# Detecting Phishing using Visual Similarity
## Outline

- Current Tools & Techniques

- Web Crawling

- Creating Datasets

- Distance & Similarity

- Grouping Images

- Testing Detection

- Beyond Images

  - HTML Similarity

  - Scraping Text from Images

  - Text Classification with LLMs

- Research

- Alerting

- Action

If you see text like this, it's a note for this pdf version of the slides so you don't have to guess what's going on.
If you see my face next to the text, it's just me "saying" the note so you can easily differentiate from other text on a busy slide.

DEEPSEC

# Current Tools & Techniques

## Services

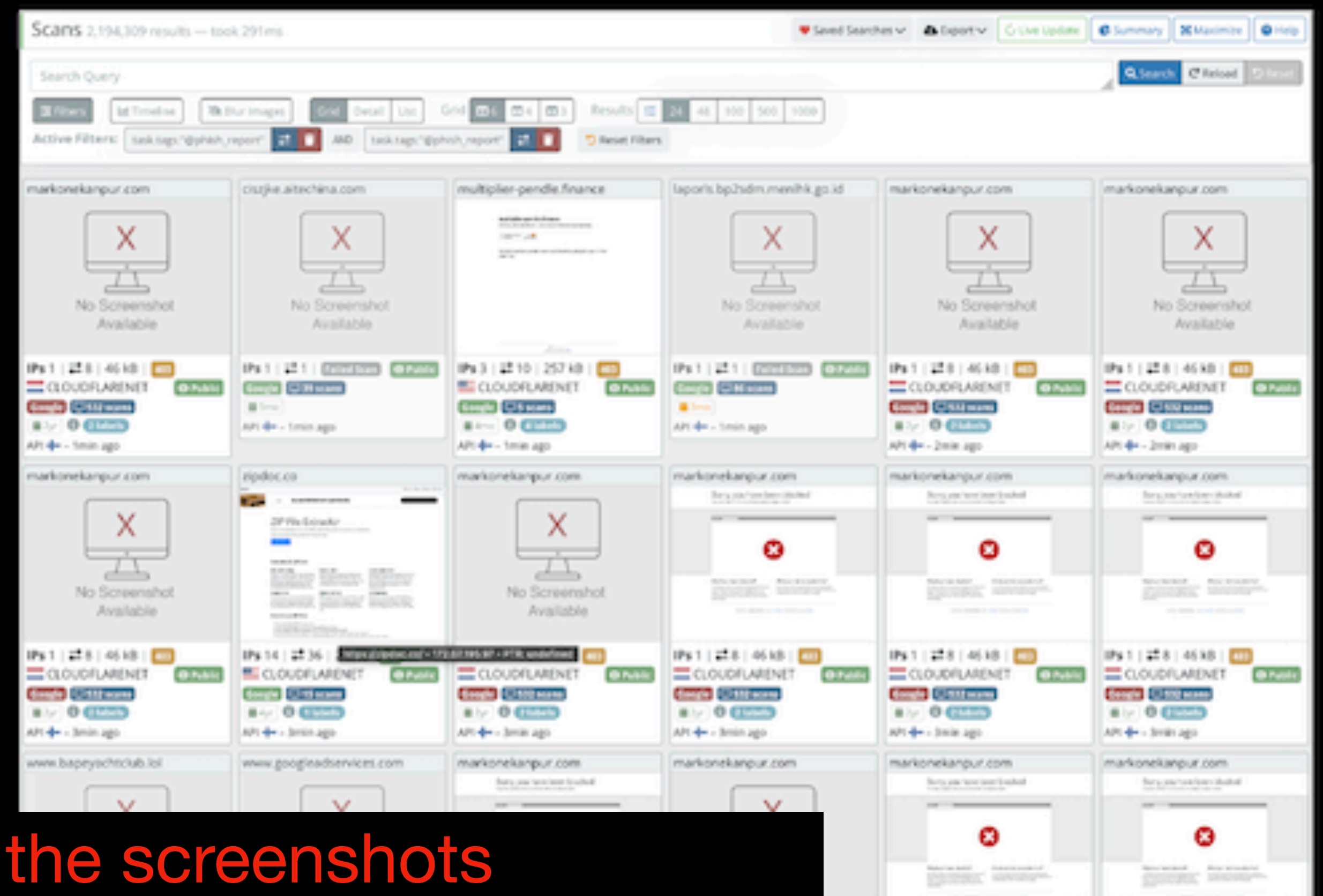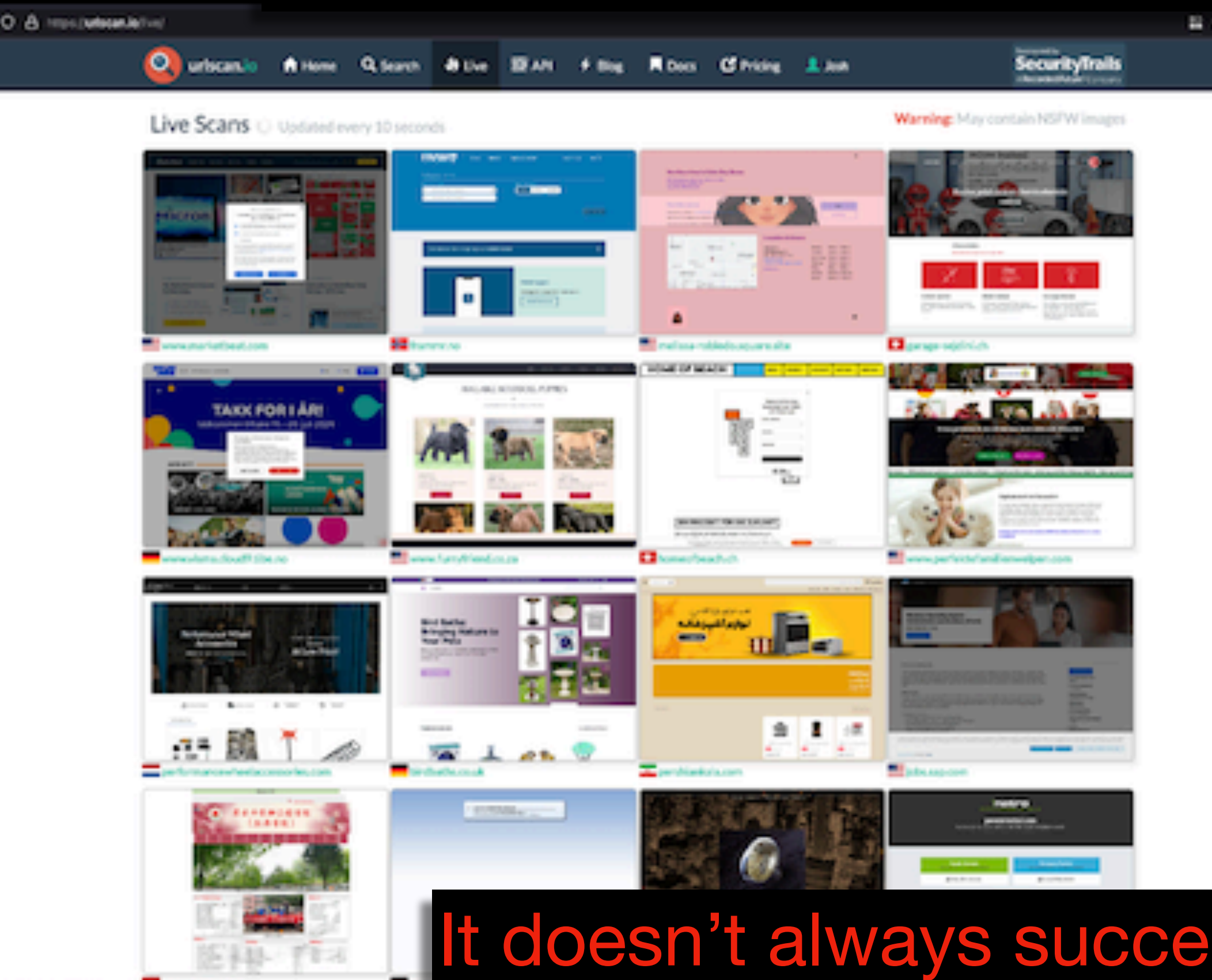# URLScan.io

## Crawling

## Threat detection

URLScan takes screenshots and provides information on websites.
There is a YARA search that can be setup for certain kinds of website behavior





It doesn't always succeed at the screenshots
And it doesn't do extensive analysis on the screenshots

# urlscan.io Phishing URL Feed

urlscan.io detects thousands of suspicious, malicious and phishing URLs every day. Our organic phishing URL detection is able to pinpoint these attacks and associate them with one of the 680+ popular brands that we track. The feed of these detected URLs is available to our customers to ingest

**URLScan has a phishing feed that's useful. In our line of work, we all likely ingest feed from various locations as a quick way to use the intelligence and work of others to improve our security. These feeds can also be used to gather data about the current threat landscape.**

- IP address & Geo-information of the IP hosting the phishing page
- ASN and ASN Name hosting the phishing URL

You can download a static 7-day sample of the feed here:

⊕ Download JSON Sample    ⊕ Download CSV Sample

## Phishing URL Feed

> 1500 detected URLs / day
> 680+ brands tracked
✔ Live URL retrieval
✔ Commercial use
✔ JSON & CSV format

**Targeted Brand Information**
Brand name
Brand country & industry vertical

**Phishing Page Information**
Phishing page metadata
Phishing domain metadata
Phishing IP metadata

# Current Tools & Techniques

## Detection Methods
Threat hunting
User reporting
Various Products

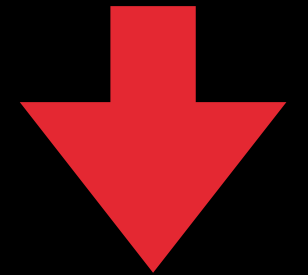# Threat Hunting

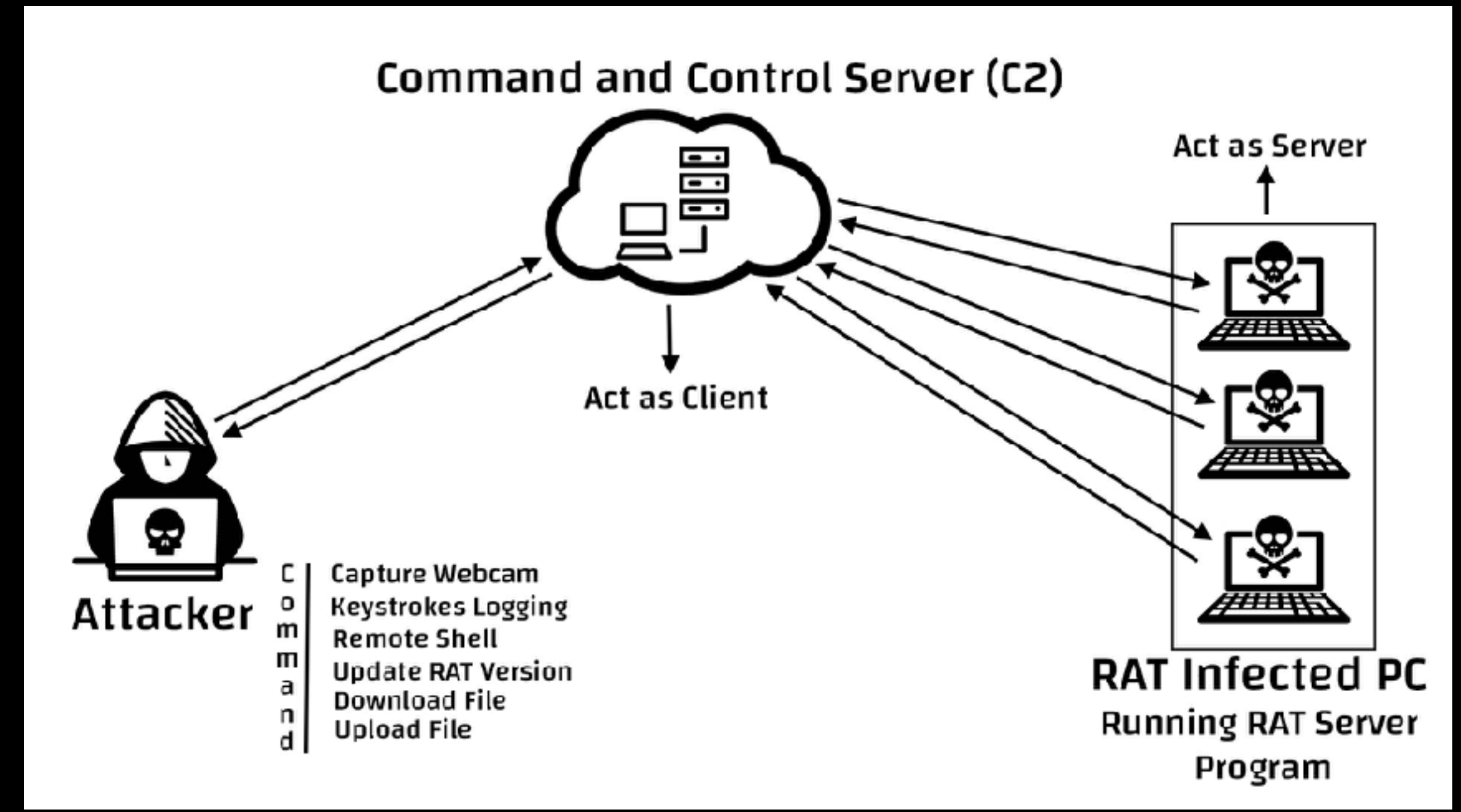## Compromised Infrastructure

```
<iframe src="http://far.IAAS.NEWS/?biw=OMITTEDURI" width="263" height="257"></iframe>
```
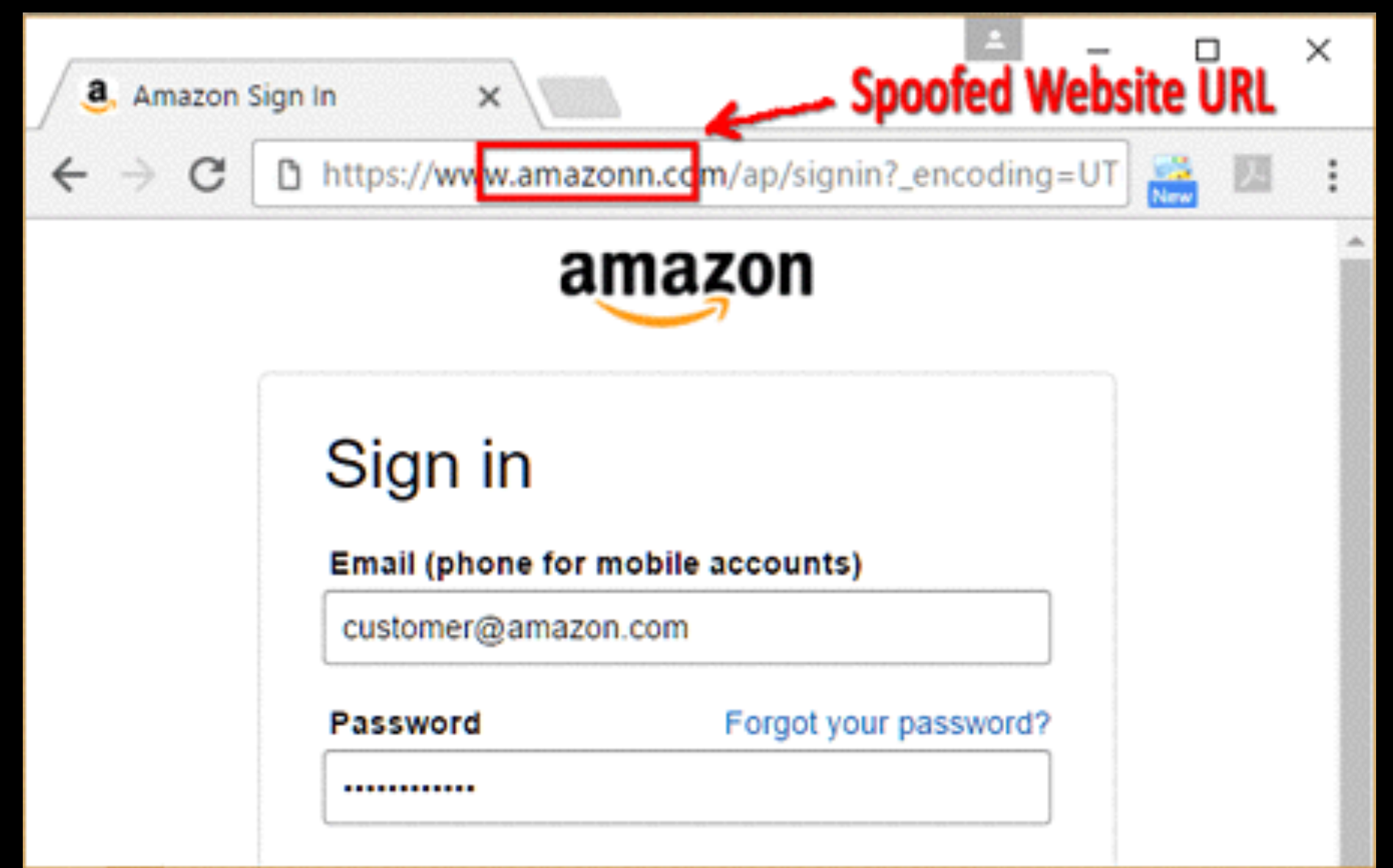
```
hXXp://www.fullcircleliterary.com/
hXXp://danielpsheehan.com/areas-of-
expertise/educator/ucsc-2016-rulers-
of-the-realm
hXXp://danielpsheehan.com/
hXXp://www.cafemuseroyaloak.com/
hXXp://kdsross.com/about-us/
hXXp://usdiagnostics.com/index.php/
certification-testing/uscreen-cup
hXXp://psychologywiththal.com/
2015/09/30/life-span-development-
personality/
hXXp://
thefecaltransplantfoundation.org/
what-is-fecal-transplant/
hXXp://optimalwellnessaz.com/about/
hXXp://optimalwellnessaz.com/about/
hXXp://chworks.org/real-estate-
```

