

Net Gains: Systematic Phishing Page Differentiation by JavaScript Behavior

Anonymous authors

Abstract

In 2024, the Anti-Phishing Work Group identified over one million phishing pages. Phishers achieve this scale by using phishing kits — ready-to-deploy phishing websites — to rapidly deploy phishing campaigns with specific data exfiltration, evasion, or mimicry techniques. In contrast, researchers and defenders continue to fight phishing on a page-by-page basis and rely on manual analysis to recognize static features for kit identification.

This paper aims to aid researchers and analysts by automatically differentiating groups of phishing pages based on the underlying kit, automating a previously manual process, and enabling us to measure how popular different client-side techniques are across these groups. For kit detection, our system has an accuracy of 97% on a ground-truth dataset of 548 kits families deployed across 4,562 phishing URLs. On an unlabeled dataset, we leverage the complexity of 434,050 phishing pages’ JavaScript logic to group them into 11,377 clusters, annotating the clusters with what phishing techniques they employ. We find that UI interactivity and basic fingerprinting are universal techniques, present in 90% and 80% of the clusters, respectively. On the other hand, mouse detection via the browser’s mouse API is among the rarest behaviors, despite being used in a deployment of a 7-year-old open-source phishing kit. Our methods and findings provide new ways for researchers and analysts to tackle the volume of phishing pages.

1 Introduction

Web-based phishing attacks, where a webpage, through mimicking or urgency, tricks the user into submitting personal information or handing over access to a machine to an attacker, have been increasing for the last 5 years [9, 15]. Phishing attacks can have high-profile targets, like the NGOs and government workers targeted in 2021 [3, 26] and lead to more sophisticated cyberattacks and data breaches [4]. Phishing kits, ready-to-deploy software packages sold at illicit markets

for launching phishing attacks, have lowered the barrier of entry for malicious actors. Sellers often market phishing kits as bundles of quality-of-life features for attackers, such as built-in evasions from automated crawlers, exfiltration to Telegram channels, and obfuscation [46, 61]. Deploying these kits can be as easy as uploading them to free hosting providers and mass sending multiple links that exfiltrate the credentials to an endpoint controlled by the attacker. As phishing kits receive software updates from the original developer or the phishers who bought them, they split into variations belonging to the same “phishing kit family”.

One of the selling points of phishing kits is evasion from researchers and analysts, extending the time between deployment and discovery, and increasing the number of victims visiting the page without a browser warning. While server-side logic of kits can only make assumptions about the system based on the IP address and user agent, through API calls to the browser, client-side JavaScript code can query the user’s system for CPU core counts, memory overhead, request user interaction, or call out to a third party bot detection like CloudFlare¹ [65, 67]. In the end, the JavaScript logic in a kit can range from a simple user-agent-based redirection to an AES-encrypted script that dynamically decrypts itself, identifies the browser through a series of API calls, and renders the page after confirming the victim is using a mobile device.

This paper aims to aid researchers and analysts by automatically differentiating groups of phishing pages based on the pages’ JavaScript behaviors. This automates a previously manual process and enables us to measure the popularity of different client-side techniques across these groups. Figure 1 shows the relative difference in volume of phishing domains compared to the number of clusters monthly active. We focus on the following 10 techniques employed by phishing pages: fingerprint exfiltration, client-side IP check, timing-based bot detection, encoding-based obfuscation, dynamic script execution, basic fingerprint, dynamic script injection, Cloudflare turnstiles, and pop-ups. These techniques harvest

¹ Similar to the phishing page that stole credentials from Troy Hunt, the maintainer of HaveIBeenPwned [2]

credentials [54], evade crawlers [45, 66], or obfuscate code to avoid analysis [50, 55], all of which are sought-after capabilities of phishing kits [46].

By clustering pages based on their browser API usage, we construct groups of pages based on a shared set of techniques. Evaluating these groups over a ground truth dataset of URLs to kit-family mappings, we find that clustering over the set of browser APIs executed yields an FMI-based accuracy of 97% (69% with a rebalanced ground truth) and a validity score of 91%. We then turn the clustering approach to 952,155 pages, which we group into 11,377 clusters.

This work stands apart from prior research, because we automatically differentiate pages using features from dynamic instead of static analysis. As these features are tied to capabilities a prospective kit buyer is looking for, they are less likely to vary between different deployments of the same kit family. Compared to the prior attempt at kit identification [12], of $F_1 = 9.03\%$ and $F_1 = 31.11\%$ with DOM clustering and URL path-based signatures, respectively, we achieve a much higher FMI of 97%.

Overall, we offer the following contributions:

- C1. We find that *browser API usage alone* is sufficient to isolate and distinguish known phishing kit families. With a ground truth dataset of 548 kit families deployed on 4,562 URLs, we achieve accuracy metrics of a 97% Fowlkes-Mallows score (69% when rebalanced for a single-instance of a mass-deployed kit) and an 91% V-measure in the clusters of these pages relative to their kits used.
- C2. We experimentally show that browser APIs common on the web (DOM APIs, property reads, etc.) serve as a valuable identifier for identifying kits, as they signal the kit developer’s choices, and that the more sophisticated a page’s client-side logic is, the more indicative it is of the underlying kit.
- C3. On the unlabeled dataset, by propagating techniques from member pages to the cluster as a whole, we find that UI-interactivity and basic fingerprinting are near-universal in the ecosystem. At the same time, mouse detection via browser APIs, Cloudflare Turnstile embedding, and dynamic script creation are still relatively rare. While we observe client-Side IP checks occur in 19,869 pages they spanned across 504 clusters, and 23% of the pages come from a single cluster, which consists of Facebook business account phishing pages. We also find that compared to prior work, there is a decreased ratio of pages that employ pop-ups as a form of bot detection.
- C4. We release a dataset of 4,562 phishing pages labeled with their underlying kit and an unlabeled dataset of browser APIs executed on over 1.3 million pages.

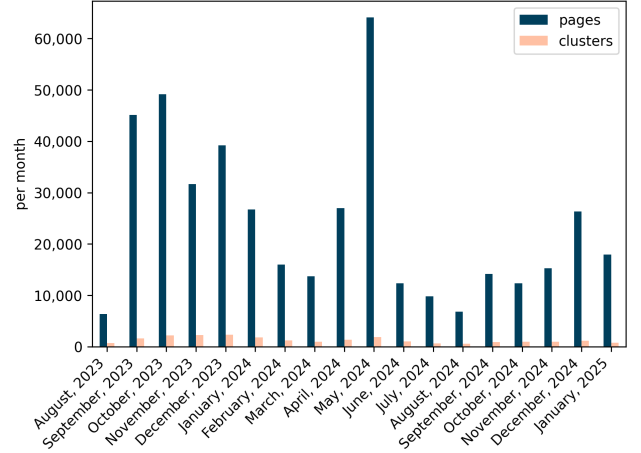


Figure 1: Comparison between monthly unique domains observed vs the monthly clusters observed based on dynamic behaviors. We see a drastic, non-linear reduction of phenomena that need to be investigated monthly.

2 Background

In this section, we provide a background on current developments in phishing as a phenomenon and adversarial JavaScript techniques, as well as an overview of the parties involved in the phishing ecosystem. We also introduce a basic overview of the clustering techniques used in this work.

2.1 The Phenomenon of Phishing

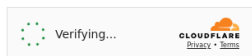
Phishing is a form of social engineering where an adversary pretends to be a trusted entity to steal a user’s credentials or gain access to a specific machine, network, or account to which the user has access. While delivery mechanisms vary, most phishing eventually leads to a webpage that requests some personal information (usually credentials) from the user. Because some legitimate websites use web fingerprints as a secondary authentication vector, phishers now also use browser fingerprinting APIs to identify real users and exfiltrate fingerprints to pair with stolen credentials [54].

Phishing is ever-evolving and still growing in prevalence. Groups that track phishing saw an increase in phishing domains in the 2020s. The most popular target sectors vary each year, but include software-as-a-service and webmail services (Q3 2020), financial services (Q3 2021), and social media (2024) [1]. In Q3 2024 alone, the Anti-Phishing Working Group (APWG) reported 900,000 phishing attacks.

As more enterprises and researchers study and combat phishing, phishers respond with new countermeasures to prevent automated crawling and phishing detection, collectively called “cloaking.” If a phishing page determines the client is not a viable victim (e.g., a crawling bot, not in a specific country, etc.), it takes actions not to serve real phishing content.

zss8ecker.com

Verifying you are human. This may take a few seconds.



zss8ecker.com needs to review the security of your connection before proceeding.

Ray ID:
Performance & security by Cloudflare

Figure 2: Example of a phishing page in our dataset that embedded Cloudflare Turnstile verification on a non-Cloudflare domain

The page may halt with an empty DOM or redirect to a benign page, a long-dead phishing page, or an affiliate marketing page.

Cloaking techniques are broadly classified into server-side and client-side techniques [46]. Server-side techniques are stealthier, but they rely on limited information about the client. Most server-side techniques rely on precompiled blocklists or allowlists of IP addresses, user agents, or referers (the page from which the link is visited, as identified by an HTTP header).

Client-side techniques allow for richer evasion strategies but are also more detectable. Phishing pages use browser APIs to trigger permission pop-ups to identify crawling browsers, which often cannot interact with the whole browser UI. Because cloaking is technically very similar to legitimate bot-detection and abuse prevention, phishers use CAPTCHA and click-through pages as client-side cloaking. Recently, phishing pages have used Cloudflare Turnstile, a popular widget for abuse prevention, to identify automated browsers. Figure 2 shows a phishing page using a Cloudflare Turnstile when we crawled it.

Even the URL features of a phishing link contain techniques that have evolved to respond to anti-phishing research. Phishing pages frequently use URL shorteners (public or private) to obfuscate the final destination, landing pages requiring a user to follow hyperlinks to the actual page, and free web hosting with trustworthy top-level domains (TLDs).

2.2 Phishing as an ecosystem

Phishing is a logistical and technical challenge because a phisher must develop an effective phishing page with cloaking, find robust hosting for it, entice a victim to browse to it

through SMS, email, or social media, exfiltrate the phished data, and then monetize the stolen credentials. This technical and logistical complexity, combined with interest from potential phishers, has created an underground economy to facilitate each step.

