815937280654907 Google/www.whiskey-memoirs.com.html**





Wellsfargo.com.amserver.ui.login.onlineaccounts.billing.account.updatemyaccount.wellsfargo.com.onlineaccounts.upgrade.online.billing.account.update.nlineaccounts.upgrade.o.cartoonnotworksoltion.org.html

Sss[.]www[.]facebook[.]com[.]uk2[.]gsr[.]awhoer[.]net

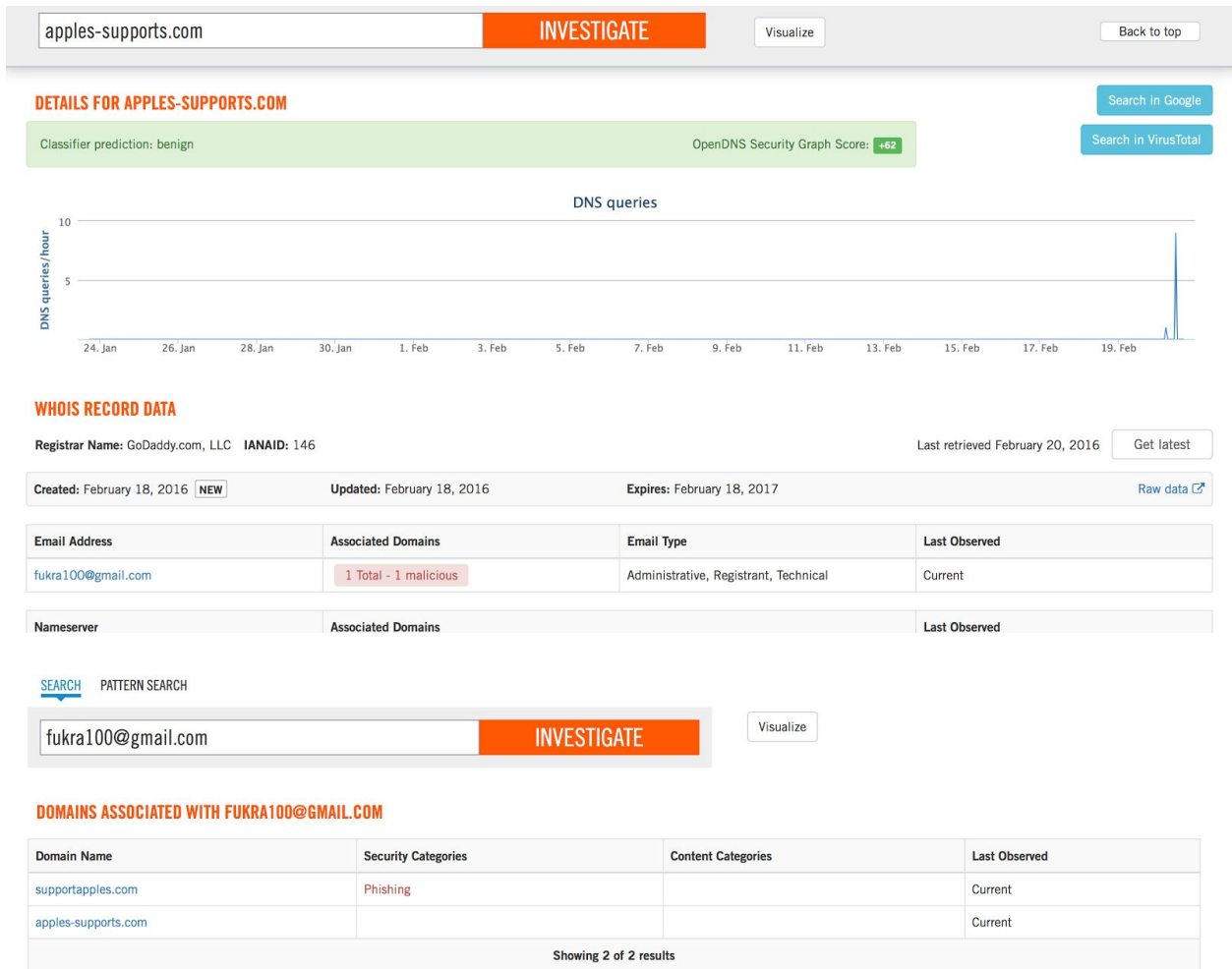
Here are some more examples of Dedicated Phishing domains:

2016-03-06T04:24:11.808Z paypal-accounts-security[.]com 0.9957063794136047

Paypal/bluespruceus.com.html

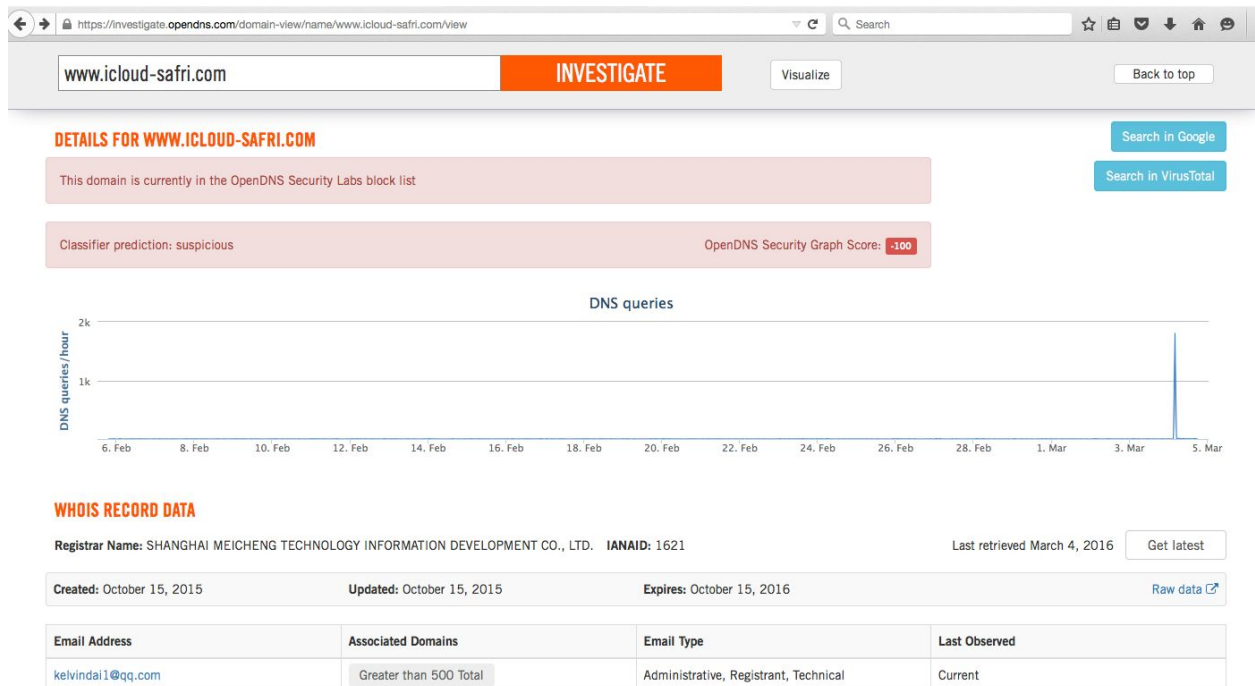
Traffic Patterns

Generally for dedicated hosting we see domains that are registered very recently and just become active for phishing. Here are a few examples:



This is different than compromised domains.

Model Stacking/Combining with Spike Detection:



As mentioned above the phishing sites we see generally have low traffic patterns, but [www.\[.\]icloud-safri\[.\]com](http://www.[.]icloud-safri[.]com) is an example that we caught that was particularly interesting in that it exhibited a large spike in traffic around March 4th. This is indicative of a large-scale phishing campaign, and we're able to detect this at OpenDNS Labs by stacking phishing and our spike model (worked on by Dhia and Thomas).

Caught by Phishing Detection today close to around 95% probability:

2016-03-04T10:39:38.372Z [www.\[.\]icloud-safri\[.\]com](http://www.[.]icloud-safri[.]com) 0.9479695558547974

Apple/applelock.ru.html

Confirmed with Thomas also caught by Spike:

icloud-safri.com. 1.0 1749 1749.0 1514 23 {{{(nyc),109},{{(ash),88},{{(chi),95},{{(yvr),17},{{(ams),110},{{(cdg),39},{{(prg),27},{{(yyz),24},{{(nrt),12},{{(sin),71},{{(otp),24},{{(fra),227},{{(jnb),18},{{(dfw),135},{{(lax),43},{{(wrw),106},{{(pao),40},{{(mia),301},{{(ber),8},{{(sea),43},{{(hkg),70},{{(syd),17},{{(lon),125}}}}}}}}}}}}}} {{{(1),1749}}}

False Positives:

Some of the False Positives we are identifying are domains which are protected by the service MarkMonitor/CSC domains which are typically typosquatting looking domains that redirect back to the legitimate website. For example: wellsforge.com → redirects back to → wellsfargo.com. This actually proved to be useful in our research because it is an indicator we are on the right track being that MarkMonitor protects against brand-infringement, and testing against the legitimate site itself will give back high cosine similarity scores. Markmonitor domains function almost like test data interspersed within our live DNS feeds.

Some of the other FPs we're having are not phishes, but are not necessarily sites people will miss. Ex. googles.ru redirecting to payperinstall.ru.

Immediate Goals/Outcomes:

We have a couple primary goals/outcomes that we are focusing on as a result of this detection technique. First and foremost is to detect dedicated and compromised domains in live DNS Traffic. Additionally now that we are acquired by Cisco we are working on integrating NLPRank with their Cloud Web Services/Proxy/HTTP logs, and also their E-Mail Corpus. The second is to apply NLPRank as a recommender system to reduce submission to verification time on Phishtank and push those out to community. Another thing we are working on is stacking models, similar to the example we gave above about combining Spike Detection (worked on by Thomas and Dhia) and Phishing Detection models to discover large-scale phishing campaigns as they are happening. Another concept we are working on trying to detect targeted attacks, while DNS is limited we are optimistic with the addition of the Proxy and Email Corpus dataset we can start detecting them.