## C2



Command and Control Server (C2)

Act as Server

Act as Client

Attacker — Command: Capture Webcam, Keystrokes Logging, Remote Shell, Update RAT Version, Download File, Upload File

RAT Infected PC
Running RAT Server Program

## Phishing

# Threat Hunting cont...

DNS



URL Analysis

# Current Tools & Techniques
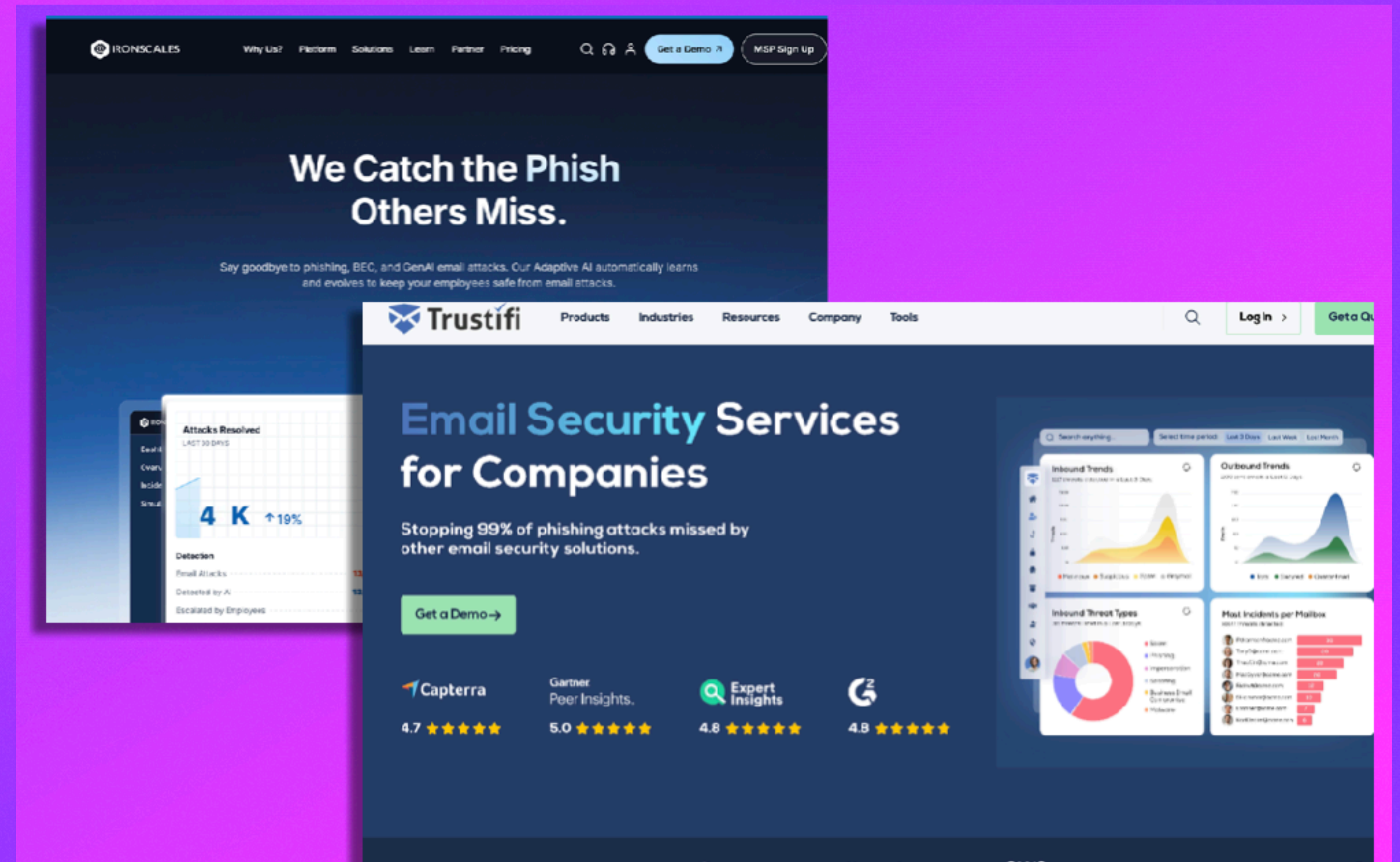
User reporting is another source of malicious activity detection

And there are multiple products you can pay for to increase detection

## User Reporting

## Products

**Security Is Made of Layers**

**Phisning Tests**

**Various Products**

The work I'm demonstrating are part of a layered approach.

There's always room for more layers to improve security

**User Reports**

**Egress Analysis**

# What do we have to do?

**Crawl Websites**
**Take Screenshots**
**Build Datasets**
**Find Similar in…**

**Images**

**HTML**

**Text**

**Infrastructure**

# Web Crawling

## Crawler Architecture



Send URLs → to Selenium → and take screenshots

```
●→  Detect Phishing with Similarity Searching python url_crawlerv2.py urls.txt
SUCCESS: https://pyosec.com
SUCCESS: https://google.com
SUCCESS: facebook.com
Results saved to output_20240828_2315/results.html
○→  Detect Phishing with Similarity Searching █
```

Demo of three websites being crawled, after which, I am left with the screenshots and a results.html file

| Name | Date Modified | Size | Kind |
|------|---------------|------|------|
| facebook.com.jpg | Today at 23:15 | 14 KB | JPEG image |
| google.com.jpg | Today at 23:15 | 15 KB | JPEG image |
| pyosec.com.jpg | Today at 23:15 | 22 KB | JPEG image |
| results.html | Today at 23:15 | 603 bytes | HTML text |

output_20240828_2315

| Thumbnail | URL | Status |
|-----------|-----|--------|
| | https://pyosec.com | Success |
| | https://google.com | Success |
| | facebook.com | Success |

Screenshots

URLs

Status

This is the start of collecting crawling data, but it will increase/improve through the presentation

| Thumbnail | URL | Status |
|---|---|---|
| Failed to generate thumbnail | https://docs.google.com/presentation/d/e/2PACX-1vQzmTfShxetobTbZx9gY00VgL-gXRC9gFeU_6RLdklj1LXKN3UTNWNtwYDlx5OQD9xTj378VsXZiErM/pub/?start=false&loop=false&delayms=3000 | Failed |
| Failed to generate thumbnail | https://taplink.cc/ahhdtttt/ | Failed |
|  | http://avionhealthcare.in/ | Success |
|  | http://paramojuntplus.com/ | Success |
| Thumbnail | https://urk... | Success |
| Failed to generate thumbnail | https://docs.google.com/presentation/d/e/2PACX-1vRSrJlzvElfXLuaw5hXN1bLC0zWfF7d5-k_408_eZJxvTGmsdU7BRpig_Vb_vW3f8qHKkDJlRJondk6/pub/?start=false&loop=false&delayms=3000&slide=id.p%3E | Failed |
| Failed to generate thumbnail | https://docs.google.com/presentation/d/e/2PACX-1vR1rwYJjTCFYDvYy8kVg9R1fBH1jg6SDbZHT4ZPSlnJ9OCVCsNNj2uc7IzIpgaktqB8vADFTqHW3Iwf/pub/?start=false&loop=false&delayms=3000 | Failed |
|  | http://ayurvationhealthcare.org/ | Success |
| Thumbnail | http://gentleskinhealthllc.org/ | Success |
| Failed to generate thumbnail | https://docs.google.com/presentation/d/e/2PACX-1vR3V6FwlkStTwTGj_qT89ZrQY6VN_Ow4uUMf8gpeMQLXMGdZ8SmHY1MbiLA0LLKxoAcONa6JdoC12Sv/pub/?start=false&loop=false&delayms=3000&slide=id.p%3E | Failed |
| Failed to generate thumbnail | https://foodblogspottingeleberate.blogspot.com/?m=1 | Failed |
| Failed to generate thumbnail | https://docs.google.com/presentation/d/e/2PACX-1vSuWITrMvwpFAiW3Jyqs_DisVlUP1BH62mBJkFHTCFzbZv0Afb4dQnD5ooWFOy3i1ASsZ1ZjOJnHmvJ/pub/?start=false&loop=false&delayms=3000 | Failed |
| Failed to generate thumbnail | http://ancient-thunder-0448.chinnabhai944.workers.dev/personalization/cl2/freeform/websitedetect/?source=wwwhead&fetchtype=css&modalview=nmlanding | Failed |
| Failed to generate thumbnail | https://docs.google.com/presentation/d/e/2PACX-1vSbzx_-G6NM9PvAnaXR-G_STx-EQ6mcoAa9CJ4VeD6_G4YAQzTBbUM7fBKK31Cev8ZtGipLO8tzTsPt/pub/?start=false&loop=false&delayms=3000&slide=id.p | Failed |
| Failed to generate thumbnail | https://enews.classicfirearms.com/q/bh-yefdf-1KDwz0X30bk7wbN8xBQdmmgRdcZcOJcmVwb3J0LcGhpc2hpbmdAYW50aXBoaXNoaW5nLm9yZ8OICri-Y8SAcL_YktPQzcdBxQgxtg/ | Failed |
| Thumbnail | http://conscioushealthcafe.com/ | Success |

**What it looks like when sending a large list of URLs**

# Issues/Obstacles when Crawling

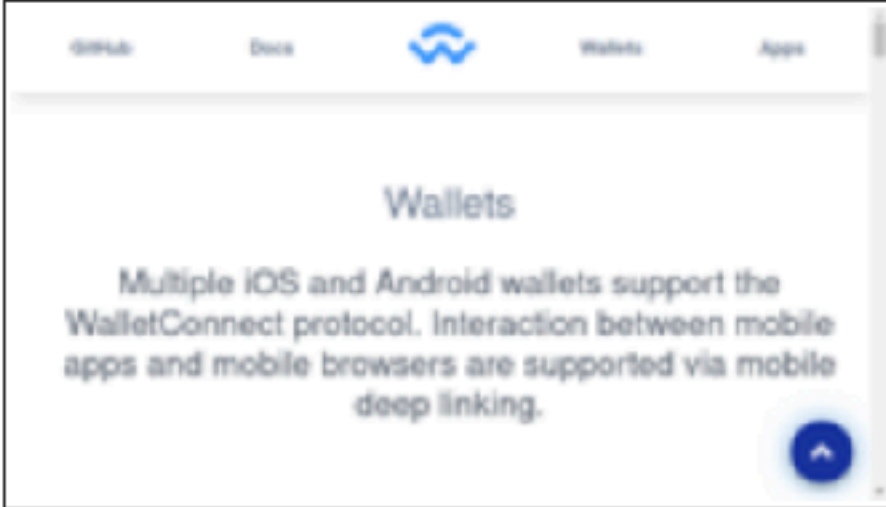# Crawl errors

Failed to generate thumbnail errors…

| Thumbnail | http://gentleskinhealthllc.org/ |
|---|---|
| Failed to generate thumbnail | https://docs.google.com/presentatio start=false&loop=false&delayms= |
| Failed to generate thumbnail | https://foodblogspottingcleberate.b |
| Failed to generate thumbnail | https://docs.google.com/presentatio |
| Failed to generate thumbnail | http://ancient-thunder-0448.chinna |
| Failed to generate thumbnail | https://docs.google.com/presentatio |
| Failed to generate thumbnail | https://cnews.classicfirearms.com/c |
| Thumbnail | http://conscioushealthcafe.com/ |

| Thumbnail | URL | Status |
|---|---|---|
| Thumbnail | http://gentleskinhealthllc.org/ | Success |
| Thumbnail | http://conscioushealthcafe.com/ | Success |

Was able to crawl, but no screenshot generated…

The crawlers run headless selenium and can be modified easily

# Fixed with more info

| Thumbnail | URL | Status | HTTP Status Code |
|---|---|---|---|
| Shaw) Webmail — Sign in to access your Shaw email | hotmail-105506.weeblysite.com | Success | 200 |
| | jpyorre.com | Success | 200 |
| Failed to generate thumbnail | emapdiwhf7.nnnn.eu.org | Failed | none |
| Wallets — Multiple iOS and Android wallets support the WalletConnect protocol. Interaction between mobile apps and mobile browsers are supported via mobile deep linking. | web3autofix.pages.dev | Success | 200 |

# Stopped at Error

**FAIL**

| Thumbnail | URL | Status | HTTP Status Code |
|---|---|---|---|
| | hotmail-105506.weeblysite.com | Success | 200 |
| | jpyorre.com | Success | 200 |
| | web3autofix.pages.dev | Success | 200 |
| | | | |

Other page crawling errors. A human would click, but the crawler stops

# Use Selenium

Read Text



Click through

PROFIT!

Tell Selenium to click as if it were a human…

**And then realize the privacy error was because of your own protection URL**



hotmail-105506.weebl

jpyorre.com

# Bypass Security

**Your connection is not private**

Attackers might be trying to steal your information from **domain.com** (for example, passwords, messages, or credit cards). Learn more

NET::ERR_CERT_COMMON_NAME_INVALID

Hide advanced          Back to safety

This server could not prove that it is **domain.com**; its security certificate is from **\*.kinsta.com**. This may be caused by a misconfiguration or an attacker intercepting your connection.

Proceed to domain.com (unsafe)

Cisco Umbrella

You can't run the crawlers on a network where you are running your threat mitigations

# Build a System and use custom DNS

Summary

| Console
| Hardware
| Cloud-Init
| Options
| Task History
| Monitor
| Backup
| Replication
| Snapshots
| Firewall
| Permissions

Hour (average)

**urlcrawler (Uptime: 39 days 21:07:18)**

Notes

Status — running
HA State — none
Node — athena

CPU usage — 0.89% of 1 CPU(s)
Memory usage — 85.63% (5.14 GiB of 6.00 GiB)
Bootdisk size — 50.00 GiB

IPs — No Guest Agent configured

In my case, I built a VM with custom networking and DNS to bypass my security

4.5
4
3.5
3
2.5
2
1.5
1
0.5
0

%

2024-11-06 2024-11-06 2024-11-06 2024-11-06 2024-11-06 2024-11-06 2024-11-06 2024-11-06 2024-11-06 2024-11-06 2024-11-06 2024-11-06 2024-11-06
21:08:00   21:14:00   21:20:00   21:26:00   21:32:00   21:38:00   21:44:00   21:50:00   21:56:00   22:02:00   22:08:00   22:14:00

# Use Selenium

Read Text

Click through

PROFIT!

…possibly for real this time

| Thumbnail | URL | Status | HTTP Status Code |
|---|---|---|---|
|  | jpyorrc.com | Success | 200 |
|  | hotmail-105506.weeblysite.com | Success | 200 |
|  | web3autofix.pages.dev | Success | 200 |

```python
48      class Webdriver(object):
124         def handle_insecure_connection(self):
126             try:
131                 if warning_present:
                        advanced_button.click()
139                     time.sleep(1)  # Sleep only after performing this action
140
141                     # Click on 'Proceed to site' link if it's clickable
142                     proceed_link = WebDriverWait(self.driver, 3).until(
143                         EC.element_to_be_clickable((By.ID, "proceed-link"))
144                     )
145                     proceed_link.click()
146                     print("Clicked through the warning.")
147                     time.sleep(1)  # Sleep only after proceeding
148
149             except Exception as e:
150                 print(f"No 'Your connection is not private' warning detected or unable to find elements: {e}")
151
152         def handle_dangerous_site_warning(self):
153             """Handle 'Dangerous site' warning (phishing, malware, etc.)."""
```

Some of the code showing that you can look for any text
and tell selenium to click the link.

```python
160                 print("Detected 'Dangerous site' warning.")
161
162                 # Click on 'Details' button if it's clickable
163                 details_button = WebDriverWait(self.driver, 3).until(
164                     EC.element_to_be_clickable((By.ID, "details-button"))
165                 )
166                 details_button.click()
167                 time.sleep(1)  # Sleep only after performing this action
168
169                 # Click on 'this unsafe site' link if it's clickable
170                 unsafe_link = WebDriverWait(self.driver, 3).until(
171                     EC.element_to_be_clickable((By.XPATH, "//a[contains(text(), 'this unsafe site')]"))
172                 )
173                 unsafe_link.click()
174                 print("Clicked through the 'Dangerous site' warning.")
175                 time.sleep(1)  # Sleep only after proceeding
176
177             except Exception as e:
178                 print(f"No 'Dangerous site' warning detected or unable to find elements: {e}")
179
180         def handle_dismiss_warning(self):
181             """Handle 'Dismiss this warning and enter site'."""
182             try:
183                 # Check if the "Dismiss this warning and enter site" button is present
184                 dismiss_button = WebDriverWait(self.driver, 3).until(
185                     EC.element_to_be_clickable((By.XPATH, "//button[contains(text(), 'Dismiss this warning and enter site')]"))
```

# You end up going from the left side results to the right side results when crawling

**FAIL**

# Stopped at Captcha

palacecirwoos.shop

Verify you are human by completing the action below.