Phishing facilitators sell bundles of customizable or ready-to-deploy phishing pages known as “phishing kits.” They vary in features, sophistication, and cost. Phishing-as-a-service providers offer phishing kits and hosting services — essentially turnkey phishing systems. Both products lower barriers to entry. Prior work has shown that phishing kits may steal credentials from their customers’ kit deployments [17, 37], adapt or “borrow” features from other kits [27], and they are sometimes tied to specific actors [6]. In this paper, when two kits share the same set of features, and only differ in minor additions to new IPs in the blocklists, or different directory structure, we refer to them as being the same “kit-family.” Phishing systems may store credentials on the same server, risking loss when the page is inevitably taken down, but a more common practice is to send them to the phisher over instant messaging channels.

Credential sales markets simplify monetization. The credential sales part of the ecosystem has also adapted to modern MFA/2FA practices. With prior work showing that a browser fingerprint is enough to trick online services into triggering an MFA bypass [34], and has an increasing effect on the costs of stolen credentials [54].

2.3 JavaScript

Originally meant as a way of adding interactivity to webpages. The JavaScript ecosystem has evolved to allow varying low-level features to webpages through Browser APIs. HTML DOM APIs enable developers to modify page appearance, while LocalStorage and IndexedDB allow write access to the browser’s internal storage buckets; meanwhile, the File System API can allow access to the user’s real machine’s storage. Browser APIs can be function calls (or constructors), property reads, and property writes. Most of the privileged functionality comes from function calls.

The dynamic nature of JavaScript enables a variety of techniques for concealing itself from analysis and detection. JavaScript obfuscation can transform a known malicious sample into an undetectable one. Webpack enables bundling benign and malicious scripts and wrangling them to make static analysis harder. The other side of obfuscation is evasions; in addition to making code comprehension (via human or machine) harder, malicious actors have deployed time bombs, offloading parts of the malicious script to be read from the DOM or via a network request, and avoid detection. [50]

2.4 Hierarchical clustering

This paper utilizes hierarchical clustering, an unsupervised learning technique for segmenting data into nested structures (clusters) to identify pages that share phishing kits from their client-side behaviors. Specifically, we use HDBSCAN, a hierarchical variant of the density-based DBSCAN clustering technique. Advantages of HDBSCAN include not requiring prior knowledge of the number of clusters, no hyper-parameter tuning required out of the box². HDBSCAN has been used by prior work for categorizing malware families [50, 58]. When we can access ground truth, we utilize the V-measure and Fowlkes-Mallows index (FMI) to validate our clusters. Both are scores of 0 and 1 that address how well the clusters map to ground truth classes. V-measure (Validity measure) is the harmonic mean between completeness (all members of the same class are clustered together) and homogeneity (clusters only contain members of a single class) [53]. FMI, on the other hand, is the geometric mean between precision and recall, providing a close analog for an F1 score in supervised learning. When we do not have ground truth for the pages, we use the silhouette score of all the clusters. Silhouette-score measures how well separated the clusters are, between -1 and 1. While usually used to fine-tune hyperparameters, we primarily use it to measure how well-formed the structures we extract out of HDBSCAN for all of the phishing pages are. Scores under zero signal overlapping structure, while scores above 0.5 and 0.7 indicate medium or firm separation, respectively.

3 Methodology

This paper provides a methodology for clustering using dynamic features, evaluating how closely the clusters resemble underlying phishing kits, and describing how widespread different adversarial techniques are in the ecosystem. The building blocks of our experimental design are browser API execution traces from phishing pages and a ground-truth dataset of pages where we know the underlying phishing kit. In the following section, we describe the experimental setup for gathering this data, the steps we took to aggregate and enrich the execution traces, annotate the clusters based on different techniques, and finally, the steps we took for clustering the data and evaluating them as analogs for phishing kits. An overview of our crawling and analysis pipeline can be seen in Figure 3.

3.1 Data Gathering

Our crawling infrastructure aims to ingest phishing URLs from upstream providers and output execution traces from the page, as well as a potential kit used for that page.

URL feeds: We gathered phishing pages by monitoring the following phishing feeds: OpenPhish [49], PhishTank [16],

URLScan [56], SMS Gateways [43], PhishDB [41], and APWG [9], based on the availability of the feed and the level of access we had at the time. Every two hours, we checked these feeds for new URLs (limited to the last 48 hours) and submitted them to two different crawlers: VisibleV8 and *KitPhishr*.

VisibleV8: To get execution traces for every script loaded when visiting the page, we used an automated Chromium-based crawler with VisibleV8 patches applied [28]. The VisibleV8 patches modify Chromium to output a log of all JavaScript APIs executed for every script loaded on the page. We automate the browser to visit the page and take screenshots with puppeteer [24], an NPM package by Google to help in UI/UX testing and browser automation. The patched Chromium crawler uses puppeteer-stealth, a set of configurations to help mask the headless Chrome and Puppeteer itself from detection tools [48]. We initiated the crawls from a network designated for research purposes, for which the ISP would register as ‘educational’ for any IP intelligence API. We used Catapult [23], a man-in-the-middle proxy, to capture the entire HTTP archive for replayability. The crawler stays on the page for 45 seconds before taking a screenshot to allow scripts to load and start executing, consistent with prior work [20, 28].

KitPhishr: While not guaranteed, some malicious actors leave the zip files of the kits used in a discoverable folder on the same server that hosts the website (for example, the Apache document root). *KitPhishr* [18] is a Go-based URL fuzzer that attempts to identify any leftover zip files on the server. Prior work [35, 46] establishes *KitPhishr* as a method for collecting and analyzing phishing kits. If successful, it will download the zip file and make a note of the domain from which *KitPhishr* acquired it.

3.2 De-duplication

Once we have collected browser API traces and potential phishing kits, we post-process the traces into a set of APIs executed in a 1st-party context per page and de-duplicate the phishing kits.

3.2.1 Trace postprocessing

VisibleV8 has a default log-postprocessor set. These programs take the raw logs generated by the patched Chromium browser and convert them into an organized database, identifying duplicate scripts via SHA3 hash, clearly marking the origin of each script that loaded, and isolating which JavaScript API calls are browser API calls defined in the WebIDL file³ generated while building the patched Chromium. For our analysis, we extracted tuples of the original page’s URL, the script’s URL, and an unordered set of APIs executed by this script.

²HDBSCAN picks the most stable clusters based on the excess of mass algorithm

³<https://developer.mozilla.org/en-US/docs/Glossary/WebIDL>

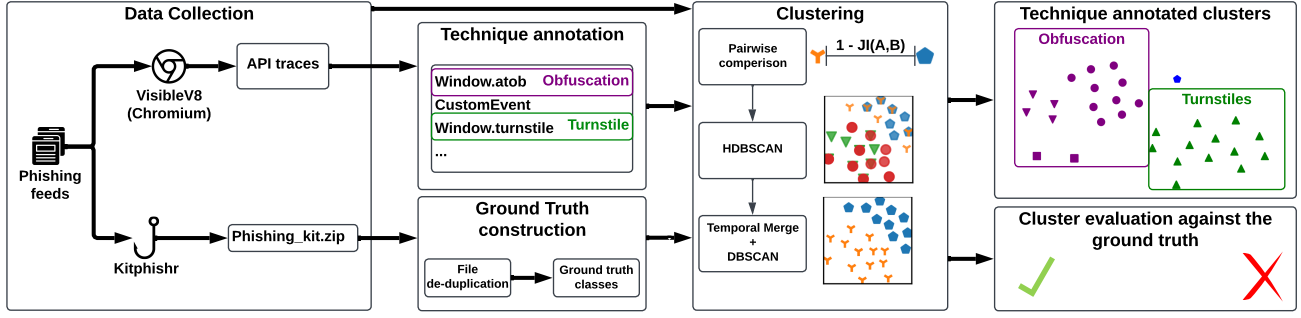


Figure 3: Crawling and clustering infrastructure

To not introduce artifacts into our clusters from cloaked pages and 3rd-party scripts like Google Analytics, known to be present in phishing pages [33], we isolate API sets executed by 1st-party scripts (from now on called 1st-party API sets). We establish the root domain as the domain submitted to the feeds and the origin as the domain from which the script is loaded. We consider a script 1st-party only if it is hosted on the same domain we acquired from our feeds (root domain). The only exception is when we identify which pages embed a Cloudflare Turnstile script. As some of the Turnstile scripts are hosted on ‘Cloudflare.com.’

Identifying Kit-families: We crawl the URLs with *KitPhishr* to establish a ground truth dataset with URLs originating from the same kit⁴. For zip files extracted via *KitPhishr* that have password-based encryption enabled, we use the SHA256 hash of the zip files to identify which domains yielded the same kit. For the remaining zip files, we further de-duplicate them into Kit-Families, examining them as a set of SHA256 hashes of source files⁵. If the Jaccard index-based similarity, where $J I(A, B) = |A \cap B| / |A \cup B|$, of these sets is equal to or greater than 90%, we consider the two kits to be from the same family, and thus group all domains from both kits into the same group.

Adjusting for Cloudflare: Many anomalies (multi-layer evals, non-deterministic behavior, dynamically updated scripts) originate from Cloudflare scripts on the same domain as the phishing pages. Since Cloudflare scripts load from a URL with a ‘cdn-cgi’ in the path, for most of our analysis, we do not consider scripts loaded from that endpoint.

3.3 Identifying kits

We hypothesize that similarity in browser API execution means that the pages originate from the same phishing kit. To test this, we ingest API sets from pages for which we were able to identify phishing kits and output a potential clustering of those pages, checking how well the clustering maps back to the ground truth information.

⁴For any domain that yielded two or more zipfiles, we discard all

⁵Source files identified via a python Magika module [22]

To establish this similarity, we use the Jaccard index on the 1st-party API sets with Hierarchical Density-Based Spatial Clustering of Applications with Noise [11] (HDBSCAN), with a minimum cluster size of 2, to cluster the pages with the Jaccard distance as our distance kernel. With HDBSCAN requiring minimal fine-tuning out of the box, and requiring no prior knowledge of the number of clusters we need to look for, we use the ground truth labels from *KitPhishr* to evaluate the clustering. When ground truth is available, we evaluate the clustering using the *Fowlkes-Mallows Index* [21] and *V-measure*. We use sklearn’s measure module to calculate all cluster evaluation metrics [51], and separate all the noise elements into singleton clusters for evaluation, as removing all noise samples would give it an unreasonable advantage; however, keeping the elements in the same noise cluster would artificially reduce the homogeneity score. For the unlabeled set of pages for which we do not have kit labels, we use silhouette score⁶, a metric for how well packed and separated the clusters are, to evaluate the clustering when we do not have ground truth.

To process 952,155 pages would require a distance matrix of distances (64-bit float) over a 1.5Tb in size. To resolve this restriction, we first divide our data using a 4-week rolling window, rolled by 2 weeks, and cluster pages within that window. We refer to these clusters as ‘local clusters’. However, these local clusters can not represent phenomena like the re-emergence of kits if it happens after 2 weeks and contains duplicates of the same page. To resolve this, we merge any two clusters with at least 1 page in common between local clusters and merge them using the representative API set for the clusters (API set intersection of every page in the cluster). We use DBSCAN with a conservative $\epsilon = 0.05$ to merge them.

Table 1: Manual mapping of phishing techniques to browser APIs

Technique	Category	Identifying markers
Fingerprinting extraction	Credential Harvesting	10 Fingerprinting API calls and an exfiltration related API call from [59]
Client-side IP check	Evasion	Window.fetch XMLHttpRequest.open
Timing bot detection	Evasion	Performance.now + Timeout
Encryption	Obfuscation	SubtleCrypto.decrypt
Encoding	Obfuscation	TextDecoder.decode window.atob
Dynamic script Evaluation	Obfuscation	eval
Basic fingerprinting	Evasion	HTMLDocument.cookie HTMLDocument.referrer Navigator.userAgent
Dynamic script creation	Evasion	HTMLScriptElement.text HTMLScriptElement.innerHTML
Cloudflare Turnstiles	Evasion	Window.turnstile
Pop-ups	Evasion	Navigator.requestMIDIAccess Clipboard.readText Geolocation.getCurrentPosition MediaDevices.getDisplayMedia HID.requestDevice Window.{confirm alert prompt} Accelerometer Gyroscope Window.showModalDialog MediaDevices.getUserMedia SyncManager.register Clipboard.read Serial.requestPort USB.requestDevice Window.queryLocalFonts Notification.requestPermission

3.4 Data enrichment

In addition to all the data gathered, our analysis references a manually crafted Browser API to phishing technique mapping, JavaScript execution traces for the top 5 brands’ login pages targeted in our dataset, and characterization of cluster lifetimes and deployment diversity.

Technique to Browser API mapping: Client-side JavaScript can engage in data harvesting (exfiltrating information dynamically and not through form submission), evasion (conditional dynamic behavior aimed at hiding functionality or contents), obfuscation (unconditional behavior meant to hinder static analysis), and mimicking (dynamic behavior to make the page more believable, for example, false loading pages, stage by stage data extraction). Based on prior work by Su *et al.* and Zhang *et al.* and manually identifying APIs from the Mozilla Developers Network (MDN) documentation, we present a table mapping standard phishing techniques to browser APIs

⁶While silhouette score is biased against non-convex clusters, based on our results, we do not see it necessary to switch to a density-specific cluster metric

in Table 1. We leverage the presence of these APIs in the execution traces as a signal of the technique being present in the page. If the page falls within a specific cluster, we mark the entire cluster as using that technique. We use this to combat the phishing pages that engage in non-deterministic behavior. For Cloudflare turnstiles embedding, we use a non-browser API as our detection metric, as Cloudflare’s native turnstile script will read the value of `Window.Turnstiles`. For Client-side IP checks, we manually looked through every `Window.fetch` and `XMLHttpRequest.open` argument URL given that argument was present in more than 50 phishing pages, and identified 15 that were IP reputation APIs.

Brand’s Original page: OpenPhish and APWG’s eCrime Exchange report the brand that a phishing page targets. We selected 5 out of the top 52 brands targeted based on popularity by page number and brands that represented seasonality targeted sectors (like the IRS or banks). We collected VisibleV8 logs for their home pages and, when applicable, their login pages, to assess how similar phishing pages are to their target pages.

Cluster lifetime and deployment diversity: Throughout this

work, we refer to cluster lifetime as the time range between when we observe the first page belonging to the cluster on a phishing feed and when we observe the last page belonging to the cluster. We measure deployment diversity of the phishing pages by looking at the effective 2nd-level domain (e2LD) for the URLs. Using the e2LD instead of the entire hostname ensures that pages deployed on 'pages.dev' or 'blogger.com' are considered a single deployment form.

4 Data characterization

This paper utilizes two distinct datasets. Labeled dataset of phishing pages deployed in the wild, their corresponding phishing kit, and an unlabeled dataset of phishing pages.

4.1 Labeled dataset

The ground truth data for this paper is a mapping of 548 kit families to 4,562 phishing pages. We collected 7,273 archive files by running all 952,155 through *KitPhishr*. Only two were encrypted, which we treated as unique kits based on their SHA256 hash. We further filtered the archive files by ensuring they contained at least a single code file, as Magika [22] identified. This reduces our archive count to 2,262. Finally, we group these kits into 2,00 families, out of which 548 were deployed on at least two distinct URLs. The most deployed kit was "2e94aff28a2c", a Wells-Fargo targeting kit (1,073 URLs), which made use of 3 separate server-side blocklists (.htaccess file and two separate PHP modules with regex rules for user-agent and IPs) and called 52 distinct browser APIs. Another kit of interests was seen throughout 545 days of the crawling (from 256 distinct URLs) was "fce61e98018d", a USPS phishing kit which executed on average 170 distinct APIs, with server-side and client-side IP checks, advance fingerprinting API calls, obfuscation, and JavaScriptgenerated DOM (Vite).

In line with prior work, 89% of these phishing kits were written in PHP, and we found 11 contained Python code. During our manual analysis of these kits, we concurred with prior work that many of these kits reused each other's code, especially regarding server-side IP blocklists. We found two module-like anti-bot detection files frequently redistributed across kits.

Finding 1: *Phishing pages from different kit families have vastly different browser API usage.* Pages from different kits have an average API similarity of 15.9% ($\sigma = 0.2$). On the other hand, on average, the pages from the same kit family have an API similarity of 98.6% ($\sigma = 0.1$).

4.2 Unlabeled Phishing pages

In total, we crawled for 523 days (2023-08-19 to 2025-01-17), collecting browser traces from 1,328,917 URLs. These included false Facebook suspension, USPS mail delivery, tech

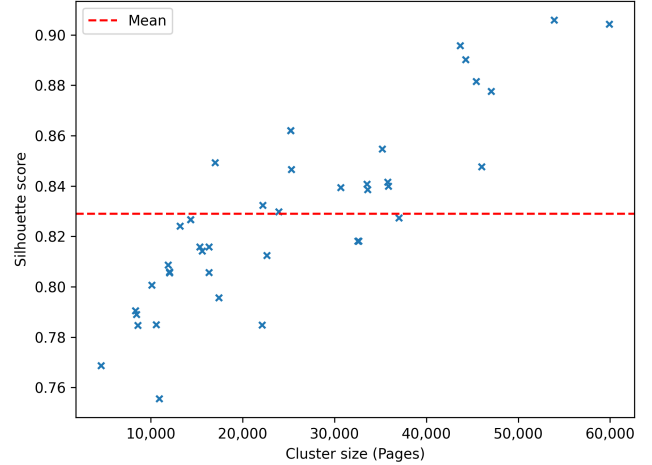


Figure 4: Distribution of the Silhouette Score of local clusters vs. their size (sum of pages)

support scam, AT&T login pages, and Bitcoin Wallet logins. All the pages collected have been labeled as phishing, as they pose as a trustworthy entity to gain credentials or network access from their victims, despite their varying tactics and credentials collected [15]. 388,536 pages did not execute any APIs in the first-party context, and only 434,050 executed at least 8 APIs in a first-party context to qualify for clustering. We first processed the pages into 27,833 local clusters. Figure 4 shows the distribution of the silhouette score of these local clusters. With an average score of 0.83, these clusters are well-formed. Once clustered together, the average cluster in our dataset contained 48 pages ($\sigma = 468$), existed for 40 days ($\sigma = 77$), and the intersection of API sets of all member pages had 63 ($\sigma = 57$) APIs in common. We remove 1,791 local-clusters from the merge algorithm, as they formed "malformed clusters" with fewer than 4-APIs in common between all pages. We note that even after the merge algorithm, 108,954 pages ended up belonging to more than one cluster; we do not merge these clusters further, instead all reporting in Section 5 is done as unique URLs.

After manually inspecting a sample of clusters, we observed that these clusters have unique pages across deployment types (AWS, Cloudflare, DigitalOcean, etc.) and languages. For example, Cluster-e325887b (shown in Figure 8) comprises 487 pages across five unique e2TLDs and contains pages engaging in voice-based phishing attacks in Japanese, English, and German, varying phone numbers, and error messages. With over 57 APIs in common, it is clear that these pages' usage of keyboard intercepting APIs, Audio APIs, and Network APIs calls (to ipwho.is) for IP intelligence caused the pages to cluster together.

Figure 5 shows two different clusters where pages have different DOM elements or landing pages being grouped, with example screenshots pulled from both. One cluster is from

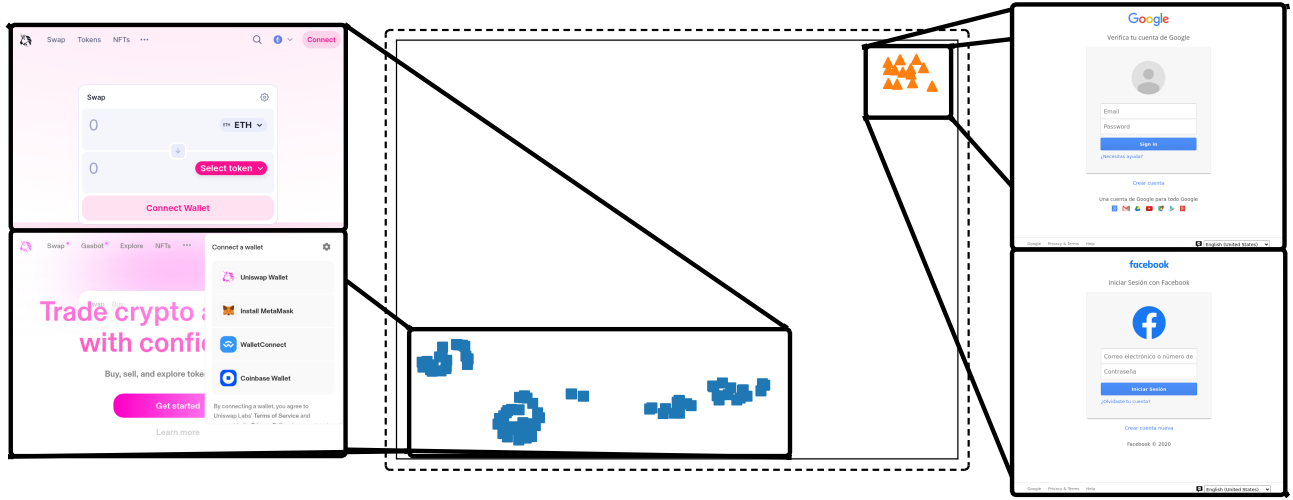


Figure 5: Example of pages from two different pages and embedding between them. We used t-distributed stochastic neighbor embedding to visualize the distance matrix between all pages in a 2-D plane. The clusters presented here are from Ethereum wallet pages with a few variations of the landing page (▲) and pages that descend from an open-source phishing kit discussed in Section 5 (■). We should note that while the pages have different logos on the right-hand side, the bottom left corner still reads "Google" on those pages.

a deployment we can tie to a public GitHub repository, and the other is a collection of crypto-wallet targeting pages. The figure shows a t-distributed stochastic neighbor (t-SDN) embedding for the distance matrix between all the pages from the two clusters.