☐ Verify you are human

CLOUDFLARE

palacecirwoos.shop needs to review the security

There are multiple products that make captcha problems easier to solve. I didn't use them in this presentation (they cost money)

Demonstrating how crawling is slow
You have to go to each page to get
the screenshot and other data

Crawl Speed

# To speed things up, I run about 20 docker-based crawlers at once, on one VM

```
josh@athena:~/multi_selenium-as-a-service-docker$ docker-compose ps
                    Name                                    Command              State                        Ports
----------------------------------------------------------------------------------------------------------------------------------
1a01a7abf2e4_multi_selenium-as-a-service-docker_chrome_1   /opt/bin/entry_point.sh      Exit 143
multi_selenium-as-a-service-docker_api_1                   flask run                    Up        0.0.0.0:5002->5000/tcp,:::5002->5000/tcp
multi_selenium-as-a-service-docker_api_10                  flask run                    Up        0.0.0.0:5015->5000/tcp,:::5015->5000/tcp
multi_selenium-as-a-service-docker_api_11                  flask run                    Up        0.0.0.0:5004->5000/tcp,:::5004->5000/tcp
multi_selenium-as-a-service-docker_api_12                  flask run                    Up        0.0.0.0:5008->5000/tcp,:::5008->5000/tcp
multi_selenium-as-a-service-docker_api_13                  flask run                    Up        0.0.0.0:5019->5000/tcp,:::5019->5000/tcp
multi_selenium-as-a-service-docker_api_14                  flask run                    Up        0.0.0.0:5007->5000/tcp,:::5007->5000/tcp
multi_selenium-as-a-service-docker_api_15                  flask run                    Up        0.0.0.0:5012->5000/tcp,:::5012->5000/tcp
multi_selenium-as-a-service-docker_api_16                  flask run                    Up        0.0.0.0:5014->5000/tcp,:::5014->5000/tcp
multi_selenium-as-a-service-docker_api_17                  flask run                    Up        0.0.0.0:5013->5000/tcp,:::5013->5000/tcp
multi_selenium-as-a-service-docker_api_18                  flask run                    Up        0.0.0.0:5018->5000/tcp,:::5018->5000/tcp
multi_selenium-as-a-service-docker_api_19                  flask run                    Up        0.0.0.0:5001->5000/tcp,:::5001->5000/tcp
multi_selenium-as-a-service-docker_api_2                   flask run                    Up        0.0.0.0:5009->5000/tcp,:::5009->5000/tcp
multi_selenium-as-a-service-docker_api_20                  flask run                    Up        0.0.0.0:5010->5000/tcp,:::5010->5000/tcp
multi_selenium-as-a-service-docker_api_3                   flask run                    Up        0.0.0.0:5005->5000/tcp,:::5005->5000/tcp
multi_selenium-as-a-service-docker_api_4                   flask run                    Up        0.0.0.0:5016->5000/tcp,:::5016->5000/tcp
multi_selenium-as-a-service-docker_api_5                   flask run                    Up        0.0.0.0:5006->5000/tcp,:::5006->5000/tcp
multi_selenium-as-a-service-docker_api_6                   flask run                    Up        0.0.0.0:5011->5000/tcp,:::5011->5000/tcp
multi_selenium-as-a-service-docker_api_7                   flask run                    Up        0.0.0.0:5017->5000/tcp,:::5017->5000/tcp
multi_selenium-as-a-service-docker_api_8                   flask run                    Up        0.0.0.0:5003->5000/tcp,:::5003->5000/tcp
multi_selenium-as-a-service-docker_api_9                   flask run                    Up        0.0.0.0:5000->5000/tcp,:::5000->5000/tcp
multi_selenium-as-a-service-docker_hub_1                   /opt/bin/entry_point.sh      Up (healthy)  4442/tcp, 4443/tcp, 0.0.0.0:4444->4444/tcp,
```
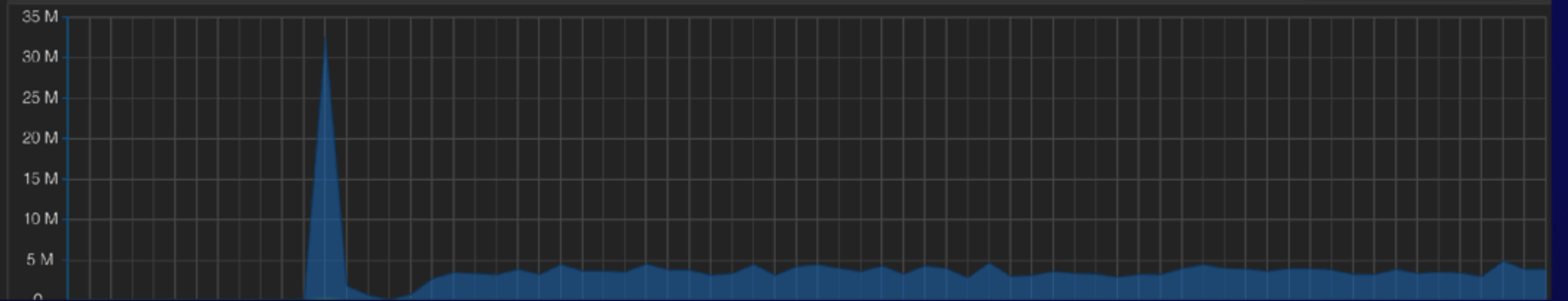
It doesn't use much of the VMs resources

1 VM with 20 crawlers is not enough. I want them all over, but where?

But it's not enough. I want

# Even More Crawlers!

In the right places

# Crawler Distribution

As an example, I took the phishtank dataset:
I collected 'all online and verified' phishes, got the domain out of the URL, and looked up the A record of the domain.

Then, I get the latitude and longitude of the IP.

PhishTank® Out of the Net, into the Tank.

Register | Forgot Password   Sign In

**Found a phishing site?** Get started now — see if it's in the Tank:

http://   Is it a phish?

**What is PhishTank?**

PhishTank is a collaborative clearing house for data and information about phishing on the Internet. Also, PhishTank provides an open API for developers and researchers to integrate anti-phishing data into their applications at no charge.
Read the FAQ...

**Recent Submissions**

You can help! Sign in or register (free! fast!) to verify these suspected phishes.

| ID | URL | Submitted by |
|----|-----|--------------|
| 8861148 | https://movil-ing-area.com/es/c/user_673d56b3be588... | BPhv |
| 8861147 | https://hjuuouvhhjjh.weebly.com/ | titus |
| 8861146 | https://www.es-ing-aviso.com | BPhv |
| 8861145 | https://magenta063943.studio.site/ | titus |
| 8861144 | https://jkdhfkjldfff.weebly.com/ | CertSecurity |
| 8861142 | https://sogiy20060.wixsite.com/home | titus |
| 8861141 | https://currentlyattfoldersrenews89000000.weebly.c... | titus |
| 8861139 | https://currentlyhomesattfolderviwes00989900.weebl... | titus |
| 8861136 | https://stem12ya.weeblysite.com/ | titus |
| 8861135 | https://ddpd.86231548.xyz/s0tqv3fi/TeaxRp/7 | Clutter |
| 8861134 | https://u.updatetrackeys.top/l/ | titus |
| 8861133 | https://vjnted.26231548.xyz/s0tqv3fi/TeaxRp/7 | Clutter |
| 8861132 | https://ing-dirct.com | BPhv |
| 8861131 | https://u.updatetrackkgp.top/l/ | titus |
| 8861130 | https://nl-template-bakker-1732014658539.onepage.w... | verifirum |

See more suspected phishes...

**New to PhishTank?**

Subscribe to the PhishTank mailing lists.

{'description': 'Unknown City, United States (Domain: q-r.to)',
 'lat': 37.751,
 'lon': -97.822},
{'description': 'Kansas City, United States (Domain: replit.app)',
 'lat': 39.1027,
 'lon': -94.5778},
{'description': 'Unknown City, United States (Domain: '
                'firebaseapp.com)',
 'lat': 37.751,
 'lat': 1.2808,
 'lon': 103.8503},
{'description': 'Unknown City, United States (Domain: '
                'google.com)',
 'lat': 37.751,
 'lon': -97.822},
{'description': 'Boardman, United States (Domain: ngrok.app)',
 'lat': 45.8234,
 'lon': -119.7257},
{'description': 'Unknown City, France (Domain: ovh.ca)',
 'lat': 48.8582,
 'lon': 2.3387},
{'description': 'Unknown City, Australia (Domain: dreamwp.com)',
 'lat': -33.494,
 'lon': 143.2104},
{'description': 'Singapore, Singapore (Domain: grefghdf.com)',
 'lat': 1.2868,
 'lon': 103.8503]]

And I map those locations (this is in the code at the end)

And get statistics on where most of the domains hosting phishing URLs are located

USA in the lead… *(cries.jpg)*

United States: 2813 URLs
Germany: 27 URLs
Hong Kong: 25 URLs
Russia: 16 URLs
France: 15 URLs
The Netherlands: 13 URLs
Brazil: 13 URLs
United Kingdom: 11 URLs
Japan: 9 URLs
Canada: 9 URLs
Australia: 8 URLs
India: 7 URLs
Indonesia: 6 URLs
Poland: 5 URLs
Spain: 5 URLs
Italy: 5 URLs
Singapore: 4 URLs
Türkiye: 4 URLs
South Africa: 4 URLs
Bangladesh: 4 URLs
Thailand: 4 URLs
Czechia: 4 URLs
Mauritius: 3 URLs
Malaysia: 3 URLs
Argentina: 3 URLs
China: 3 URLs
Finland: 3 URLs
Bulgaria: 2 URLs
Taiwan: 2 URLs
Iran: 2 URLs
Vietnam: 2 URLs
Belgium: 2 URLs
Portugal: 2 URLs
Nepal: 1 URLs
Slovakia: 1 URLs
Peru: 1 URLs
Switzerland: 1 URLs
South Korea: 1 URLs
British Virgin Islands: 1 URLs
Ukraine: 1 URLs
Hungary: 1 URLs

Quick sidebar about staying hidden as a threat hunter

Staying hidden?

**America's Cyber Defense Agency**
NATIONAL COORDINATOR FOR CRITICAL INFRASTRUCTURE SECURITY AND RESILIENCE

Search

Topics ⌄   Spotlight   Resources & Tools ⌄   News & Events ⌄   Careers ⌄   About ⌄

Home / News & Events / Cybersecurity Advisories / Alert

**Generally, you don't want to touch systems from locations that can be attributed to you or your organization.**

ALERT

# Threat Actors Targeting Cybersecurity Researchers

**Last Revised:** April 14, 2021

Google ↗ and Microsoft ↗ recently published reports on advanced persistent threat (APT) actors targeting cybersecurity researchers. The APT actors are using fake social media profiles and legitimate-looking websites to lure security researchers into visiting malicious websites to steal information, including exploits and zero-day vulnerabilities. APT groups often use elaborate social engineering and spear phishing schemes to trick victims into running malicious code through malicious links and websites.

There are services that provide access to distributed proxies
that you can use for crawling



They cost money (I didn't use them for this presentation)

# Your own Proxy

## Privoxy - Home Page

Privoxy is a non-caching web proxy with advanced fil~~...~~
controlling access, and removing ads and other obnoxious Internet junk. Privoxy has a flexible configuration and can be customized to suit individual needs and tastes. It has application for both stand-alone systems and multi-user networks.

Privoxy is Free Software and licensed under the GNU GPLv2 or later.

Privoxy is an associated project of Software in the Public Interest (SPI).

Helping hands and donations are welcome:

- https://www.privoxy.org/participate
- https://www.privoxy.org/donate

The most recent release is 3.0.34 (stable).

**But I did set up my own proxies in various networks**

Docker-Compose

```
privoxy:
  build:
    context: ./privoxydocker
  container_name: privoxy
  restart: always
  volumes:
    - ./privoxydocker/logs/privoxy:/var/log/privoxy
    - ./privoxydocker/privoxy.conf:/etc/privoxy/config:ro
  ports:
    - "1080:1080"
  dns:
    - 192.168.1.5
    - 192.168.1.3
```

Now that crawling is more or less sorted,
let's move onto creating the malicious dataset
that will be used to compare unknown URLs against.

# Creating Datasets
## Building a Malicious Dataset

# Data Sources

## Some are better than others

# URLhaus
## Community Sourced
## Malware Activity

There are lots of datasources. You can use your own corpus of known-bad URLs, or you can use free services (or you can pay for feeds).
This example shows using URLHaus

You can often grab a CSV, json, or other file with URLs

| generate thumbnail | http://p6.zbjimg.com/task/2010-12/03/519808/4cf8... |
| --- | --- |

Total Bandwidth Used: 0.02 KB

# Summary of HTTP Status Codes

- HTTP 404: 9405 occurrences
- HTTP 200: 4 occurrences

This is a malware feed

# Other feeds are better, such as urlscan.io and phishintank phishing feeds

## URLScan.io



PhishTank
Community Sourced

Once you have your list of URLs and you've crawled them to get screenshots, you may notice that a lot of screenshots are not worthy of being used as malicious screenshot to compare other screenshots to. At this point, you have to go through the screenshots, removing anything that doesn't apply.

# Removing non-relevant



In this example of screenshots, some are good for a malicious dataset, but you can see that others aren't so great.

To quickly remove screenshots that won't work, I will introduce some of the main ideas behind locating similar screenshots. The information in this portion of the presentation will be used in various ways throughout the many different similarity analysis Techniques.

# Distance & Similarity

The primary algorithm and technique behind all similarity detection is Levenshtein. I'll be demonstrating using it as well as various hashing comparison mechanisms that use it under the hood of their own processes.