5 Results

5.1 Kit identification

Finding 2: *In the majority of the cases, pages that execute at least two distinct browser APIs can be related to one another based on the underlying phishing kit; as the sophistication of the page grows, so does the accuracy of the clustering.* We find that clustering all pages that execute at least two browser API, yields an FMI based accuracy of 0.92. Figure 6 shows the V-Measure and FMI for our clusters as we increase the requirement of distinct APIs in the execution trace. We ultimately chose four browser APIs as the requirement for further experimentation on the ground truth data, as they provided a good tradeoff between V-measure and the number of pages used. For the unlabeled pages, we used 8 APIs, as they provided zero malformed clusters on our ground truth data (clusters with no API sets in common). Clustering pages from 4,562 pages across 548 kits, yields 654 clusters. Evaluating these clusters against the ground truth labels for each page, we find that our clusters have an FMI-based accuracy of 97% and a V-measure of 91% (with completeness and homogeneity scores of 90% and 92 respectively).

Table 2: Comparison of the evaluation metric when Script hashes are used instead of APIs executed

Method	FMI	V-measure
dynamic	0.97	0.91
Scripts (No eval)	0.88	0.85
Scripts (No eval, 1st party)	0.80	0.81
Scripts (1st party)	0.80	0.81
All script	0.89	0.85

Pages from kit "2e94aff28a2c" targeting Wells Fargo make up 24% of the ground truth dataset. We evaluate the clustering on a rebalanced dataset, where we sample 104 pages (the same number of pages as the second most popular kit) instead of using all 1,073. This reduces our FMI score to 0.69 (V-Measure of 0.9), maintaining higher accuracy than prior work. We discuss the potential ways to improve this accuracy in Section 8.

Finding 3: *Browser API sets better separate phishing kits than script hashes, even when we include scripts dynamically extracted from eval statements.* We use sha256 hashes of executed scripts as features in HDBSCAN to compare against our approach. Table 2 shows that browser API sets maintain a higher accuracy than SHA256 hashes of scripts, even for scripts that are extracted out of eval statements, and would require dynamic or static analysis to acquire.

Finding 4: *DOM APIs and property reads are a valuable signal in kit differentiation.* We find that removing DOM-

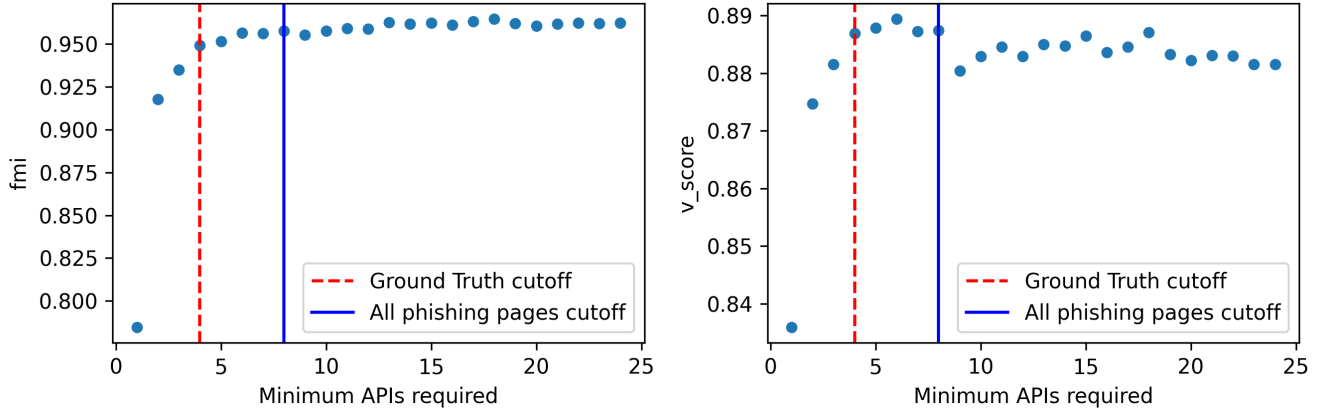


Figure 6: Validity measure for clusters vs. minmum distinct APIs required for clustering

related APIs or property reads out of consideration drastically reduces the number of pages we can consider for ground truth evaluation without increasing our overall accuracy. The key insight here is that these APIs signal particular choices the kit author made about the UI library they used, if they chose to hide the DOM as an evasion, or not draw it to begin, or which selectors, IDs, or classes they prefer to use. Clustering evaluation metrics for the ground truth dataset with DOM, SVG, and CSS APIs removed results in an FMI of 0.93 and a V-measure of 0.86. When all property reads were removed, we saw an FMI and V-measure of 0.93 and 0.85, respectively. In both cases, we can cluster fewer pages and thus identify fewer kits, but we do not see a significant improvement in clustering accuracy.

5.2 Clusters in the wild

Finding 5: 69% of clusters contain URLs only marked by a single target brand by our threat intel sources⁷. This phenomenon is observed in the ground truth dataset, as URLs for 90% of the kits were targeting a singular brand. Together, this provides strong circumstantial evidence that deployed kits are increasingly becoming brand-specific. 1,253 clusters (19%) of the clusters had two brand labels. However, the most popular combination of these was "Meta/Facebook", "Facebook/Instagram", "National Police Agency JAPAN/Facebook", and "holiganbet/jojobet", keeping the parent organization of the target the same in the majority of the cases. Manual examination of clusters with "National Police Agency JAPAN/Facebook" brand labels revealed shopping pages in Japanese to be mis-marked with that label from our data feeds. The cluster with the most diverse set of brand labels had 14 unique brand labels, which was a cluster with 12,467 pages with simple sign-in pages that exfiltrated information using client-side registered

⁷We did not include clusters in this count that had no brand-labeled URLs in them

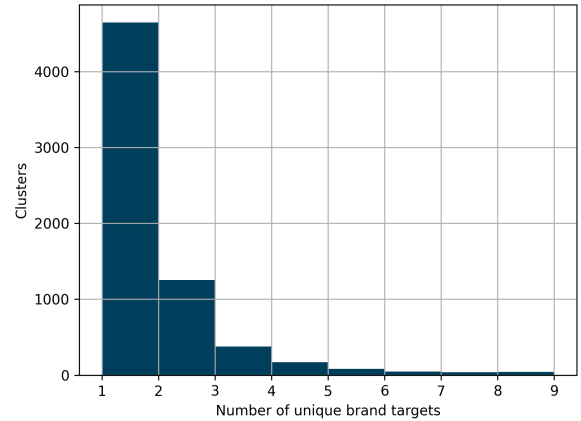


Figure 7: Distribution of unique brand labels per cluster (only counting clusters that had at least one page with a brand label)

event listeners. Furthermore, we find that 16,100 pages (3% of the pages observed spanning 313 clusters) come from phishing kits collected using KitPhisher; however, we could not pull the kit from the URL in all pages. Figure 7 shows a histogram of the unique brand labels observed per cluster.

5.3 Phishing Techniques across clusters

Finding 6: UI interactivity and fingerprinting are a near-universal behavior across clusters. Multi-stage phishing pages are very well documented in prior work, and we find that most clusters (91%) register a click event listener using JavaScript. Though this could be as simple as submitting credentials using JavaScript, this highlights the need for researchers to augment their crawlers in the future to extract better and more complete execution traces from websites. We split fingerprinting into two categories, basic and advanced. Basic fingerprinting,

Table 3: Breakdown of the obfuscation techniques observed in our dataset

Obfuscation techniques	Pages	Clusters
Window.atob	61,125	1,455
eval	14,561	982
Textdecoder.decode	11,113	534
SubtleCrypto.decrypt	1,185	36

which follows the list of APIs identified by Zhang *et al.* in [65] was present in 80% of the clusters (287,983 unique pages), and Advance fingerprinting (measured by at least 5 APIs identified by Su *et al.* in [59]) show up in 69% of the clusters. Together, 85% of clusters (9,572 clusters, 313,212 pages) exhibit some form of fingerprinting.

Finding 7: *Fingerprint exfiltration, dynamic script creation, obfuscation, and bot detection are uncommon across clusters.* While fingerprinting is near universal, we find that a smaller fraction of the clusters employ obfuscation, fingerprint exfiltration, and timing for bot detection. Prior work has shown interest in these behaviors [34, 50, 54, 66], meaning kits that forgo this may be rudimentary either by negligence or design, to avoid detection. Dropping anti-bot detection features has been observed before, with an Office-365 phishing kit (dubbed Tycoon2FA), opting to remove CloudFlare Turnstile integration, as it was being used as a feature for detection [39].

Another common tactic for bot detection is timing-based checks; by calling `Performance.now` right before and inside of a `Window.setTimeout` statement to measure the time differential between setting the timeout and its triggering. We find that 22% of clusters call `Performance.now` in conjunction with `setTimeout`.

On the other hand, 21% of clusters (2,395) employ some form of obfuscation, and 1,271 of the clusters (15,980 pages) dynamically generate an HTML script element. Table 3 lists all obfuscation techniques, with `eval` and Base64 encoding as the most popular methods. Despite the best recommendations to web developers to avoid using ‘eval’ [5], JavaScript’s `eval` function remains a favorite for obfuscation and evasions [50]. We observe script executed inside an ‘eval()’, evaluating yet another script; we measure this phenomenon as a level in *eval-depth*. We find that 48 clusters have pages that go to *eval-depth* 3. However, this seems to be a side-effect of embedding the phishing pages (mainly ones targeting Facebook) in Blogger.com pages.

Finding 8: *While rare, client-side IP reputation checks are present across multiple clusters.* While only present in 504 clusters (19,869 pages), we identify 15 unique IP reputation APIs used by phishing pages as soon as the page loads. We present a complete breakdown in Table 4 While not the most popular, *api.ipregistry.co* presents an interesting case study, as it enables the identification of educational networks. Manual examination of pages from these clusters reveals snippets sim-

```
await this.$http({
  method: "get",
  url: "https://api.ipregistry.co/?key=" +
    ~ this.key
}).then(e => {
  const s = e.connection.type,
    c = ["cdn", "hosting", "education"];
  /*omitted for brevity*/
  if (c.indexOf(s) !== -1)
    /*Redirect*/
})
```

Listing 1: Example of IP-based cloaking by using a 3rd-party reputation API. The following example specifically cloaks away from educational networks like ours.

Table 4: List of all API endpoints that client-side code reaches out for IP intelligence.

API url	clusters	pages
api.db-ip.com	40	3,124
api.geoapify.com	3	63
api.ipapi.com	5	143
api.ipgeolocation.io	21	96
api.ipify.org	177	7,417
api.ipregistry.co	9	3,995
freeipapi.com	38	6,444
geolocation-db.com	14	492
geolocation.onetrust.com	37	364
get.geojs.io	14	753
ipapi.co	106	735
ipinfo.io	65	1,963
ipwho.is	47	798
pro.ip-api.com	20	101

ilar to Listing 1, which conditionally chooses to redirect away from cloud hosting providers, content delivery networks, and educational networks, like the one we performed the crawls through. However, due to them employing other browser APIs, before the cloaking behavior, we can still cluster the pages based on the initial logic of the landing page.

Finding 9: *Pop-UP APIs are declining in usage.* We see only 104 clusters (1,323 pages) call out to pop-up requesting APIs. Geolocation.getCurrentPosition (55 clusters) was the most popular among these. While requiring a pop-up to interact with, this API can also be crucial in cloaking, as any VPN or proxy does not mask the results. We observe a smaller fraction of the ecosystem (16 clusters, 148 pages) than [65] employs this cloaking technique, especially when it comes to triggering a notification pop-up to verify user interaction. This could be a result of Firefox, citing low engagement with the notifications, starting to require user interaction to trigger the

Table 5: Comparison of browser APIs executed by phishing pages against the original login page of the brand they are targeting

Brand	Phishing pages observed	Avg % Similar	Std Dev (%)	Perfect Subset	50%+ Match
Facebook	338,686	11.6	9.6	0	24,039
USPS	66,692	10.3	10.5	0	53
Meta	51,969	12.4	7	0	2,605
Microsoft	3,768	11.8	8.7	0	419
IRS	1,207	11.5	9	0	61

popup [42] at the end of November 2019, when crawlphish’s data collection ended. Chrome has since discussed modifying the notification API to make the request less disruptive to the user experience [38]. The Lack of pop-up requests could also be explained by our *Finding 9* regarding the usage of IP intelligence APIs or by the overwhelming amount of pages (67%) registering at least one HTML element event handler, which would be classified as a *Click-through* by Crawlphish’s taxonomy. We report a full breakdown of APIs related to the Crawlphish categorization of client-side cloaking in Table 6.

Finding 10: *Mouse Detection API calls and Cloudflare Turnstile embedding are specific to a small group of clusters.* While supported by most modern browsers, we see an infrequent use of Mouse Detection APIs. Only 35 clusters employ mouse detection-related APIs. Two of these clusters⁸ are from an open-source phishing kit, leveraging botguard, and are from a public GitHub repository, which was last updated in 2017⁹. We see these clusters deploy across 17 unique domains, starting from 2023-10-07 all the way to 2024-07-19. 7 clusters (181 pages) embed a Cloudflare turnstile check in their page; some domains are not hosted on Cloudflare. It should be noted that in the case of redirection. At the same time, we discussed the presence of WebAssembly-based captchas for bot detection. Embedding a Cloudflare Turnstile check allows the phishing kit authors to offload bot detection to a well-established ecosystem. Recently, analyses have identified high-value phishing kits with this behavior [44]. However, subsequent analysis of the same threat actor identified a shift from turnstiles to HTML-Canvas drawn captchas.

Finding 11: *The phishing ecosystem consists of clusters that utilize both cutting-edge, experimental browser APIs, and extremely deprecated APIs.* We find 421 clusters (8,270 pages) that utilize *Navigator.UAData.getHighEntropyValues* for fingerprinting and *Keyboard.lock* to restrict user input, APIs not fully supported by Firefox and Safari. We identified 10 clusters across 4,433 pages that used *Scheduling.isInputPending*; however, upon closer inspection, these were not pages using a novel kit, but pages that used Google Sheets to construct their landing page. On the other hand, 25% of the clusters spanning (47,001) pages use a deprecated API.

While not experimental, WebAssembly is still a relatively modern web practice. We find 199 clusters (5,107 pages) that

use WebAssembly related APIs. Upon manual inspection of 33 unique WASM modules on these pages, we identify bot-detection, in most cases, by using FriendlyCaptcha, as the most common use case for WebAssembly in phishing.

Figure 10 shows the confusion matrix between the techniques we enumerated and the size of the clusters. We see no noticeable difference in techniques regarding cluster size, except for a higher percentage of large clusters using client-side IP checks. We also observe that pages that employ experimental APIs also tend to include a deprecated API call, which aligns with their usage for browser fingerprinting, rather than novel cloaking logic.

Finding 12: *Phishing pages vary wildly from the brand that they are mimicking.* We collect browser API traces from Facebook, USPS, Meta, Microsoft, and IRS login pages and compare them to their phishing counterparts. The average similarity of the APIs executed by the phishing page and the original page is 11%, indicating that browser APIs do not relate to the target page. We found no pages where the original page’s API set was a subset of the phishing page’s API set, and in 5% of the cases, the phishing page executed at least half of the APIs from the original page. We report per-brand findings in Table 5 and note that the least similar brand was USPS, which could be the result of USPS phishing pages being multi-stage pages requiring user interaction and targeting credit card information [13].

6 Discussion

The complexity of a phishing page’s client-side code helps analyze the massive volume of phishing pages in the wild caused by the proliferation of phishing kits. This section discusses what makes kits identifiable, why, and what the clusters represent.

6.1 What makes a kit identifiable?

The more sophisticated the phishing kit becomes, the easier it is to spot by the browser APIs it uses. Everything from exfiltrating a browser fingerprint to sophisticated evasions (e.g., Canvas captchas) makes mass deployments of the kit stand. There is also an economic incentive to sophistication, as novel use of browser APIs leads to better evasions from detectors [45, 65], more valuable credentials harvested [34, 54],

⁸Split due to infrastructure inability causing crashes in early crawls

⁹https://github.com/ashanahw/Gmail_Phishing

or resilience to analysis [50, 55]. For example, the APIs *Keyboard.lock*, *HTMLDocument.onkeydown* for keyboard locking, *Window.atob* for obfuscation, and a handful of fingerprint APIs and DOM APIs for dynamic content generation set the cluster shown in Figure 8 apart from other pages.

6.2 Advantages to behavioral aggregation

Phishing feeds are a noisy data source for studying the ecosystem, from e-commerce pages to mass-spammed USPS and EZ-parking phishing pages. As client-side code for phishing pages grows to be more complex, behavioral aggregation, akin to what has been done by prior work to identify exploit kits [58], is necessary. The methodology in this paper is aimed at researchers and analysts. For research, identifying shared kits in a dataset of phishing pages helps control for easily obtainable or mass-deployed phishing kits, measuring the prevalence of different techniques across kits, rather than pages. *At worst, we overestimate how popular different techniques are across kits, meaning we provide an upper bound for the less popular techniques.* For analysts, our methodology acts as a quick way to aggregate and share phishing kits-related threat intelligence between pages. Things like server-side cloaking technique, preferred exfiltration method, ties to APTs, and data exfiltrated. While kit-families can vary in which IPs they denylist, and what user-agents they allow, fingerprinting the underlying kit can allow analysts to deduce if a page employs these techniques in the first place.

7 Related Work

7.1 Phishing

Phishing detection: There is a wealth of research on phishing detection as the tug-of-war between adversaries and security professionals continues. Recently, Liu *et al.* and Abdelnabi *et al.* have deployed vision-based techniques to detect phishing pages [7, 35, 36]. [36] also presents over 6,000 phishing kits analyzed as part of the work. With adversarial attacks ensuring that the page looks different to crawlers and analysis, some have turned to extracting features from the URLs themselves, more recently via LLMs in [14, 29] and earlier via statistical models and machine learning in [31, 52, 57, 64]

Studying and combating adversarial techniques: Divakaran *et al.* in [19] reaffirms the need to keep up with the latest adversarial techniques to build better detection systems for phishing. Prominent work in this area includes [65] by Zhang *et al.*, which uncovered and categorized many novel client-side techniques by forcing the execution of phishing pages to trigger the cloaking behavior. Acharya *et al.* in [8] uncovered that phishing pages can successfully evade blocklists by knowing how to identify their crawlers, and Oest *et al.* in [45] demonstrated that cloaking from non-mobile based

devices as a phishing page can ensure that the attacker’s page goes unmarked by the blocklists for more than 48 hours.

Kondracki *et al.* in [30] uncovered a massive blindspot of the phishing detection ecosystem that was Man-in-the-middle phishing kits. Kits that would transparently forward the victim’s connection to the target page, mimicking brand logos on pages like Outlook without any configuration. Fortunately, the authors addressed the blind spot by demonstrating that these proxies remain fingerprintable using TLS fingerprinting. [63] proposed a similar attack, however, one that used JavaScript and NoVNC to trick the user into signing into their account through a VNC session in their browser.

With adversaries becoming creative with their evasions and obfuscation techniques, some novel defenses have also opted to think outside the box. Zhang *et al.* in [67] proposed a phishing defense solution that leverages the high likelihood of a phishing page cloaking away from a crawler to the defender’s advantage. They demonstrated that a web browser configured to look like a crawler triggers a cloaking response from phishing pages, ensuring that victims never see the page while maintaining compatibility with all of Alexa’s Top One Million websites. Meanwhile, using CAPTCHA [62] utilized vision-based models to combat phishing pages.

To better understand why, other than cloaking, phishing pages may choose to fingerprint. Lin *et al.* in [34] showed that browser fingerprints could be successfully used to bypass multi-factor authentication, a system meant to be a last line of defense against stolen credentials, for 10 out of 16 websites that provide popular services.