# Levenshtein Distance

Levenshtein is used best on short strings of text.
In this example, how many changes to you need to make
to turn the word 'bats' into the word 'cats'. It's just one letter, so
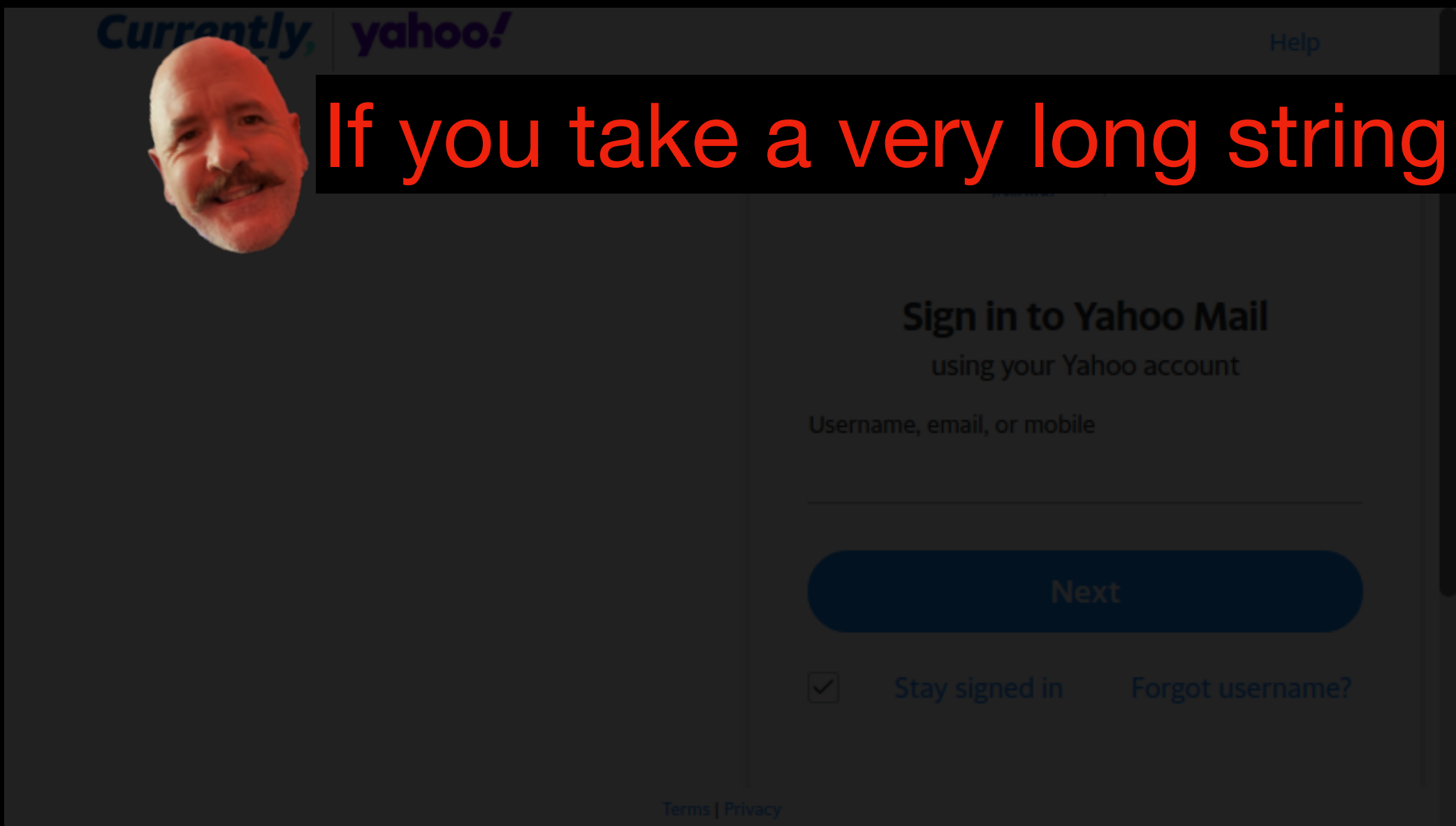this would be a Levenshtein distance of 1
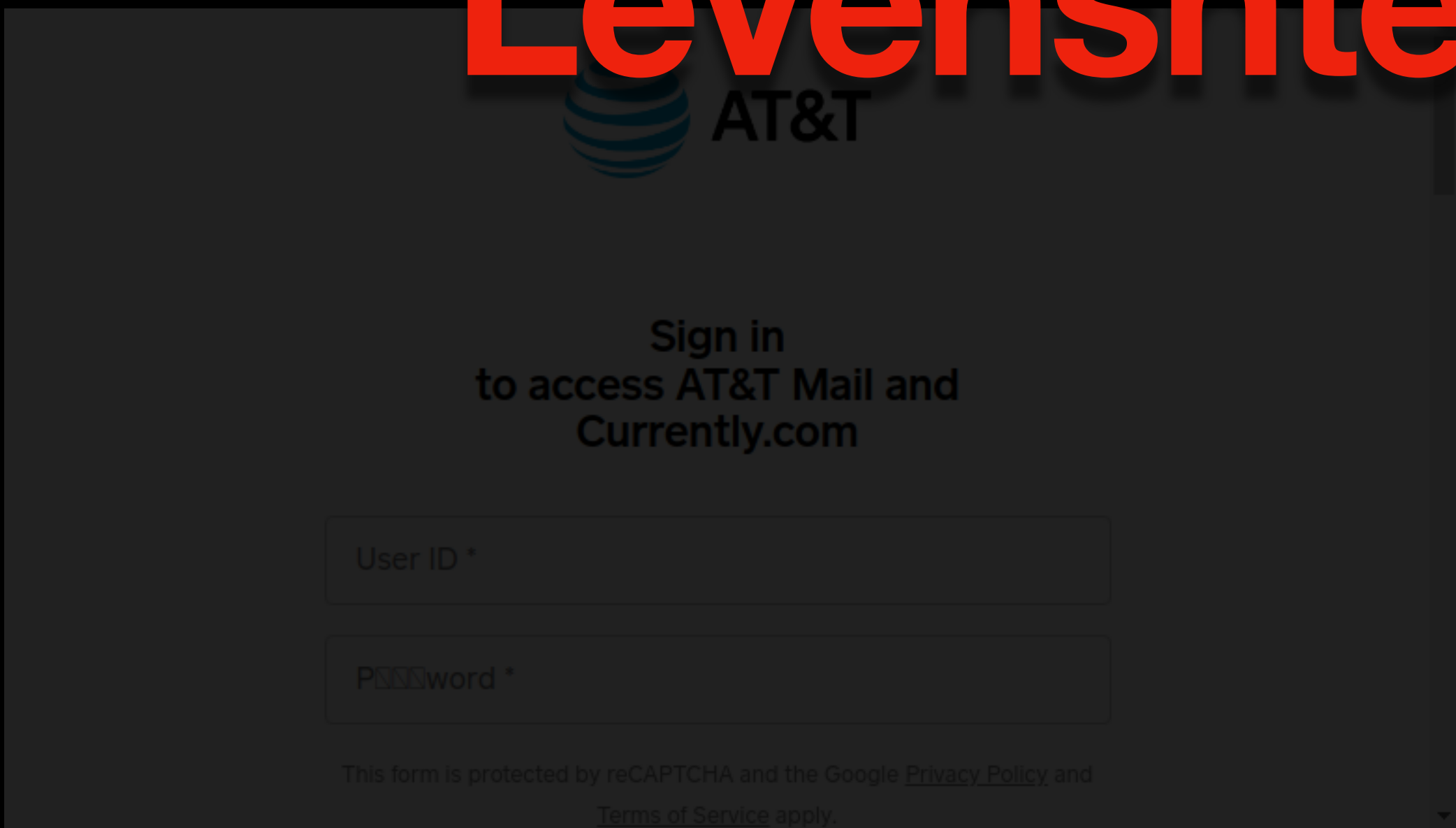


BATS

CATS

Levenshtein Distance of 1

If you take a very long string, you will get a much larger Levenshtein distance.

## Sign in to Yahoo Mail
using your Yahoo account

Username, email, or mobile

Next

☑ Stay signed in    Forgot username?

Terms | Privacy

yahoo.com gmail.com outlook.com aol.com
Stay signed in
Forgot username?
Create an account
Yahoo makes it easy to enjoy what matters most in your world.
Best in class Yahoo Mail, breaking local, national and global news, finance, sports, music, movies and more
Enter password to finish sign in
Password Next
TermsPrivacy
x
Yahoo works best with the latest versions of the browsers. You're using an outdated or unsupported browser
browser version now.
More Info

# Levenshtein Distance of 717

AT&T

Sign in
to access AT&T Mail and
Currently.com

User ID *

P░░░word *

This form is protected by reCAPTCHA and the Google Privacy Policy and
Terms of Service apply.

Processing file: ./009-288-49.weeblysite.com.txt
Home | 009-288-49
Shopping Cart
You don't have any items in your cart.
Checkout
Continue Shopping
Accepted here
Sign into access AT&T Mail andCurrently.com
This form is protected by reCAPTCHA and the Google Privacy Policy and Terms of Service apply.
Sign in
Back to Cart
009-288-49
Secure checkout by Square

If you take a shorter string, like these two URLs, you get a smaller one

http://61.52.12.185:54720/bin.sh

🔥 Online (spreading malware for 14 minutes)

61.52.12.185

2024-10-21 05:56:05 UTC

🐛 Malware download

http://61.52.12.185:54720/bin.sh

# Levenshtein Distance of 14

http://61.52.95.103:33803/i

🔥 Online (spreading malware for 16 minutes)

61.52.95.103

2024-10-21 05:54:06 UTC

🐛 Malware download

http://61.52.95.103:33803/i

Before continuing with Levenshtein, I want to talk about a hashing algorithm designed for images called PHASH, then we will apply Levenshtein to it.



Using

# PHASH
To Classify Images

Using
**PHASH**
To Classify Images

These are NOT the same
but they are similar

DEEPSEC

PHASH is different than SHA256.
SHA256 makes a longer hash, but a very small
change in the hash means a drastically different
file.

SHA256:
84926feadda7dfeb59777ab0e4b4cc60a977b6e1

SHA256:
3b4cd30a1f1d85b28601641afde04375169854f8

PHASH is different than SHA256.
SHA256 makes a longer hash, but a very small change in the hash means a drastically different file.

84926feadda7dfeb59777ab0e4b4cc60a977b6e1
3b4cd30a1f1d85b28601641afde04375169854f8

SHA256:
84926feadda7dfeb59777ab0e4b4cc60a977b6e1

Levenshtein Distance

36

SHA256:
3b4cd30a1f1d85b28601641afde04375169854f8

PHASH is shorter, but a difference in the hashes don't suggest a drastic change in the images.



PHASH:
aa95c1d5d595d485



PHASH:
aad595c1d595d590

PHASH is shorter, but a difference in the hashes don't suggest a drastic change in the images. Levenshtein distance can be applied to these shorter hashes with promising results.

aa95c1d5d595d485
aad595c1d595d590

PHASH:
aa95c1d5d595d485

Levenshtein Distance

7

PHASH:
aad595c1d595d590

These two images are obviously different, but they have some similarities. The PHASH's are different with a lev distance of 14

e59832669b9ac699

Levenshtein: 14

b3517399334ccccc

Showing multiple images and the lev distances

Levenshtein: 14

e59832669b9ac699

aa95c1d5d595d485

Levenshtein: 15

Levenshtein: 14

Levenshtein: 14

Levenshtein: 7

Levenshtein: 16

b3517399334ccccc

aad595c1d595d590

The images that look the most alike, but which are still different have a lev distance of 7. I can use this analysis to start setting a threshold.

e59832669b9ac699

aa95c1d5d595d485

Levenshtein: 15

# **Setting a Threshold**

Levenshtein: 7

Levenshtein: 16

b3517399334ccccc

aad595c1d595d590

Now I can get back to removing non-relevant images in the malicious dataset. I'm going to group them using a threshold and delete non-relevant images as groups.

# Grouping Images

## Using the threshold

To demonstrate what grouping all the crawled screenshots looks like, I'll be taking a random bunch of screenshots like this:

**Group 1**

**Group 2**

**Group 3**

**Group 4**

Outlook

**Group 5**

**Group 6**

And group similar together.

signup-swellnetwork.net

advertisedszp.shop

quialitsuzoxm.shop

DocuSign

Here's the initial process of grouping that I did using Python.
It just output a text file with the file locations that matched a group
all together.
However, all images are in groups, including images that don't
look like anything else. You can see that I have 1,507 groups.

Group 1:
../malicio
../malicio
../malicious dataset/urlscan_output_20240906_2335/cash67.com.jpg: 9a5ad69bea687848
../malicious dataset/urlscan_output_20240906_2335/getmoney3.com.jpg: 9a5ad69bea687848
../malicious dataset/urlscan_output_20240906_2335/www.findcash7.com.jpg: 9a5ad69bea687848
../malicious dataset/urlscan_output_20240906_2335/www.cash67.com.jpg: 9a5ad69bea687848
../malicious dataset/urlscan_output_20240906_2335/bigwealth7.com.cash67.com.jpg: 9a5ad69bea687848
../malicious dataset/urlscan_output_20240906_2335/www.getmoney3.com.cash67.com.jpg: 9a5ad69bea687848
../malicious dataset/urlscan_output_20240906_2335/www.onlinework7.com.jpg: 9a5ad69bea687848

Group 2:
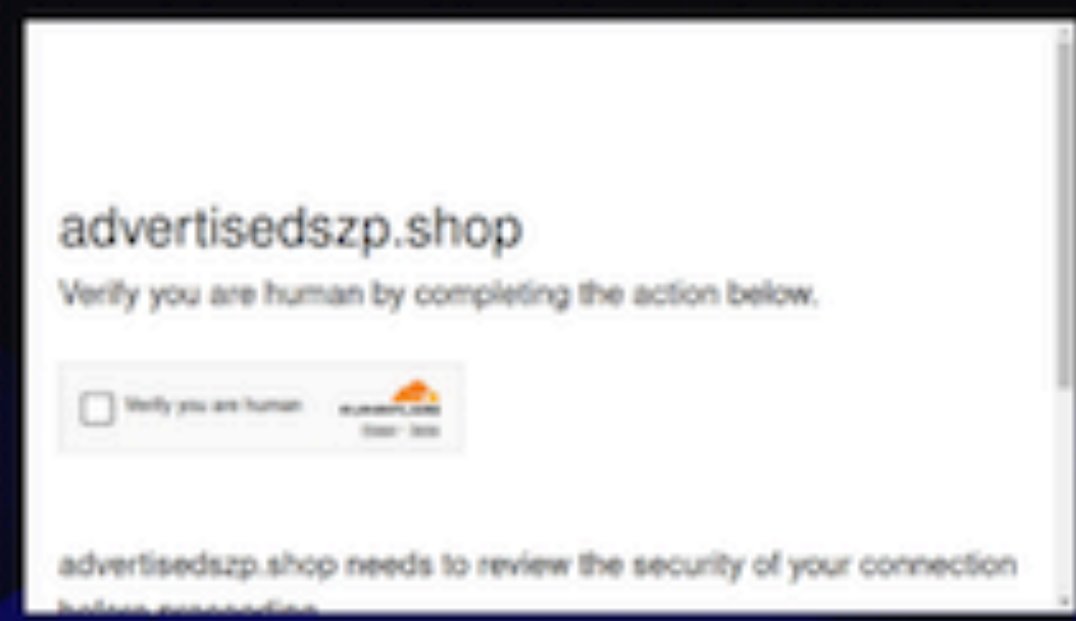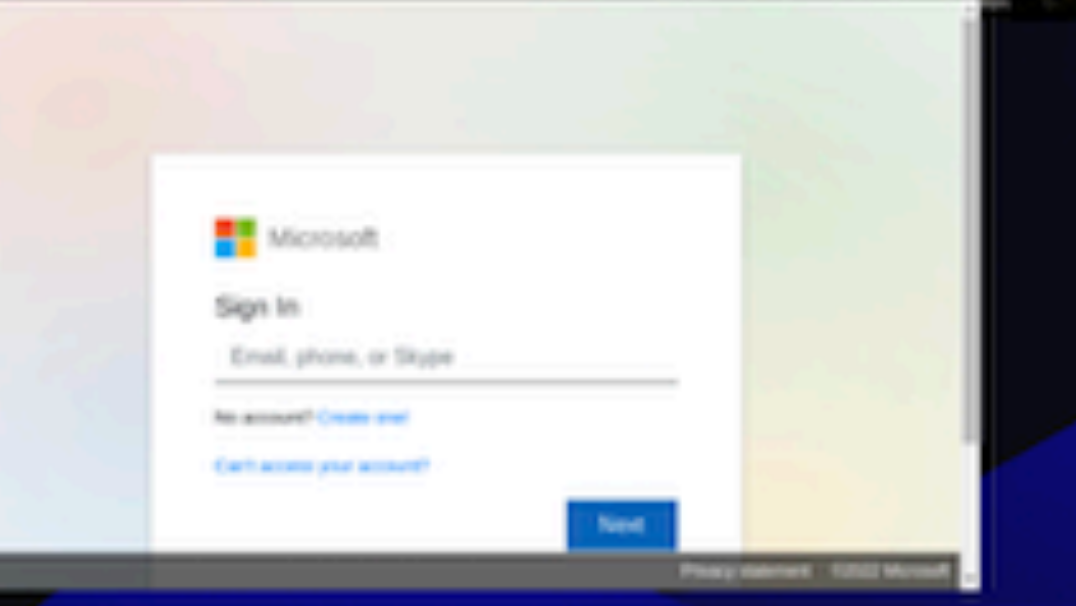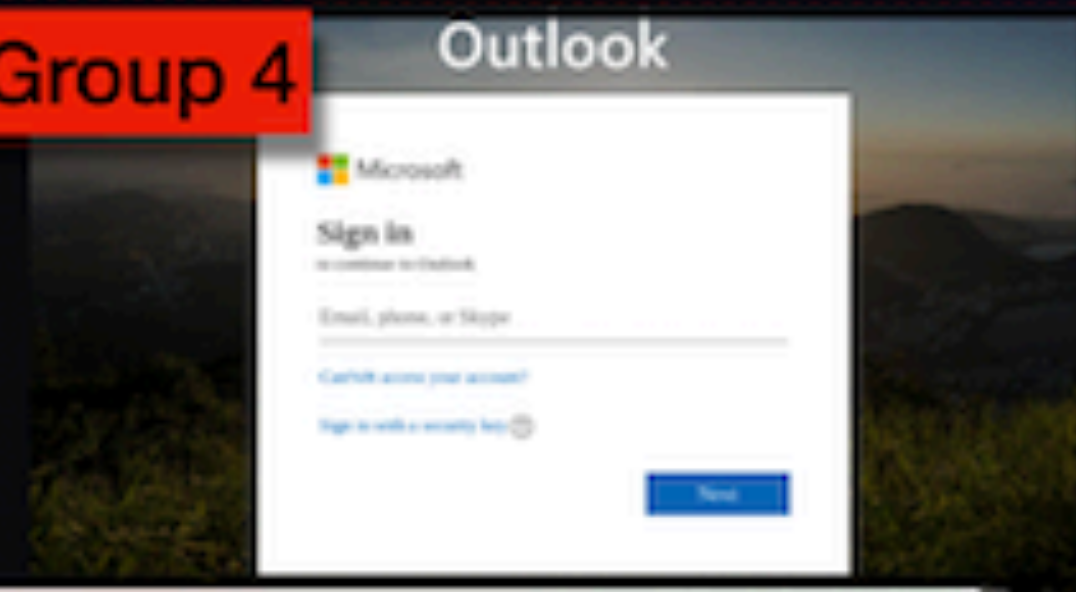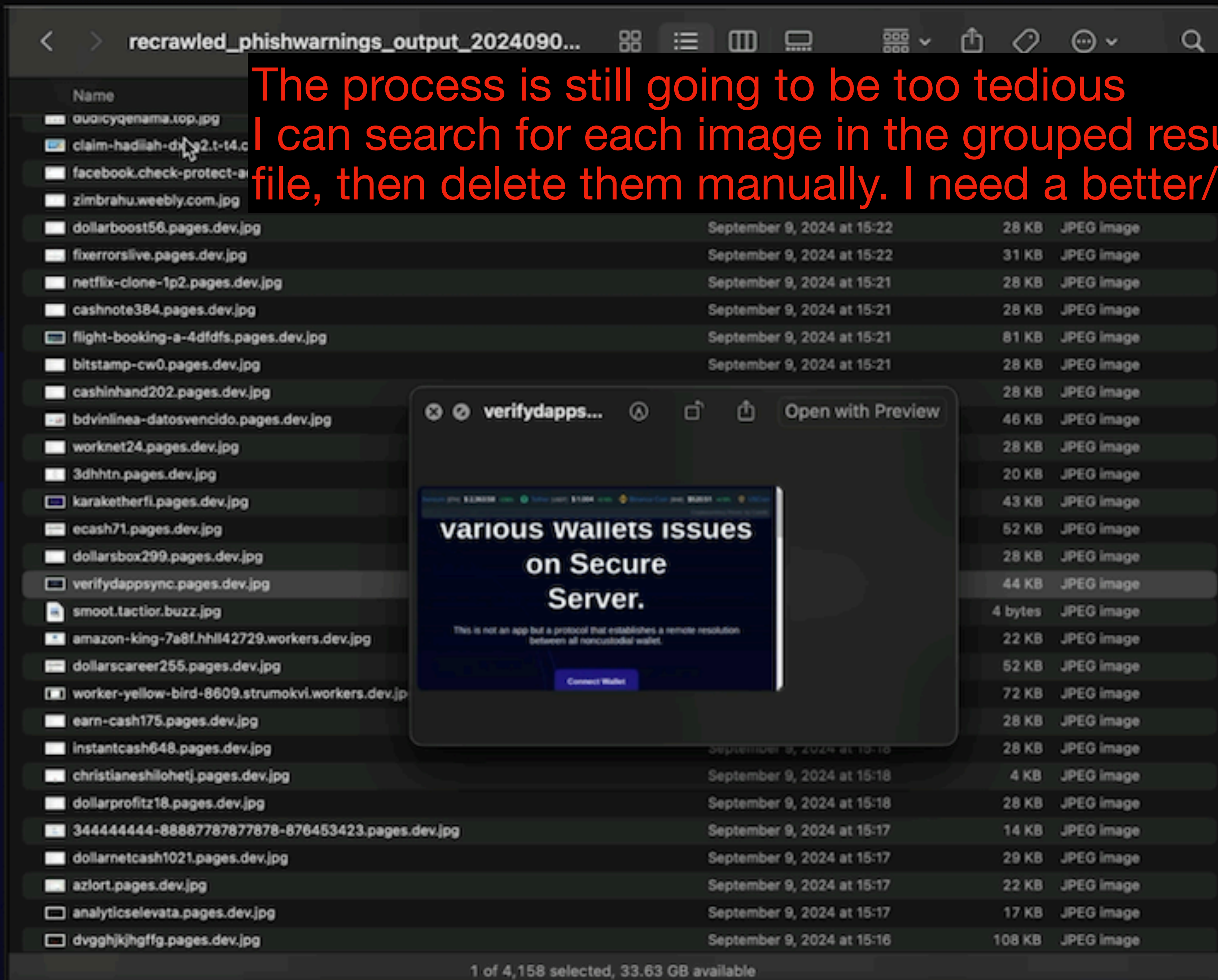../malicious dataset/urlscan_output_20240906_2335/fafafa0322.com:8989.jpg: f056d77125a98574