Phishing kits: Much can be studied about the phishing ecosystem via phishing kits. Cova *et al.* in [17] uncovered that most “free” phishing kits contain a backdoor, effectively serving as a way to offload the deployment of a campaign to a third party while siphoning off their stolen credentials.

Similar to our goals, PhishKitA [12] uses a dataset of phishing kits gathered through *KitPhishr* and a collection of features extracted from the HTML DOM to classify websites into their matching kit. Their multi-class classifiers for identifying the kit only achieved an F1 score of 39%, 31%, and 9% based on three different algorithms. Merlo *et al.* in [40] further expanded on our understanding of phishing kit lineages by looking at over 20,000 phishing kits and identifying, via token similarity, most of them as clones of one another or previously encountered kits. Prominent work in extending our understanding of phishing attacks includes Han *et al.* in [25], where they monitor the deployment of phishing kits by adversaries that compromise vulnerable web servers by hosting a well-sandboxed honeypot. They collected 643 phishing kits and established that kits take minutes to install and test, and can remain undetected for weeks. Using these kits, they were also able to identify evasion techniques used by these kits, like path randomization per visit, which, back then, was enough to bypass Google SafeBrowsing.

In [46], Oest *et al.* manually analyzed phishing kits to

establish the taxonomy for server-side cloaking, and in [10], Bijmans *et al.*, after collecting phishing kits by watching TLS transparency logs to identify Dutch bank phishing domains, manually created a fingerprint from static features to analyze their prevalence in the wild.

More recently, Lee *et al.* in [32] provides a server-side script (PHP) level analysis of phishing kits, finding that dynamically generated URLs are still standard in the ecosystem and observing seasonality in the kits they were able to obtain.

Characterizing the phishing ecosystem: Similar to our methods, Rola *et al.* in [54] deployed a modified Chromium browser to gather data and analyze phishing website browser APIs utilizing a pre-selected API list focused on first and third-party scripts for phishing pages. They find that most of the most visited phishing pages (identified via browser telemetry data) deploy fingerprinting scripts, sometimes varying from the ones of the original brand they portray. At the same time, they accessed this at a script level, reinforcing our finding that phishing pages vary vastly from their original page.

Oest *et al.* in [47] demonstrates the full lifecycle of a phishing campaign by employing the fact that phishing pages often copy assets from the target domain and refer the victim back to the original page afterward. By collaborating with a significant financial institution, they developed a framework for leveraging this data to track a phishing page from its deployment to blocklists, flagging the page as phishing. [47] observed all techniques highlighted by prior work: cloaking, user-specific URL generation, man-in-the-middle proxies, and short-lived bursty attacks. Expanding our understanding of the victim experience on a phishing website, Subramani *et al.* in [60] developed a crawler.

7.2 Dynamic analysis of webpages

Our work shares a methodology for dynamic analysis enabled by web measurement frameworks like OpenWPM [20] and VisibleV8 [28]. Su *et al.* used VisibleV8 traces and taint analysis in [59] to discover emerging fingerprinting techniques. Sarker *et al.* used VisibleV8 to create an oracle for detecting obfuscation [55]. Such an oracle was made possible by the observation that VisibleV8 marks the execution of an API at a given source line. At the same time, obfuscation techniques ensure that the API is not textually available there. Furthermore, Pantelaios *et al.* used a combination of VisibleV8 and force execution modifications to the Chromium engine to identify and defeat JavaScript evasion techniques while also leveraging API traces and clustering to identify previously unlabeled malicious extensions [50]. Iqbal *et al.* used OpenWPM to capture execution traces from franco top 100K URLs of Tranco, training a classifier on a mixture of dynamic and static features extracted from the JavaScript’s AST and execution traces, respectively, to achieve a 99.8% accuracy in identifying fingerprinting scripts online.

Our work differs from prior work in multiple ways. To date, we are the most successful in identifying phishing kits of a page; moreover, our differentiation is entirely automated, requiring no prior rule-based identification of kits. While we integrate many of the techniques annotated from prior work, we contribute an up-to-date understanding of their distribution in their ecosystem, and we do so at the cluster level, which is more likely to control for multiple deployments of the same sophisticated kit.

8 Threats to validity

Incorrect Ground truth mapping: While we manually validate our construction of kit families, we do not verify if the kit acquired is the kit deployed on the page. This remains an unexplored problem in the literature; however, this is expected to happen in sporadic cases. A mismatch between the deployed and the hosted kits would happen in cases of rare deployment errors by the phishers or shared infrastructure. As phishing infrastructure is ephemeral, averaging a below 2-day lifetime [45], cases of shared infrastructure can be expected to be rare; in most cases, this would be a result of two different entities compromising the same webhost to serve their phishing kits.

Unexplored Page states: While we visit and loiter on phishing pages, we do not explore and interact with the pages. Addressing this limitation would only improve our cluster evaluation metrics. Upon manual inspection, this is a source of the inaccuracy of clustering on the ground truth data, as some URLs are different paths on the same kit (including specific stages of the phishing page, or the root domain, which cloaks away). With recent work in LLM-powered crawlers, and ML-guided browser automation [60], we leave this to future work to use to collect a more complete list of browser APIs executed by a phishing page, at all stages of the phishing page. **Multi-page memebr clusters:** On the unlabeled dataset, our merge algorithm produces multi-member pages; these do not impact our findings for all techniques we categorize as rare, as they remain an upper estimate. For techniques we report as widespread (Fingerprinting and UI interactivity), we report distinct page numbers alongside the clusters to show that this assumption holds, as they appear in most of our pages clustered.

9 Conclusion

In this paper, we provided a workflow for researchers and analysts to automatically differentiate between a collection of phishing pages, based on a common underlying kit or shared techniques, if the behaviors are too generic. We show an accuracy of 97% against a dataset of pages and kits collected from the wild. With a curated mapping of techniques to browser APIs and 434,050 pages in which we identify 11,377 clusters,

we explore what techniques are universal, widespread across kits, or kit-specific.

10 Ethical consideration

This research relied on publicly and commercially available URLs detected using proprietary methods to be phishing websites. All of the crawling traffic originated from a network designated for web measurement, with the researchers monitoring the abuse contact for that IP space. We did not collect or store any identifiable information about the individuals behind phishing kits or the pages, and we did not conduct any live testing of their infrastructure, as we only visited at most twice upon ingestion. The phishing kits, on the other hand, contain identifiable information about their victims and Telegram API keys to access exfiltration group chats, and can be trivially manipulated to bypass current signature-based detection tools and redeployed in the wild. For this reason, we require a data-sharing agreement before we share the kits collected with any future researchers.

11 Open Science policy

We support the use of our data and tools by other researchers for any follow-up study of the phishing ecosystem or reproducibility work. We have made all of our analyses and data collection code available on anonymouse4open.science/r/NetGains for review. However, due to ethical concerns highlighted above, the data is available for researchers *upon request*. Besides containing URLs shared through data agreements with the feed operators, collected screenshots, HAR archives, and VisibleV8 logs amounting to 6.3TB of data, which would contain identifiable information about the reviews in packed binary archives (catapult’s HAR files), and create a heavy load on any data-sharing platform used.

References

- [1] APWG | Phishing Activity Trends Reports. [APWG Trends](#).
- [2] A Sneaky Phish Just Grabbed my Mailchimp Mailing List.
- [3] Sophisticated Spearphishing Campaign Targets Government Organizations, IGOs, and NGOs | CISA.
- [4] Threat Actor Leverages Compromised Account of Former Employee to Access State Government Organization | CISA.
- [5] Eval mdn. [MDN Docs](#), January 2025.
- [6] Unit 42. Threat Actor Groups Tracked by Palo Alto Networks Unit 42, June 2024.
- [7] Sahar Abdelnabi, Katharina Krombholz, and Mario Fritz. VisualPhishNet: Zero-Day Phishing Website Detection by Visual Similarity. In *Proceedings of the 2020 ACM SIGSAC Conference on Computer and Communications Security, CCS ’20*, pages 1681–1698, New York, NY, USA, November 2020. Association for Computing Machinery.
- [8] Bhupendra Acharya and Phani Vadrevu. {PhishPrint}: Evading Phishing Detection Crawlers by Prior Profiling. pages 3775–3792.
- [9] Anti-Phishing Working Group (APWG). eCrime Exchange (eCX), 2025.
- [10] Hugo Bijmans, Tim Booij, Anneke Schwedersky, Aria Nedgabat, and Rolf van Wegberg. Catching Phishers By Their Bait: Investigating the Dutch Phishing Landscape through Phishing Kit Detection. In *30th USENIX Security Symposium (USENIX Security 21)*, pages 3757–3774, 2021.
- [11] Ricardo J. G. B. Campello, Davoud Moulavi, and Joerg Sander. Density-based clustering based on hierarchical density estimates. In Jian Pei, Vincent S. Tseng, Longbing Cao, Hiroshi Motoda, and Guandong Xu, editors, *Advances in Knowledge Discovery and Data Mining*, pages 160–172, Berlin, Heidelberg, 2013. Springer Berlin Heidelberg.
- [12] Felipe Castaño, Eduardo Fidalgo Fernández, Rocío Alaiz-Rodríguez, and Enrique Alegre. PhiKitA: Phishing Kit Attacks Dataset for Phishing Websites Identification. *IEEE Access*, 11:40779–40789, 2023.
- [13] SANS Internet Storm Center. USPS Phishing Scam Targeting iOS Users.
- [14] Daiki Chiba, Hiroki Nakano, and Takashi Koide. DomainLynx: Leveraging Large Language Models for Enhanced Domain Squatting Detection. *ArXiv*, abs/2410.02095, 2024.
- [15] CISA. Recognize and report. [Secure Our World](#).
- [16] Cisco Talos Intelligence Group (Talos). Phishtank, 2025.
- [17] Marco Cova, Christopher Kruegel, and Giovanni Vigna. There is No Free Phish: An Analysis of “Free” and Live Phishing Kits.
- [18] cybercdh. Kitphishr: A tool designed to hunt for phishing kit source code. [Github](#), 2023.