Group 3:
../malicious dataset
../malicious dataset
../malicious dataset
../malicious dataset

Group 4:
../malicious dataset
../malicious dataset
../malicious dataset
../malicious dataset
../malicious dataset
../malicious dataset
../malicious dataset
../malicious dataset
../malicious dataset
../malicious dataset
../malicious dataset
../malicious dataset
../malicious dataset/urlscan_output_20240906_2335/virginmoneyhelpchat.com.jpg: 97f0e00778f80f47
../malicious dataset/urlscan_output_20240906_2335/multicoinresolve.pages.dev.jpg: 97f0e00778f80f47
../malicious dataset/urlscan_output_20240906_2335/aib-mobileservices.com.jpg: 97f0e00778f80f47
../malicious dataset/urlscan_output_20240906_2335/palacecirwoos.shop.jpg: 97f0e00778f80f47
../malicious dataset/urlscan_output_20240906_2335/earn-cash204.pages.dev.jpg: 97f0e00778f80f47
../malicious dataset/urlscan_output_20240906_2335/easylifepro78.pages.dev.jpg: 97f0e00778f80f47
../malicious dataset/urlscan_output_20240906_2335/pdffilesinv.pages.dev.jpg: 97f0e00778f80f47
../malicious dataset/urlscan_output_20240906_2335/clickintowealth629.pages.dev.jpg: 97f0e00778f80f47
../malicious dataset/urlscan_output_20240906_2335/dollarnetcash1477.pages.dev.jpg: 97f0e00778f80f47

```
Group 1500: 1
Group 1501: 1
Group 1502: 1
Group 1503: 1
Group 1504: 1
Group 1505: 1
Group 1506: 1
Group 1507: 1
Total number of all items in the groups: 8666
```

The process is still going to be too tedious
I can search for each image in the grouped results text file, then delete them manually. I need a better/faster way

# The JavaScript library for bespoke data visualization

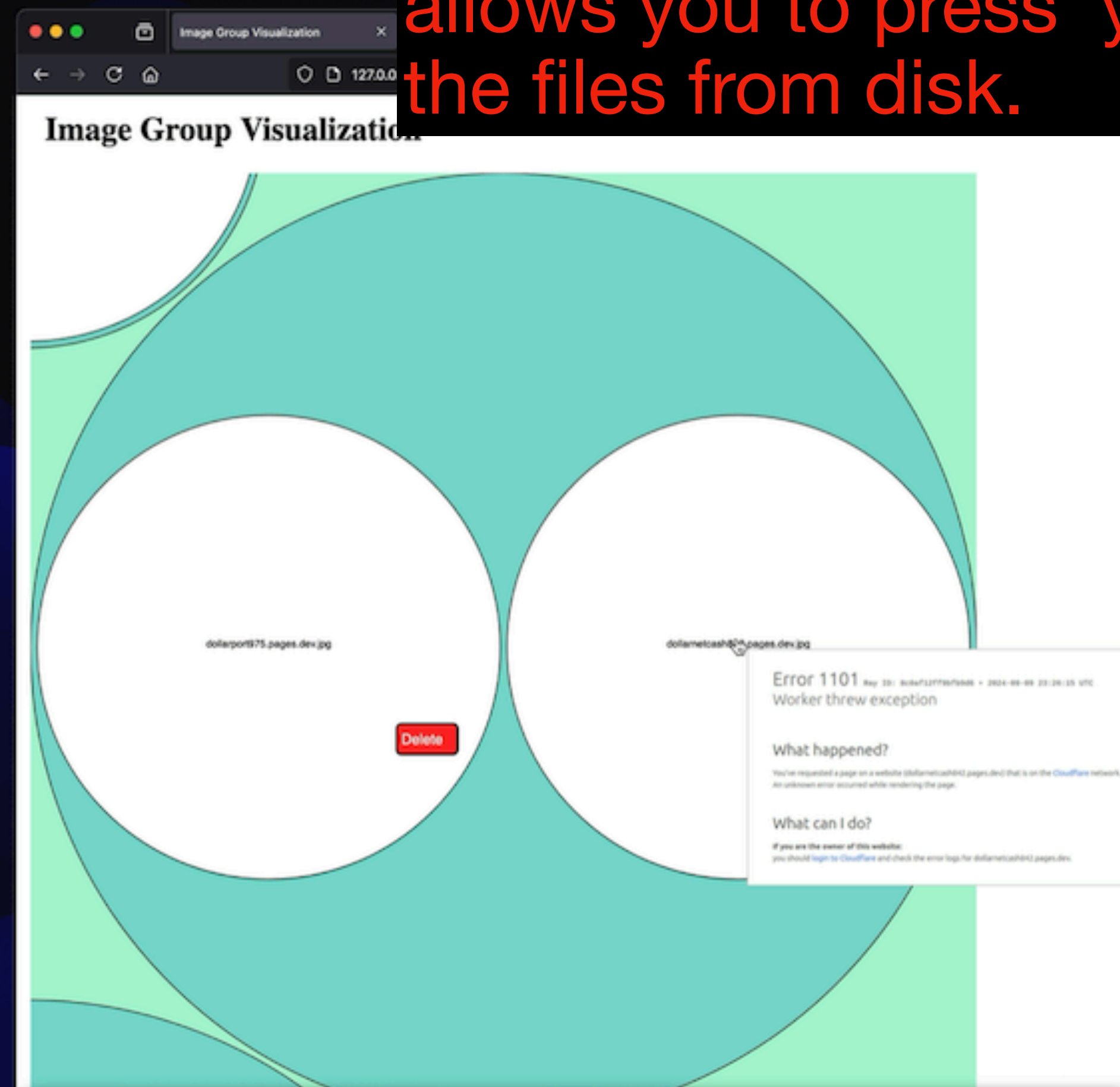Create custom dynamic visualizations with unparalleled flexibility

Get started     What is D3?     Examples

**I'm going to use D3 to make a visualization that allows me to quickly work with groups of images.**

Grouping a small test set
I added a delete button.
The version in my code (at the end)
allows you to press 'y' or 'n' to keep or delete
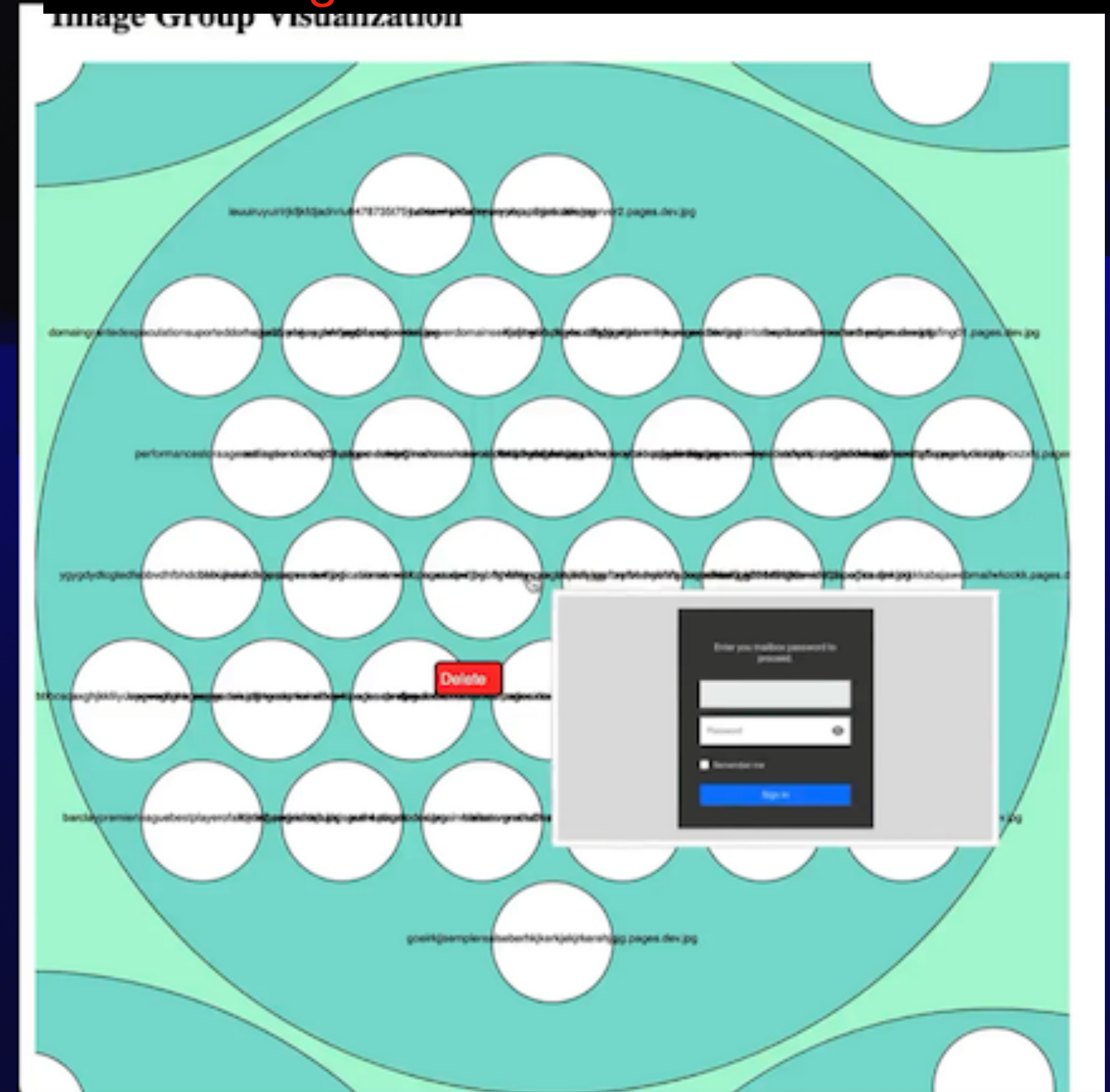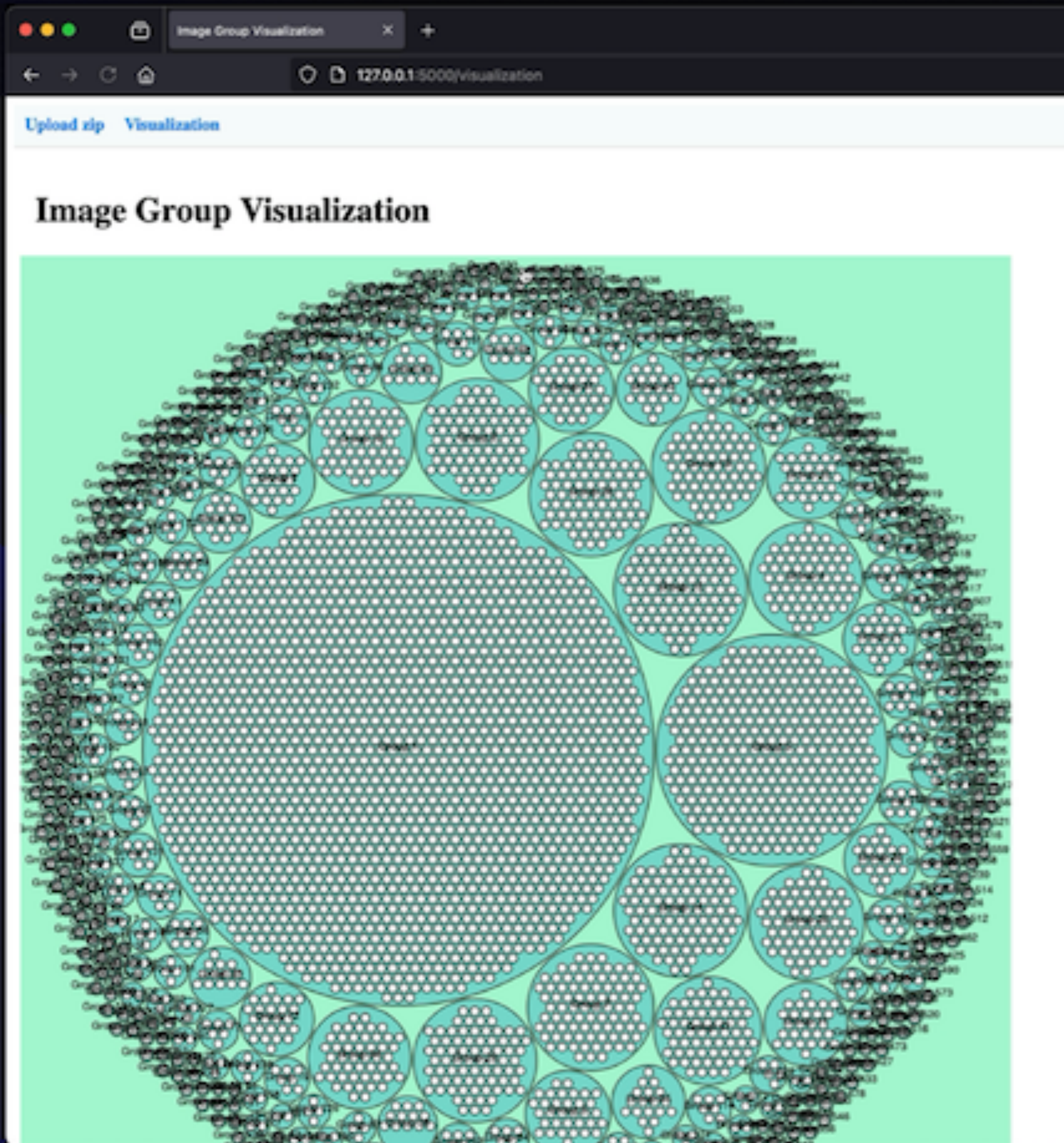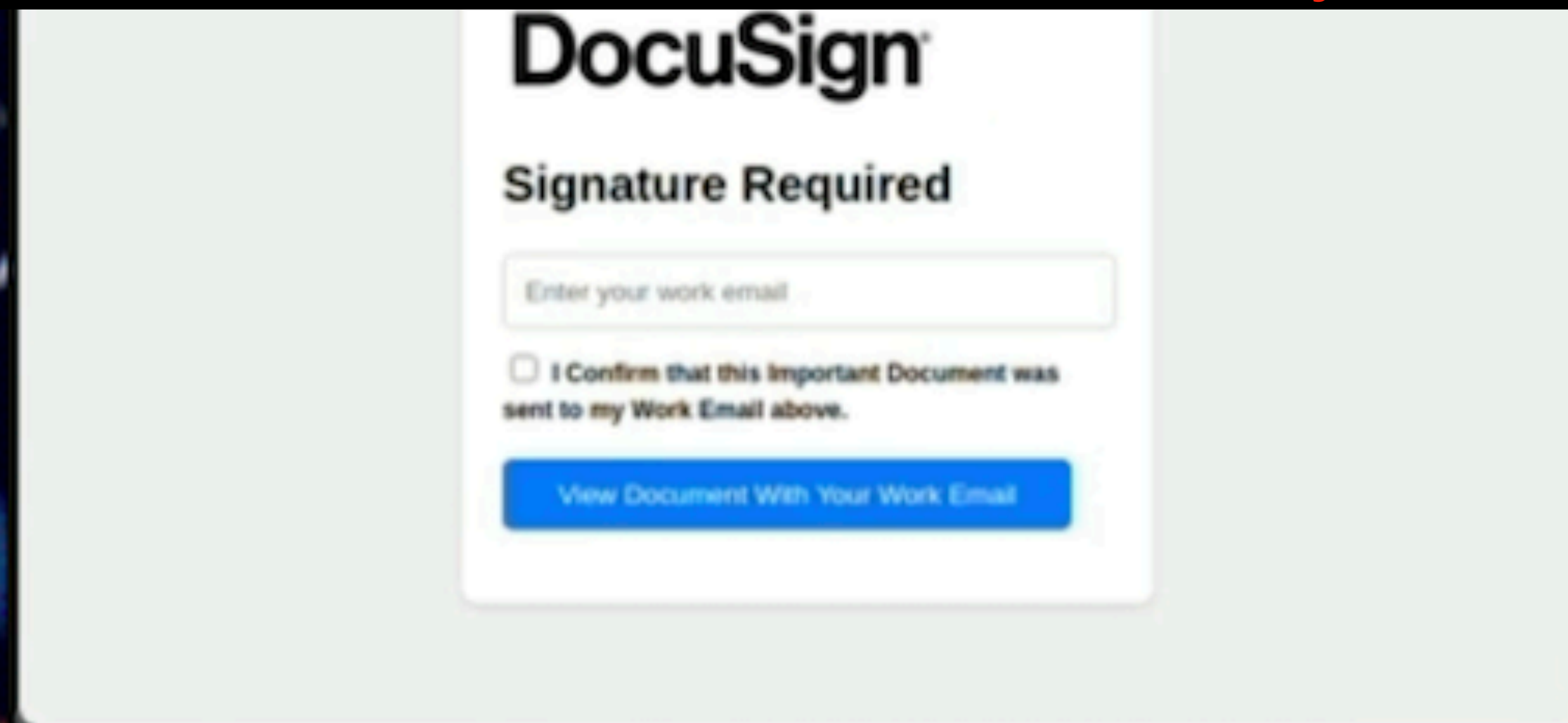the files from disk.