- [19] Dinil Mon Divakaran and Adam Oest. Phishing Detection Leveraging Machine Learning and Deep Learning: A Review. *IEEE Security & Privacy*, 20(5):86–95, September 2022.
- [20] Steven Englehardt and Arvind Narayanan. Online Tracking: A 1-million-site Measurement and Analysis. In *Proceedings of ACM CCS 2016*, 2016.
- [21] Edward B Fowlkes and Colin L Mallows. A Method for Comparing Two Hierarchical Clusterings. *Journal of the American statistical association*, 78(383):553–569, 1983.
- [22] Yanick Fratantonio, Luca Invernizzi, Loua Farah, Kurt Thomas, Marina Zhang, Ange Albertini, Francois Galilee, Giancarlo Metitieri, Julien Cretin, Alexandre Petit-Bianco, David Tao, and Elie Bursztein. Magika: AI-Powered Content-Type Detection. In *Proceedings of the International Conference on Software Engineering (ICSE)*, April 2025.
- [23] Google Inc. Catapult. [Googlesource](#), 2025.
- [24] Google Inc. Puppeteer, 2025.
- [25] Xiao Han, Nizar Kheir, and Davide Balzarotti. PhishEye: Live Monitoring of Sandboxed Phishing Kits. In *Proceedings of the 2016 ACM SIGSAC Conference on Computer and Communications Security, CCS '16*, pages 1402–1413, New York, NY, USA, October 2016. Association for Computing Machinery.
- [26] Microsoft Threat Intelligence. New sophisticated email-based attack from NOBELIUM.
- [27] Microsoft Threat Intelligence. Franken-phish: TodayZoo built from other phishing kits. <https://www.microsoft.com/en-us/security/blog/2021/10/21/franken-phish-todayzoo-built-from-other-phishing-kits/>, October 2021.
- [28] Jordan Jueckstock and Alexandros Kapravelos. VisibleV8: In-browser Monitoring of JavaScript in the Wild. In *Proceedings of the ACM Internet Measurement Conference (IMC)*, 2019.
- [29] Koide, Takashi and Fukushi, Naoki and Nakano, Hiroki and Chiba, Daiki. Phishreplicant: A language model-based approach to detect generated squatting domain names. In *Proceedings of the 39th Annual Computer Security Applications Conference, ACSAC '23*, page 1–13, New York, NY, USA, 2023. Association for Computing Machinery.
- [30] Brian Kondracki, Babak Amin Azad, Oleksii Starov, and Nick Nikiforakis. Catching Transparent Phish: Analyzing and Detecting MITM Phishing Toolkits. In *Proceedings of the 2021 ACM SIGSAC Conference on Computer and Communications Security, CCS '21*, pages 36–50, New York, NY, USA, November 2021. Association for Computing Machinery.
- [31] Hung Le, Quang Pham, Doyen Sahoo, and Steven C. H. Hoi. URLNet: Learning a URL Representation with Deep Learning for Malicious URL Detection. *ArXiv*, abs/1802.03162, 2018.
- [32] Woonghee Lee, Junbeom Hur, and Doowon Kim. Beneath the phishing scripts: A script-level analysis of phishing kits and their impact on real-world phishing websites. In *Proceedings of the 19th ACM Asia Conference on Computer and Communications Security*, pages 856–872. ACM.
- [33] Kyungchan Lim, Jaehwan Park, and Doowon Kim. Phishing Vs. Legit: Comparative Analysis of Client-Side Resources of Phishing and Target Brand Websites. In *Proceedings of the ACM Web Conference 2024, WWW '24*, New York, NY, USA, 2024. Association for Computing Machinery.
- [34] Xu Lin, Panagiotis Ilia, Saumya Solanki, and Jason Polakis. Phish in Sheep’s Clothing: Exploring the Authentication Pitfalls of Browser Fingerprinting. pages 1651–1668.
- [35] Ruofan Liu, Yun Lin, Xianglin Yang, Siang Hwee Ng, Dinil Mon Divakaran, and Jin Song Dong. Inferring phishing intention via webpage appearance and dynamics: A deep vision based approach. pages 1633–1650.
- [36] Ruofan Liu, Yun Lin, Yifan Zhang, Penn Han Lee, and Jin Song Dong. Knowledge expansion and counterfactual interaction for {Reference-Based} phishing detection. pages 4139–4156.
- [37] Heather McCalley, Brad Wardman, and Gary Warner. Analysis of Back-Doored Phishing Kits. In Gilbert Peterson and Sujeet Sheno, editors, *Advances in Digital Forensics VII*, pages 155–168, Berlin, Heidelberg, 2011. Springer.
- [38] PJ McLachlan. Introducing quieter permission UI for notifications. [Chromium Blog](#).
- [39] Rodel Mendrez. Tycoon2FA New Evasion Technique for 2025. [spiderlabs blog](#).
- [40] Ettore Merlo, Mathieu Margier, Guy-Vincent Jourdan, and Iosif-Viorel Onut. Phishing kits source code similarity distribution: A case study. In *2022 IEEE International Conference on Software Analysis, Evolution and Reengineering (SANER)*, pages 983–994. ISSN: 1534-5351.
- [41] Mitchell Krog and Nissar Chababy. PhishingDB, 2025.

- [42] Mozilla. Restricting notification permission prompts in firefox. <https://blog.mozilla.org/futurereleases/2019/11/04/restricting-notification-permission-prompts-in-firefox> Mozilla Blog, November 2019.
- [43] Aleksandr Nahapetyan, Sathvik Prasad, Kevin Childs, Adam Oest, Yeganeh Ladwig, Alexandros Kapravelos, and Bradley Reaves. On SMS Phishing Tactics and Infrastructure. In *2024 IEEE Symposium on Security and Privacy (SP)*, pages 169–169. IEEE Computer Society, 2024.
- [44] NJCCIC. Recent tycoon 2fa phishing campaigns target government entities. [NJCCIC News](#).
- [45] Adam Oest, Yeganeh Safaei, Adam Doupe, Gail-Joon Ahn, Brad Wardman, and Kevin Tyers. PhishFarm: A Scalable Framework for Measuring the Effectiveness of Evasion Techniques against Browser Phishing Blacklists. In *2019 IEEE Symposium on Security and Privacy (SP)*, pages 1344–1361, San Francisco, CA, USA, May 2019. IEEE.
- [46] Adam Oest, Yeganeh Safei, Adam Doupe, Gail-Joon Ahn, Brad Wardman, and Gary Warner. Inside a Phisher’s Mind: Understanding the Anti-Phishing Ecosystem through Phishing Kit Analysis. In *2018 APWG Symposium on Electronic Crime Research (eCrime)*, pages 1–12, San Diego, CA, May 2018. IEEE.
- [47] Adam Oest, Penghui Zhang, Brad Wardman, Eric Nunes, Jakub Burgis, Ali Zand, Kurt Thomas, Adam Doupé, and Gail-Joon Ahn. Sunrise to Sunset: Analyzing the End-to-end Life Cycle and Effectiveness of Phishing Attacks at Scale. pages 361–377.
- [48] Berstend on Github. Puppeteer stealth plugin, 2025.
- [49] OpenPhish. OpenPhish, 2025.
- [50] Nikolaos Pantelaio and Alexandros Kapravelos. FV8: A Forced Execution JavaScript Engine for Detecting Evasive Techniques. In *Proceedings of the USENIX Security Symposium*, August 2024.
- [51] F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher, M. Perrot, and E. Duchesnay. Scikit-learn: Machine learning in Python. *Journal of Machine Learning Research*, 12:2825–2830, 2011.
- [52] Ch. Chakradhara Rao, A.V.T. Raghav Ramana, and B. Sowmya. Detection of phishing websites using hybrid model. 2018.
- [53] Andrew Rosenberg and Julia Hirschberg. V-measure: A conditional entropy-based external cluster evaluation measure. In *Conference on Empirical Methods in Natural Language Processing*, 2007.
- [54] Iskander Sanchez-Rola, Leyla Bilge, Davide Balzarotti, Armin Buescher, and Petros Efstathopoulos. Rods with Laser Beams: Understanding Browser Fingerprinting on Phishing Pages. pages 4157–4173.
- [55] Shaown Sarker, Jordan Jueckstock, and Alexandros Kapravelos. Hiding in Plain Site: Detecting JavaScript Obfuscation through Concealed Browser API Usage. In *Proceedings of the ACM Internet Measurement Conference (IMC)*, page 648–661, October 2020.
- [56] SecurityTrails. Urlscan, 2025.
- [57] Hossein Shirazi, Bruhadeshwar Bezawada, and Indrakshi Ray. "kn0w thy doma1n name": Unbiased phishing detection using domain name based features. *Proceedings of the 23rd ACM on Symposium on Access Control Models and Technologies*, 2018.
- [58] Ben Stock, Benjamin Livshits, and Benjamin Zorn. Kizle: A signature compiler for detecting exploit kits. In *2016 46th Annual IEEE/IFIP International Conference on Dependable Systems and Networks (DSN)*, pages 455–466, 2016.
- [59] Junhua Su and Alexandros Kapravelos. Automatic Discovery of Emerging Browser Fingerprinting Techniques. In *Proceedings of the ACM Web Conference 2023, WWW ’23*, page 2178–2188, New York, NY, USA, 2023. Association for Computing Machinery.
- [60] Karthika Subramani, William Melicher, Oleksii Starov, Phani Vadrevu, and Roberto Perdisci. PhishInPatterns: measuring elicited user interactions at scale on phishing websites. In *Proceedings of the 22nd ACM Internet Measurement Conference, IMC ’22*, pages 589–604. Association for Computing Machinery.
- [61] Bhaskar Tejaswi, Nayanamana Samarasinghe, Sajjad Pourali, Mohammad Mannan, and Amr Youssef. Leaky Kits: The Increased Risk of Data Exposure from Phishing Kits. In *2022 APWG Symposium on Electronic Crime Research (eCrime)*, pages 1–13, Boston, MA, USA, November 2022. IEEE.
- [62] Xiwen Teoh, Yun Lin, Ruofan Liu, Zhiyong Huang, and Jin Song Dong. {PhishDecloaker}: Detecting {CAPTCHA-cloaked} phishing websites via hybrid vision-based interactive models. pages 505–522.

- [63] Jonas Tzschoppe and Hans Löhr. Browser-in-the-middle - evaluation of a modern approach to phishing. In *Proceedings of the 16th European Workshop on System Security*, EUROSEC '23, pages 15–20. Association for Computing Machinery.
- [64] Rakesh M. Verma and Avisha Das. What's in a url: Fast feature extraction and malicious url detection. *Proceedings of the 3rd ACM on International Workshop on Security And Privacy Analytics*, 2017.
- [65] Penghui Zhang, Adam Oest, Haehyun Cho, Zhibo Sun, Rc Johnson, Brad Wardman, Shaown Sarker, Alexandros Kapravelos, Tiffany Bao, Ruoyu Wang, Yan Shoshitaishvili, Adam Doupé, and Gail-Joon Ahn. CrawlPhish: Large-scale Analysis of Client-side Cloaking Techniques in Phishing. In *2021 IEEE Symposium on Security and Privacy (SP)*, pages 1109–1124, San Francisco, CA, USA, May 2021. IEEE.
- [66] Penghui Zhang, Adam Oest, Haehyun Cho, Zhibo Sun, RC Johnson, Brad Wardman, Shaown Sarker, Alexandros Kapravelos, Tiffany Bao, Ruoyu Wang, Yan Shoshitaishvili, Adam Doupé, and Gail-Joon Ahn. CrawlPhish: Large-scale analysis of client-side cloaking techniques in phishing. In *2021 IEEE Symposium on Security and Privacy (SP)*, pages 1109–1124. ISSN: 2375-1207.
- [67] Penghui Zhang, Zhibo Sun, Sukwha Kyung, Hans Walter Behrens, Zion Leonahenahe Basque, Haehyun Cho, Adam Oest, Ruoyu Wang, Tiffany Bao, Yan Shoshitaishvili, Gail-Joon Ahn, and Adam Doupé. I'm SPARTACUS, No, I'm SPARTACUS: Proactively Protecting Users from Phishing by Intentionally Triggering Cloaking Behavior. In *Proceedings of the 2022 ACM SIGSAC Conference on Computer and Communications Security*, CCS '22, pages 3165–3179, New York, NY, USA, November 2022. Association for Computing Machinery.

Appendix A

Table 6: Summary of all API calls that match to a category identified by Crawlphish

Category	Clusters	Pages
User-Interaction		
Pop-up (Total)	104	1,323
Accelerometer	7	63
Clipboard.readText	1	2
Geolocation.getCurrentPosition	55	479
Gyroscope	4	10
MediaDevices.getUserMedia	5	55
Notification.requestPermission	16	148
Window.alert	22	581
Mouse	35	62
DomEvents	10,402	362,700
Fingerprint		
Total	9,077	287,983
Navigator.userAgent	8,641	266,126
HTMLDocument.cookie	4,664	124,110
Bot Detection		
Timing	2,528	78,202

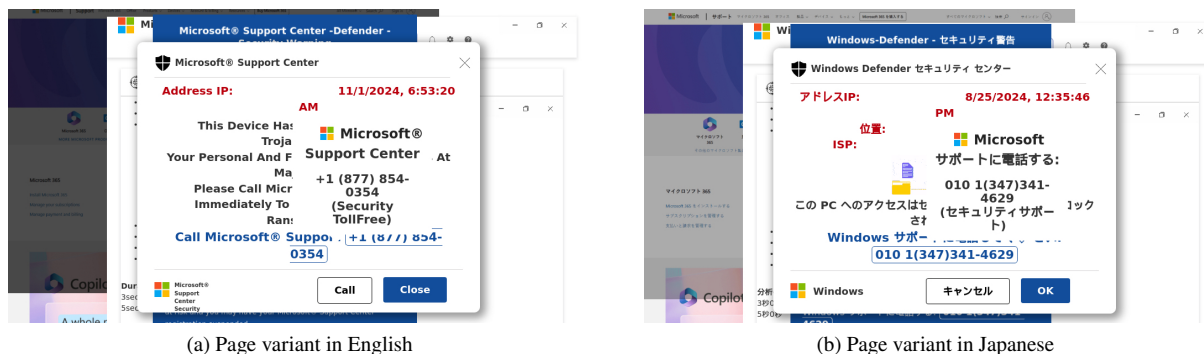


Figure 8: Cropped screenshots from Cluster-53d5c420, IP addresses and location redacted to ensure anonymity of the authors.

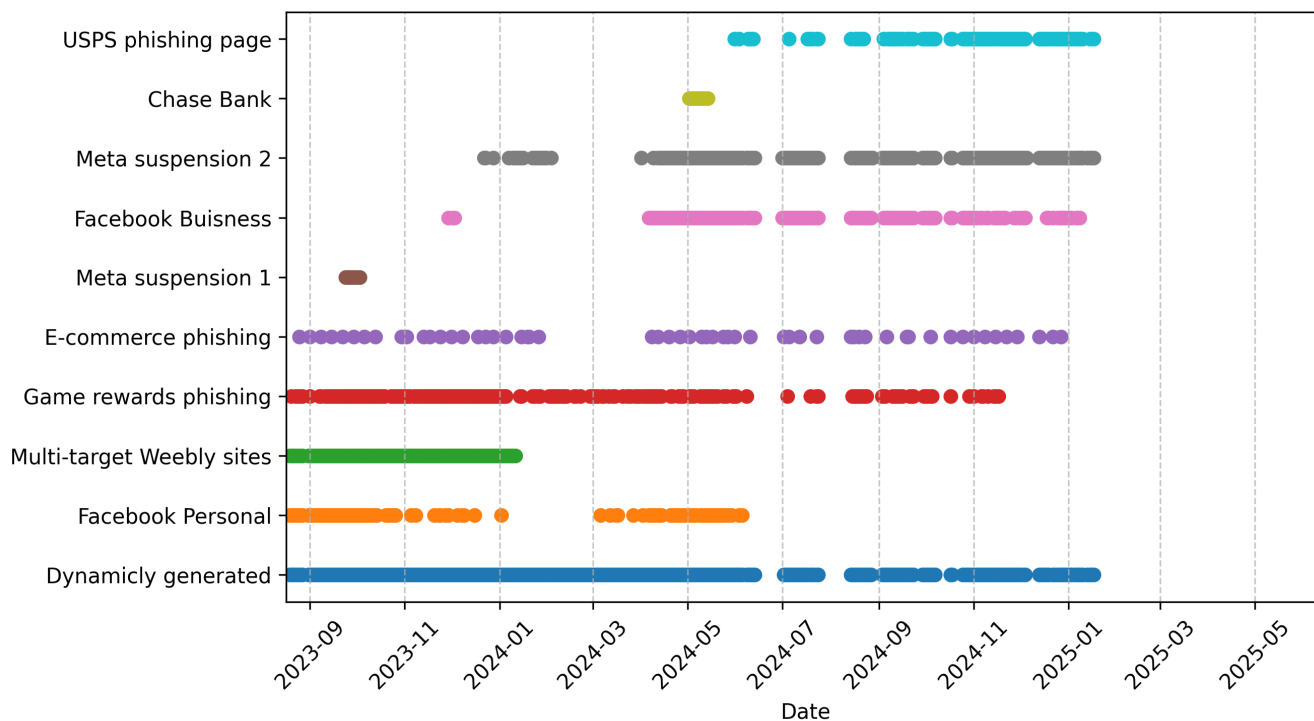


Figure 9: Timeline of the top 10 clusters based on the number of pages

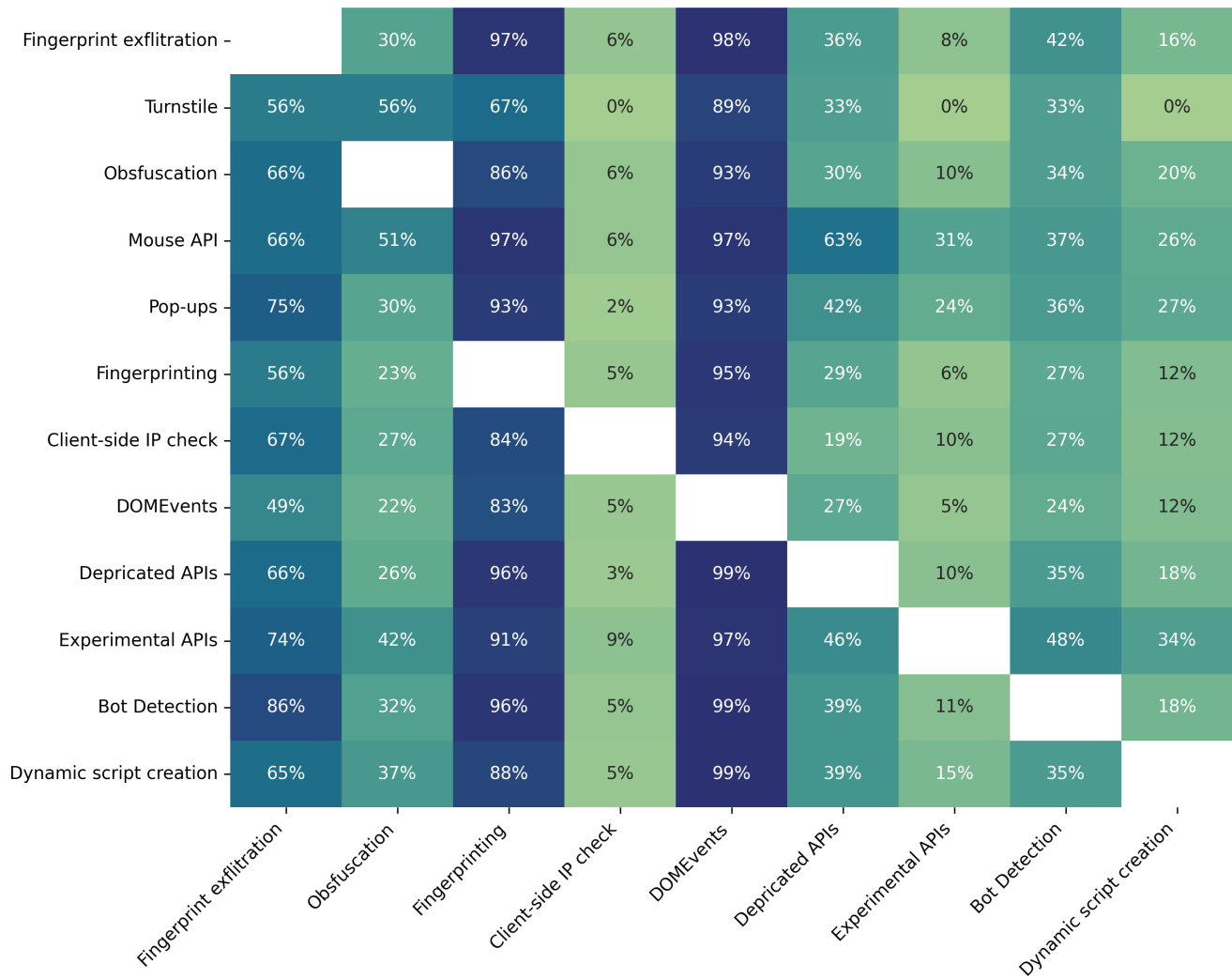


Figure 10: Confusion matrix between all of the techniques enumerated and cluster lifetime characteristics, normalized by row.