Now I can throw all the images in

# Let's build the malicious dataset!

The big circle is the most commonly seen similar images grouped together. Then, the other circles are groups.

I can click into each circle, hover over the individual circles, see the images, then delete the entire group if it's not a good screenshot.

Eventually, it's done and I'm left with a bunch of only malicious images!

# DocuSign®

## Signature Required

Enter your work email

☐ I Confirm that this Important Document was sent to my Work Email above.

View Document With Your Work Email

■ 344444444-8888778787878-876453423.pages.dev.jpg
■ 685845093490349034-489567457478478.pages.dev.jpg
■ 8?5587347645786894765.pa...dev.jp...

# Dataset is Populated!

■ aau.web-whatsappv.icu.jpg
■ account-service.navy-resourcesupdates.workers.dev.jpg
■ ad-link.pages.dev.jpg

Video showing a quick look at the images that are left.

■ adobevues.pages.dev.jpg
■ afro-china.pages.dev.jpg
▭ aj9scchg29929.krzecyq-iofnaskn-oasfiohnafs.workers.dev.jpg
■ aloentscheme.pages.dev.jpg
■ amzshopid.vip.jpg
■ applecare24x7help.pages.dev.jpg
■ apples-helps-third-cares.pages.dev.jpg
■ apps-9z2.pages.dev.jpg
■ apwebmail-acount-serve-update.pages.dev.jpg

# Testing Detection

Now that the malicious dataset is available, I can start sending other screenshots
as they're crawled to be checked against the dataset, finding similar images, which indicate similar phishing campaigns.
But first, I am going to test with a small set of data.

For the test set, I'm using these two images. One is the 'malicious' one, and the other
is from a URL that is similar. I want to see if, when crawled and compared, it
identifies them as similar.

URL: centrum-100492.weeblysite.com
PHASH: b51f9e6061c3c5e1

URL: attttt-104045.weeblysite.com
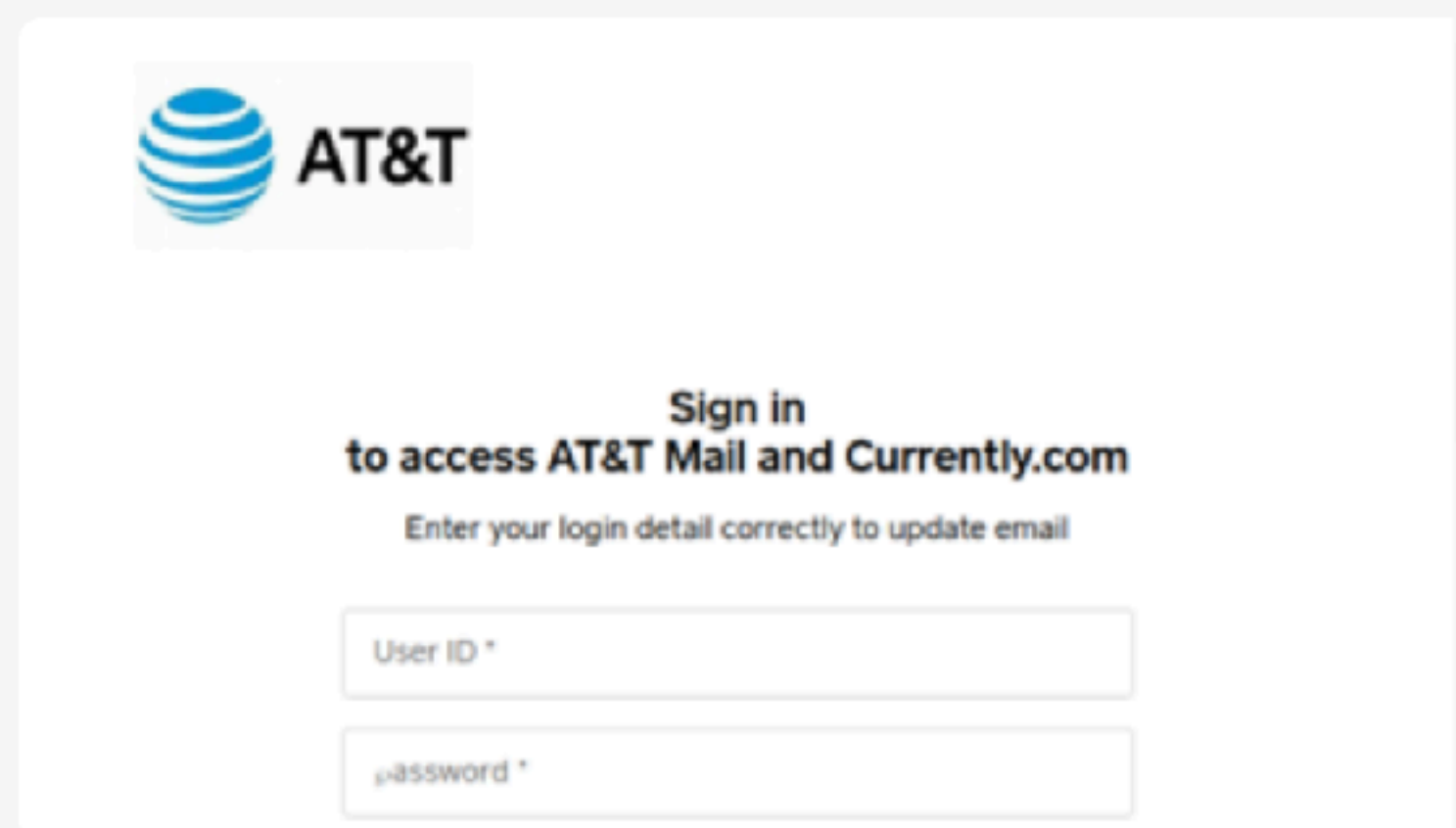PHASH: bd1f9f2d61c3c060

**Different PHASH**
**Similar Look**

# I put one image into the malicious dataset.

## Malicious Dataset

Upload ZIP File: [Browse...] No file selected.    [Upload]

| Thumbnail | URL | PHASH | A Record | City | Country |
|---|---|---|---|---|---|
|  | attttt-104045.weeblysite.com | bd1f9f2d61c3c060 | 74.115.51.54 | Unknown City | United States |

# Then, I tested two screenshots. The other one that's similar and another non-similar screenshot.

# Crawl & Detect Test

Now I'm going to actually do a crawl of a few websites, testing against the malicious image.

I've got the URL for the similar-looking site along with the URL for the screenshot that's already in the malicious dataset, and two other URLs which result in different screenshots.
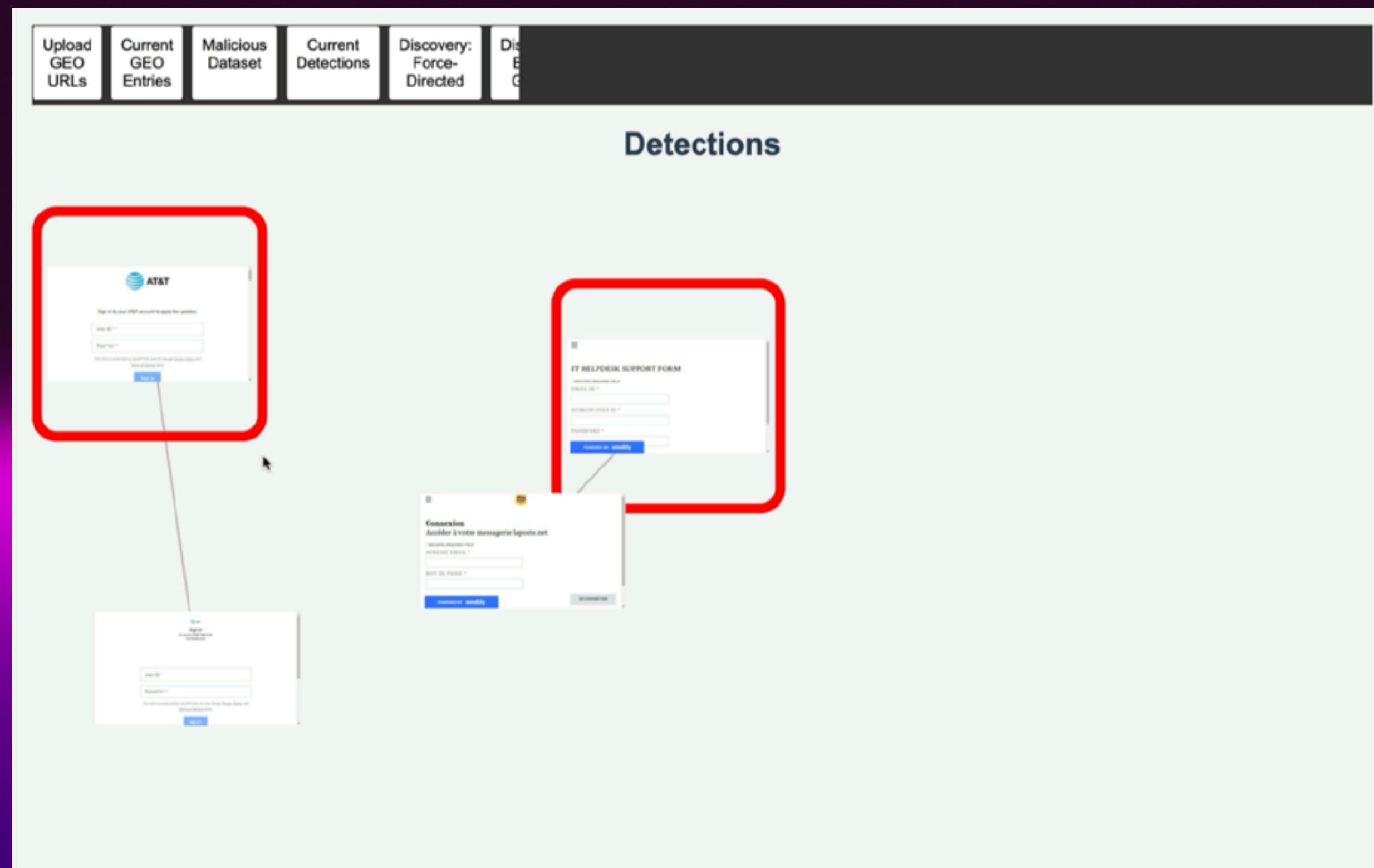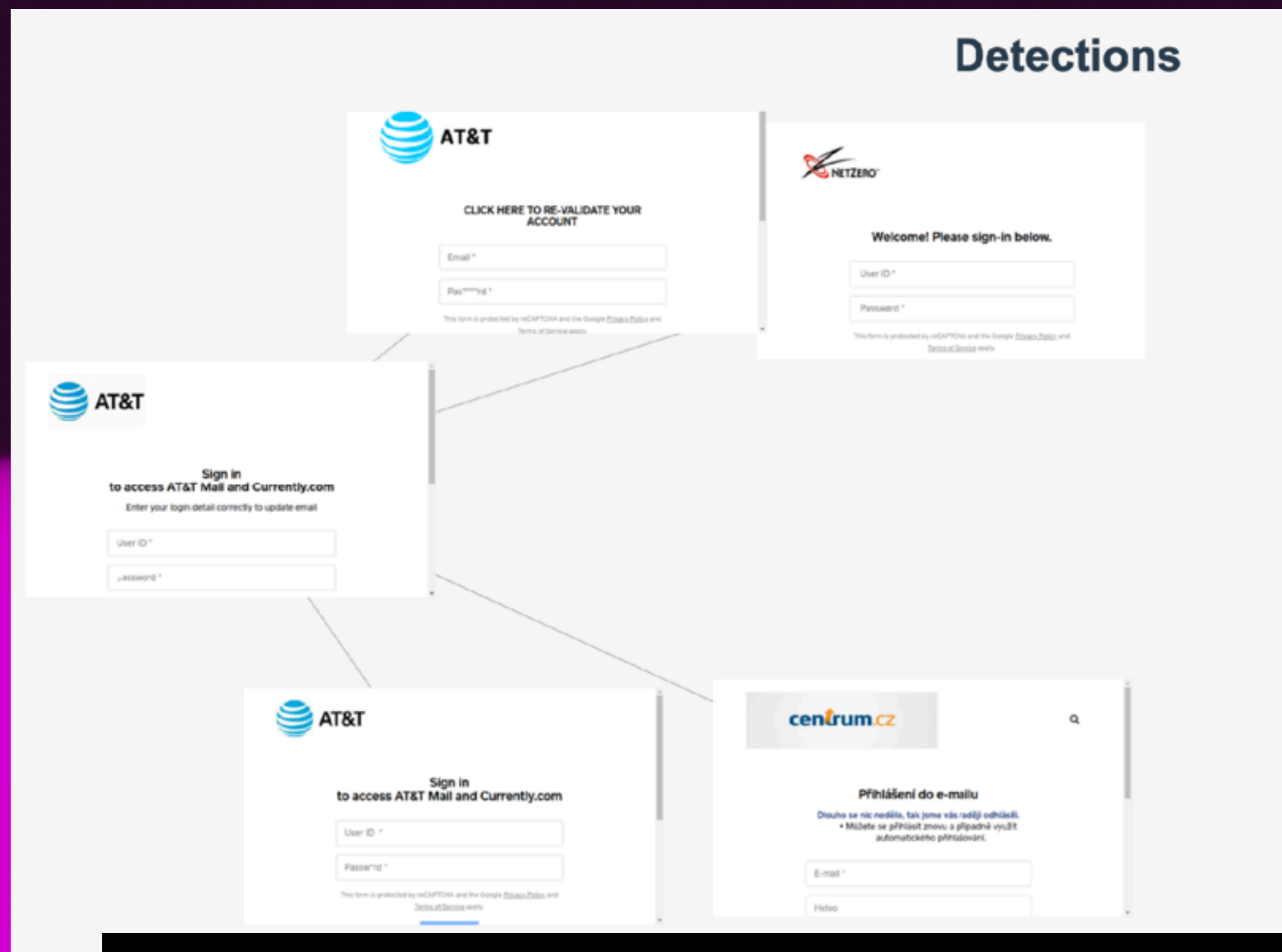The video in this slide demonstrated that only the original URL and the similar screenshot URL matched.

# Visualizing Image Similarity Detection

Now I'm going to demonstrate various methods for visualizing detections
We typically are always looking at rows and columns (log files, SIEM results, etc).
I'm interested in exploring other ways to view 'hits'

Using a force-directed graph, we can view detections



And it can be customized to make more sense, like having the malicious reference screenshot bordered with a red line or whatever (security is creative)
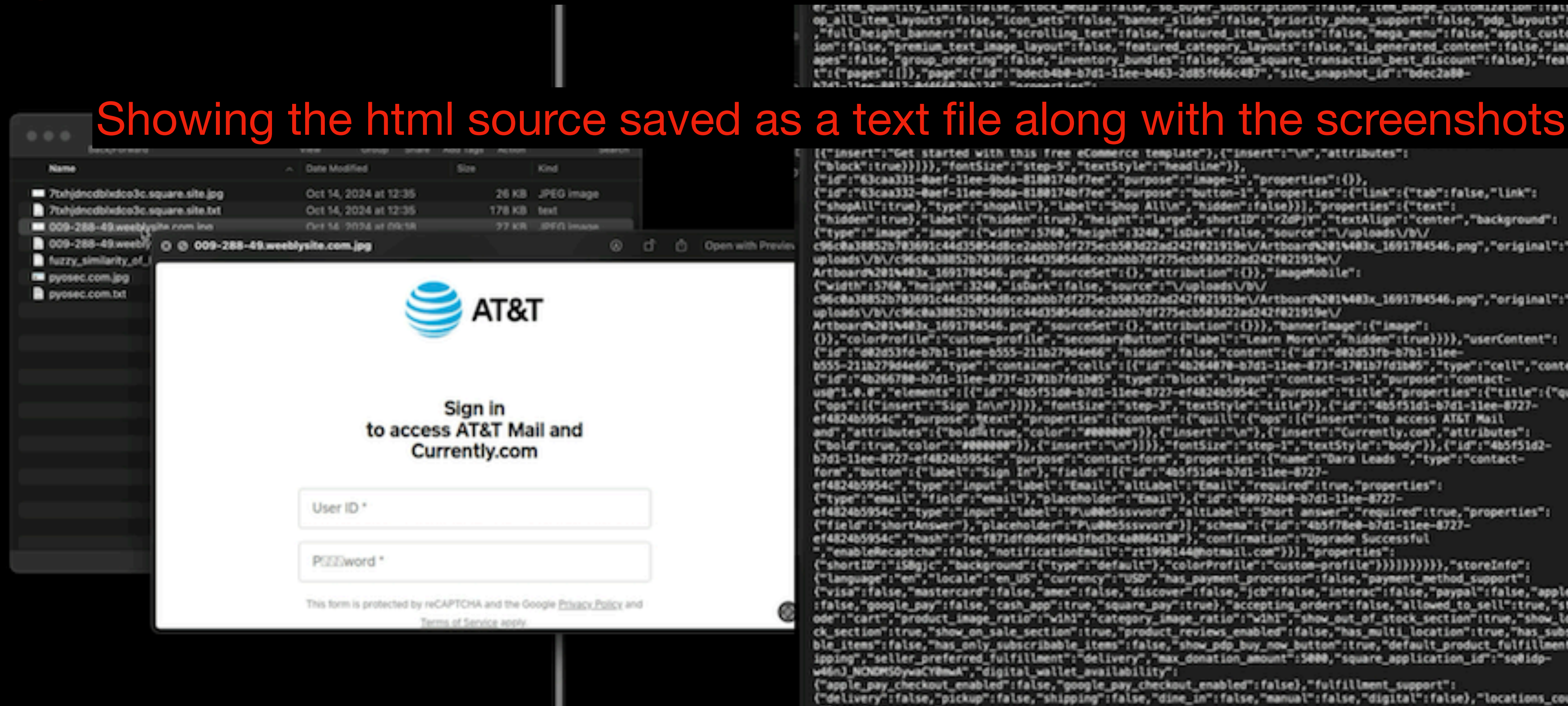
# Beyond Images

Images are good, but there's more to work with, and there are more technologies to experiment with

- **HTML Similarity**
- **OCR**
- **Using LLMs**
- **Text Content Similarity**

# HTML Similarity

When crawling, I also get the html source. In the Levenshtein section, I showed that applying Levenshtein on the entire source is not going to be very productive, but there are options…

Showing the html source saved as a text file along with the screenshots

```
→ html_comparison_test |
```

🔍 Search

🐍 fuzzy_similarity_of_html.py ✕

```python
1   from simhash import Simhash
2
3   def generate_simhash(file_path):
4       with open(file_path, 'r') as f:
5           content = f.read()
6       return Simhash(content)
7
8   # Generate Simhash for each file
9   file1_simhash = generate_simhash('7txhjdncdblxdco3c.square.site.txt')
10  # file2_simhash = generate_simhash('009-288-49.weeblysite.com.txt')
11  file2_simhash = generate_simhash('pyosec.com.txt')
12
13  distance = file1_simhash.distance(file2_simhash) # Compare the Hamming distance between the two Simhashes
```

This video shows comparing two html files with 'simhash'.
Simhash takes large amounts of text and makes a small hash out of it that is
similar to PHASH. It also has a hamming function, which basically just uses
Levenshtein to check the distance between the hashes. It's effective though at finding
files that have similar content.

This is a video showing text output of groups based off the html source code.
In the video, I search a few images that are listed within each group of html files that are similar and display the screenshots, showing that they look alike.

Malicious Item: DB: 9414069657275232887

Malicious Item: DB: 2255457264052142660

Malicious Item: DB: 980915487212584965

Malicious Item: DB: 10898285367999205273

Malicious Item: DB:

**I need a better way to show the relationships between the html files, so I put them on a timeline. The timeline also can be useful to keep track of when certain kinds of phishing occurred**

Malicious Item: DB: 14005425564567667930

Malicious Item: DB: 10610054991847526297

Malicious Item: DB: 15239755546223349657

2023-04 2023-05 2023-06 2023-07 2023-08 2023-09 2023-10 2023-11 2023-12 2024-01 2024-02 2024-03 2024-04 2024-05 2024-06 2024-07

AT&T

AT&T Sing In

15239755546223349657
2024-07-14T19:03:53
Simhash: 15239755546223349657
URLs: 1

Malicious Item: DB: 10610055000437428121

# Scraping Text from Images

## Optical Character Recognition (OCR)

Another method that can be used is to scrape the text from images and use that for similarity searching.

Here, I'm showing the text that has been scraped from the images:

| Thumbnail | URL | Image PHASH | A Record | City | Country | Added | Text |
|---|---|---|---|---|---|---|---|
| | attmailaccount22.weebly.com | ed97960396496996 | 146.112.61.108 | San Jose | United States | 2024-10-25 06:23:54 | ATT S ater Sign into access ATT Mail and Currentlyco m Indicates required field |
| | atdikg.taplink.ws | d7658c9a939a989c | 146.112.61.108 | San Jose | United States | 2024-10-25 06:23:54 | S ATT Signin to access ATT Mail and Currentlycom User ID Password |
| | | | | | | 5 | a SS ATT Sign in to access ATT Mail and Currentlycom User ID Address 9 |
| | internationaldomaierepresentationmainserstexts4.pages.dev | e59832669b9ac699 | 146.112.61.108 | San Jose | United States | 2024-10-28 05:57:30 | Enter you mailbox password to proceed Remember me |
| | ionos-webmail-105551.weeblysite.com | f7732a98cc8ccca2 | 74.115.51.54 | Unknown City | United States | 2024-10-28 05:57:31 | 1ONOS webmail login Email Address Passwrd This form is protected by reCAPTCHA and the Google Privacy Policy and |
| | papaya20688090.brizy.site | f5d320d88b7c8b2c | 146.112.61.108 | San Jose | United States | 2024-10-28 05:57:32 | VERIFICACION P Made With Brizy |
| | mail-105085.weeblysite.com | f77323331833cc8c | 74.115.51.55 | Unknown City | United States | 2024-10-28 05:57:32 | REAGAN Email Address Passwrd This form is protected by reCAPTCHA and the Google Privacy Policy and apply |

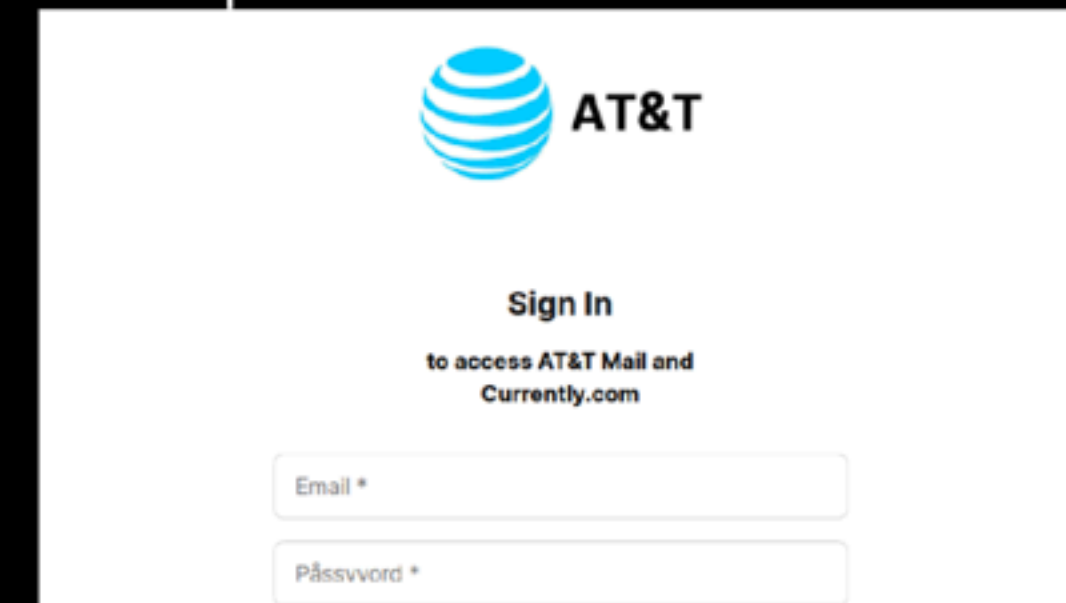Comparing the html source and the OCR text, the OCR is a bit rough. But the html has more that you have to filter out.

009-288-49.weeblysite.com

**Sign in**
to access AT&T Mail and
Currently.com

User ID *

P█████word *

This form is protected by reCAPTCHA and the Google Privacy Policy and
Terms of Service apply.

OCR: ATT Q Sign in to access ATT Mail and Currentlycom User ID

HTML: Home | 009-288-49
Shopping Cart
You don't have any items in your cart.
Checkout
Continue Shopping
Accepted here
Sign into access AT&T Mail andCurrently.com
This form is protected by reCAPTCHA and the Google Privacy Policy and Terms of Service apply.
Sign in
Back to Cart
009-288-49
Secure checkout by Square

7txhjdncdblxdco3c.square.site

**Sign In**
to access AT&T Mail and
Currently.com

Email *

Påssvvord *

OCR: SS ATT Sign In to access ATT Mail and Currentlycom Email  Passvvord

HTML: Home | .
Shopping Cart
You don't have any items in your cart.
Checkout
Continue Shopping
Accepted here
Sign In
to access AT&T Mail andCurrently.com
You may receive marketing and promotional materials. Contact the merchant for their privacy practices. This reCAPTCHA and the Google Privacy Policy and Terms of Service apply.

# Text Classification with LLMs
# 1. Easy mode

I'm going to send an entire html file and ask for a score to a locally hosted LLM that behaves much like ChatGPT (sending it via python)

A juicy piece of HTML code!

As a cybersecurity expert, I'll analyze the content and provide a phishing likelihood score based on various factors.

**Phishing Likelihood Score: 6/10**

Here's why:

1. **Unusual path data**: The SVG path data contains unusual characters and formatting, which could be an attempt to obfuscate or hide malicious code.
2. **Multiple unknown classes**: Classes like `cko__footer-text`, `payment-methods`, and others are not well-known or widely used, which might indicate a custom-made phishing template.
3. **Unusual div nesting**: The HTML structure is complex, with multiple nested div elements that could be hiding malicious content or JavaScript code.
4. **Unknown script tags**: The presence of unknown s                    eckout/js/wcko.049619101192b0140d13.js"></script>`) coul
e the presence of malicious code.
5. **Lack of transparency**: The code does not provid                    og.

However, I must note that some of these factors might                    uld need to analyze the entire website, including the s
e code, network traffic, and user interactions.

**Recommendations:**

1. **Verify the website's authenticity**: Ensure you'                    familiar URLs, and legitimate branding.
2. **Be cautious with unknown scripts**: Avoid execut
3. **Use a reputable antivirus software**: Run regula                    ode.

The response is akin to a new analyst who is overly excited and under-experienced in this work

When I remove the superfluous text commentary to keep just the score and run it again, I get a different score. This LLM is too generalized and can't produce consistent results.

I need to try using a more specific LLM

I'm going to demo zero-shot classification and a tool called OmniParser

# Text Classification with LLMs
# 2.Better modes

Zero-Shot Classification

OmniParser

🤗 Hugging Face    🔍 Search models, datasets, users...          📦 Models

---

facebook/**bart-large** 🗗    ♡ like  175

| ⊞ Feature Extraction | 🤗 Transformers | ⏻ PyTorch | ⬆ TensorFlow | 𝒥𝒳 JAX | ⓡ Rust | 🌐 English | bart |

📦 Model card    ⋮≡ Files and versions    🤗 Community 7

---

🔗 **BART (large-sized model)**

BART model pre-trained on English language. It was introduced in the paper <u>BART: Denoising Sequence-to-Sequence Pre-training for Natural Language Generation, Translation, and Comprehension</u> by Lewis et al. and first released in <u>this repository</u>.

**This LLM uses something called zero-shot classification to see how much a blob of text looks like words I'm interested in**

**Model description**

BART is a transformer encoder-decoder (seq2seq) model with a bidirectional (BERT-like) encoder

BART is particularly effective when fine-tuned for text generation (e.g. summarization, translation) but also works well for comprehension tasks (e.g. text classification, question answering).

# Zero-shot-classification

```
sequence_to_classify = "one day I will see the world"
candidate_labels = ['travel', 'cooking', 'dancing']
classifier(sequence_to_classify, candidate_labels)
#{'labels': ['travel', 'dancing', 'cooking'],
# 'scores': [0.9938651323318481, 0.0032737774308770895, 0.0028610343357920289],
# 'sequence': 'one day I will see the world'}
```

This demo, from hugging face shows three keywords ('travel', 'cooking', 'dancing'). It receives a sentence ('one day I will see the world') and then it analyzes the text, resulting in a score for each keyword. The score represents how close the text is to each keyword.

For my test, I gave it 5 keywords. They aren't super specific, but just to use against any text I send to it. You can change this in the code (at end of presentation)

```python
# Load the zero-shot classification model
classifier = pipeline('zero-shot-classification', model='facebook/bart-large-mnli', device=device)
labels = ["login", "banking", "security", "phishing", "phone"]


# Define request body model
class TextData(BaseModel):
    text: str


# Route to classify text
@router.post("/classify")
async def classify_text(data: TextData):
    classified_data = classifier(data.text, labels)
    return {"labels": classified_data["labels"], "scores": classified_data["scores"]}
```
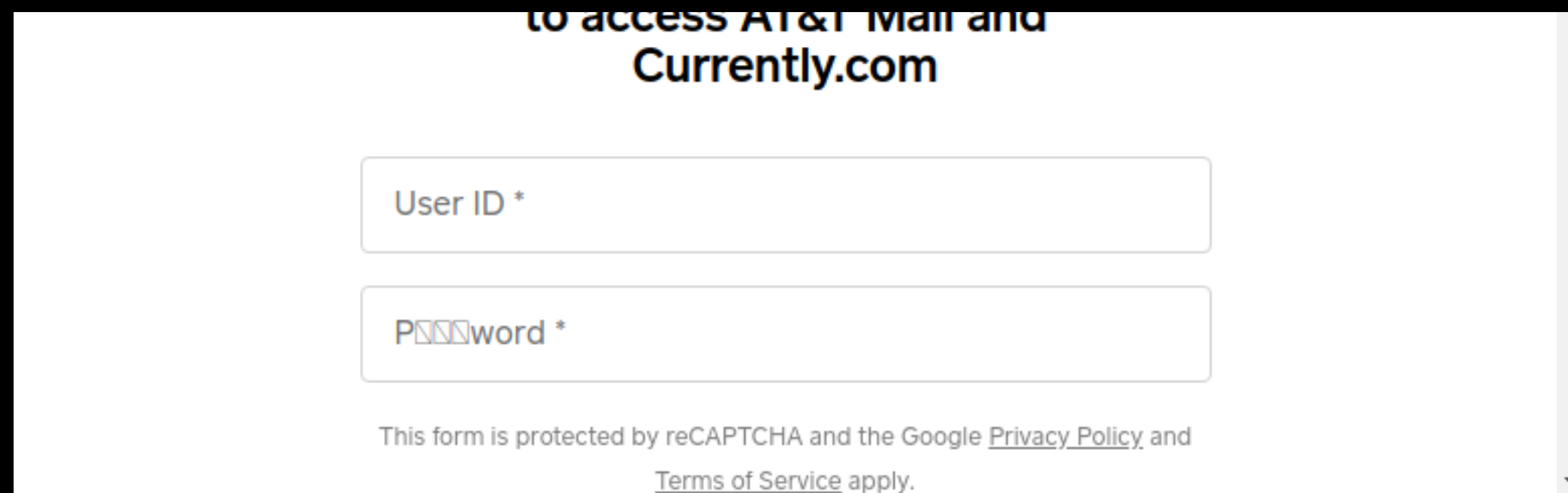
Zero-shot-classification results

```
'OCR text': 'VERIFICACION ss ',


'login': 0.6300058364868164
'phone': 0.20109270513055771
'banking': 0.10679445415735245
'phishing': 0.015628859400749207


'date_added': '2024-10-16T06:26:55.021364'
'phash': 'f5d320d88b7c8b2c'
'a_record': '146.112.61.108'
'city': 'San Jose'
'country': 'United States'
'url': 'papaya20688090.brizy.site'
```

009-288-49.weeblysite.com

Here, I'm seeing a comparison of how close the html source and the OCR text of this phishing page is to each of my specified keywords

to access AT&T Mail and
Currently.com

| User ID * |

| P⬚⬚⬚word * |

This form is protected by reCAPTCHA and the Google Privacy Policy and Terms of Service apply.

OCR Classification:

login: 0.8000938296318054
security: 0.12101650983095169
phone: 0.05189700424671173
banking: 0.016330692917108536
phishing: 0.010662009939551353

HTML Source Classification:

login: 0.8146952986717224
security: 0.09933193027973175
phone: 0.053744859993457794
banking: 0.017227377742528915
phishing: 0.015000510029494762

I want a better text scraping method. Thankfully, I found an LLM just for this.

# Parsing the UI
## And getting better OCR Results

# OmniParser: Screen Parsing tool for Pure Vision Based GUI Agent

Microsoft released an LLM that finds user interface elements of a screenshot. I set it up as a fast-api service that I can call from my web app.

Paper | License MIT

📣 [Project Page] [Blog Post] [Models]

**OmniParser** is a comprehensive method for parsing user interface screenshots into structured and easy-to-understand elements, which significantly enhances the ability of GPT-4V to generate actions that can be accurately grounded in the corresponding regions of the interface.

microsoft/**OmniParser** 🗐   ♡ like 897   Follow ▉ Microsoft 4,502

Image-Text-to-Text | 🤗 Transformers | 🌮 Safetensors | blip-2 | visual-question-answering | Inference Endpoi

📦 Model card   ≔ Files and versions   🤝 Community 12

📣 [Project Page] [Blog Post] [Demo]

## Model Summary

OmniParser is a general screen parsing tool, which interprets/converts UI screenshot to structured format, to improve existing LLM based UI agent. Training Datasets include: 1) an interactable icon detection dataset, which was curated from popular web pages and automatically annotated to highlight clickable and actionable regions, and 2) an icon description dataset, designed to associate each UI element with its corresponding function.

This model hub includes a finetuned version of YOLOv8 and a finetuned BLIP-2 model on the above dataset respectively. For more details of the models used and finetuning, please refer to the paper.

# This is OmniParsers demo of how this works. You upload an image (left) and get a result (right).



## OmniParser for Pure Vision Based General GUI Agent 🔥

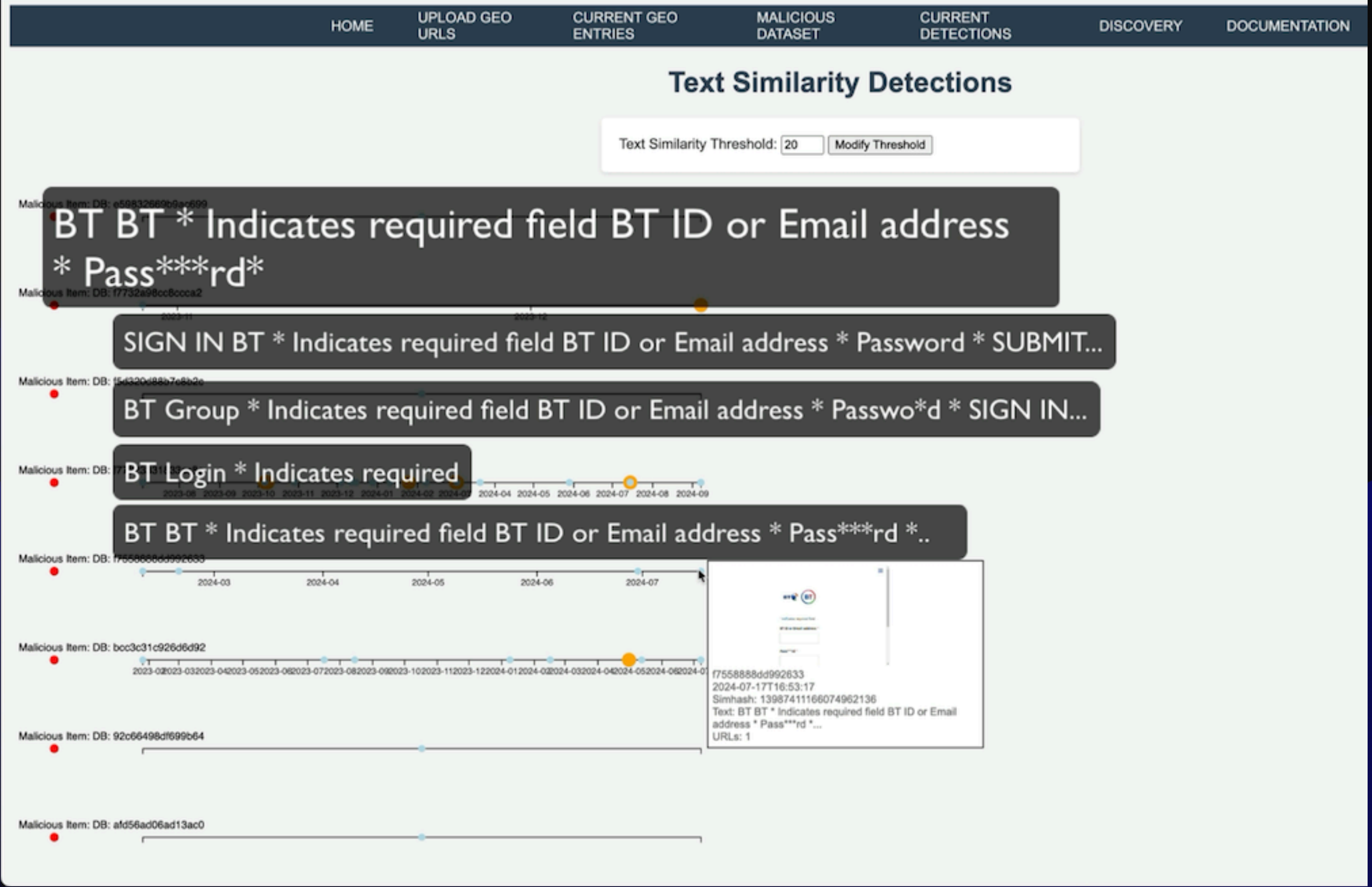OmniParser is a screen parsing tool to convert general GUI screen to structured elements.

**Upload image**

AT&T

**Sign in**

Complete the process by sign-in to AT&T Mail Any form of incorrect email/passcode will be marked as spam and email will be deleted

Email *

Password

**Box Threshold**          0.05

0.01 ————————————————— 1

**IOU Threshold**          0.1

0.01 ————————————————— 1

**Submit**

**Image Output**

0
AT&T
1
Sign in
2
Complete the process by sign-in to AT&T Mail Any form
3
of incorrect email/passcode will be marked as spam and email will be deleted
4
5
Email *
6
Password

**Parsed screen elements**

Text Box ID 0: AT&T
Text Box ID 1: Sign in
Text Box ID 2: Complete the process by sign-in to AT&T Mail Any form
Text Box ID 3: of incorrect email/passcode will be marked as spam and
Text Box ID 4: emall will be deleted
Text Box ID 5: Email
Text Box ID 6: Password

| URL | OCR text | OmniParse | Notable Differences |
|-----|----------|-----------|---------------------|
| attmailaccount22.weebly.com | login (0.551)<br>security (0.262)<br>phone (0.115) | login (0.667)<br>phone (0.226)<br>security (0.070) | Higher login and phone scores in OmniParse |
| atdikg.taplink.ws | login (0.706)<br>security (0.206)<br>phone (0.058) | login (0.574)<br>security (0.236)<br>phone (0.162) | Lower login, higher security and phone in OmniParse |
| att-mail-109459.weeblysite.com | login (0.641)<br>security (0.263)<br>phone (0.062) | login (0.638)<br>phone (0.171)<br>security (0.163) | Similar login, higher phone and security in OmniParse |

# Text similarity from the OCR text



Showing text similarity of scraped text on a timeline (timeline shows the screenshot, but is grouped from the OCR text)

# Research
## Campaigns over Time

I started building this as an inline detection tool, but it actually might have more use in researching phishing campaigns.



Looking at activity of specific phishing campaigns over a time period

# Alerting

# Taking Action

If using for detection, you're probably going to take some sort of action when a new phishing page is seen

With alerts going to your SIEM, it's easy to automate the blocking process

And that action is typically to block the domain or URL



📂 **Expanded document**

**Table**  JSON

| Actions | Field | Value |
|---------|-------|-------|
| | t _id | lO9MUZIB6OTviqU869RF |
| | t _index | phash_similarity_detections |

**Block** | auto1matiquelaposte.weebly.com |

| | t city | San Jose |
| | t country | United States |
| | 📅 date_added | Oct 3, 2024 @ 00:34:08.416 |
| | t phash | 97e718e118e316e3 |
| | t thumbnail_path | detections/auto1matiquelaposte.weebly.com.jpg |
| | t url | auto1matiquelaposte.weebly.com |

Oct 3, 2024 @ 00:13:42.819 - Oct 14, 2024 @ 00:00:00.000

| | |
|---|---|
| 3, 2024 @ 01:13:35.247 | support-team-comunidad.weebly.com |
| 3, 2024 @ 01:13:35.098 | fjgrhgrhgej.weebly.com |
| 3, 2024 @ 01:13:34.959 | atlascz-105067.weeblysite.com |
| 3, 2024 @ 01:13:34.862 | at767uyjhb.square.site |
| 3, 2024 @ 01:13:34.807 | signin-att-mail.weebly.com |
| 3, 2024 @ 01:13:34.584 | btinternet-108617.weeblysite.com |
| 3, 2024 @ 01:13:34.410 | reagan-105881.weeblysite.com |
| 3, 2024 @ 01:13:34.351 | sky-105224.weeblysite.com |
| 3, 2024 @ 01:13:34.176 | bt-103520.weeblysite.com |
| 3, 2024 @ 01:13:34.089 | bt-102082.weeblysite.com |
| 3, 2024 @ 01:13:34.032 | wpnhcj.weebly.com |

The next several slides will show various aspects of the interface at the time these slides were created (the code is still in development, so it might change in the future)

# Current Version of Interface

**Detections page**

# Current Detections

| Added | Image PHASH | URL | A Record | Location | Extracted Text | Classification | HTML Source Simhash |
|---|---|---|---|---|---|---|---|
| 2024-07-09 22:13:34 | b86fc59090c73b3a | rackspace18839994004885995.weebly.com | 146.112.61.108 | San Jose, United States | Webmail Login SUSPICIOUS Email address EMAIL Learn to quickly identify and report P@ssword A | login: 0.45 phishing: 0.37 security: 0.17 banking: 0.01 phone: 0.00 | 10337228896236707554 |
| 2024-01-11 00:42:47 | f464239a99999b39 | my-site-101155.weeblysite.com | 146.112.61.108 | San Jose, United States | A6B Q BJIE3TE, 3A AA IPObJD+KUTE ABB norpeбилen * Hapona This form is protected by reCAPTCHA and the Google Privacy Policy and Terms of Service apply. | security: 0.76 login: 0.10 phone: 0.08 phishing: 0.03 banking: 0.02 | 9931658834035383472 |
| 2024-07-02 21:04:28 | e59832669b9ac699 | inshaoiuyftdtcfyguijcdowiuhmkjnbhvgcfdxsedrftgyuhjlkopoiu.pages.dev | 146.112.61.108 | San Jose, United States | Enter you mailbox password to proceed. Password Remember me Sign In | security: 0.64 login: 0.30 phone: 0.02 banking: 0.02 phishing: 0.02 | 14423626393920664642 |
| 2024-02-08 16:02:53 | f777cc883164c998 | worker-home-att-767e.rljolo7229.workers.dev | 146.112.61.108 | San Jose, United States | AT&T Sign In to access AT&T Mail and Currently.com User ID Forgot user ID? Password | login: 0.51 security: 0.33 phone: 0.13 phishing: 0.02 banking: 0.01 | 3930323897056612998 |
| 2024-07-05 12:10:01 | b7f58c8c603333d8 | wefvbjhwevbhiwcdghwghiwerhrfhghbvhcdvd.weebly.com | 146.112.61.108 | San Jose, United States | SNetZero* Welcome! Please sign-in below. * Indicates required field Email * Passvvord * | login: 0.82 security: 0.12 phishing: 0.02 phone: 0.02 banking: 0.02 | 9944788745945115358 |

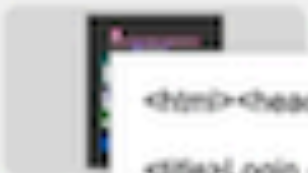# Malicious Dataset

Hovering over an OmniParser result screenshot

Upload

| Added | Image PHASH | URL | A Record | Location | Extracted Text | Classified Data | HTML Source Simhash |
|---|---|---|---|---|---|---|---|
| 2024-11-07 09:07:42 | e59832669b9ac699 | internationaldomaierepres... | | | | security: 0.64<br>login: 0.30<br>phone: 0.02<br>banking: 0.02<br>phishing: 0.02 | 1442362639392064642 |
| 2024-11-07 09:07:43 | f7732a98cc8occa2 | ionos-webmail-105551.wee... | | | | security: 0.58<br>login: 0.39<br>phishing: 0.02<br>phone: 0.01<br>banking: 0.01 | 3024539918263099024 |
| 2024-11-07 09:07:43 | f5d320d88b7c8b2c | papaya20688090.brizy.site | | | | security: 0.51<br>login: 0.39<br>phone: 0.05<br>phishing: 0.04<br>banking: 0.02 | 1949342313809748612 |
| 2024-11-07 09:07:44 | f77323331833cc8c | mail-105085.weeblysite.com | 74.115.51.55 | Unknown City, United States | REAGAN Email Address * Passw*rd * This form is protected by reCAPTCHA and the Google Privacy Policy and Terms of Service apply. Log In | security: 0.81<br>login: 0.16<br>phishing: 0.02<br>phone: 0.01<br>banking: 0.01 | 3160703433817627825 |
| 2024-11-07 09:07:45 | | mazulyta.pages.dev | 146.112.61.108 | San Jose, United States | | login: 0.84<br>security: 0.11<br>phishing: 0.03 | 6441788642322681002 |

0
Enter you mailbox password to
1 proceed.
5
2
Password
7
6 3
Remember me
4
Sign In

DEEPSEC

Hovering over the html

PHASH Detections

Timeline based off screenshot similarities

# Detecting Phishing using Visual Similarity

https://gitlab.pyosec.com/

**Thank you!**

Find me:

pyosec.com

bsky.app/profile/joshpyorre.bsky.social

infosec.exchange/@joshpyorre

Music (dark electronic): dievortex.com