# PhishReplicant: A Language Model-based Approach to Detect Generated Squatting Domain Names
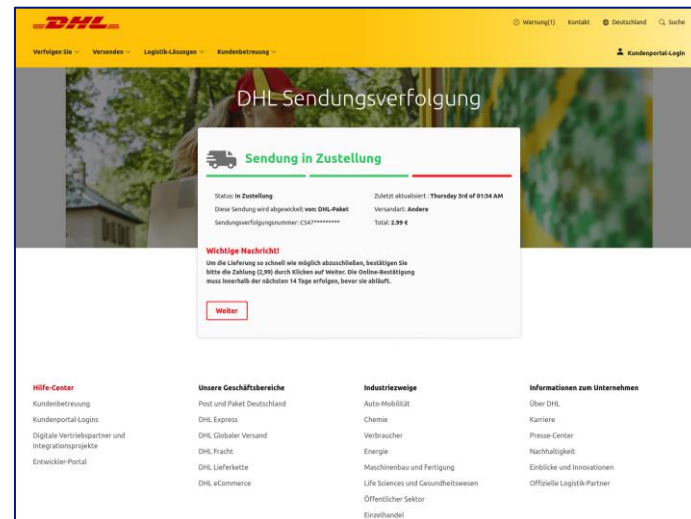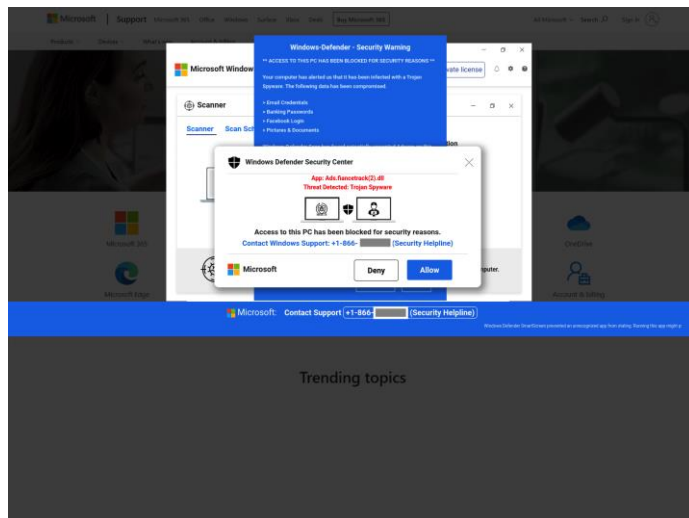
**Takashi Koide,** Naoki Fukushi, Hiroki Nakano, Daiki Chiba

NTT Security (Japan) KK

# Phishing Sites

- Employ social engineering techniques

- Impersonate brands

- **Disguise legitimate domain names**

# Domain Squatting

**Typosquatting:** Typing errors

e.g., examp**k**e[.]com

**Combo-squatting:** Including brand names

e.g., example**-login**[.]com

**Deceptive subdomain:** Using the legitimate domain name in a subdomain

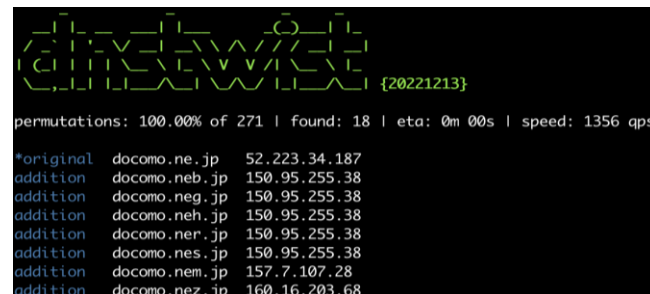e.g., example.com.**malicious[.]example**

**Homograph attack:** Replacing characters with similar ones

e.g., exampl**é**[.]com

# **Mitigating Domain Squatting**

1. Generating domain squatting candidates

   - Rule-based systems (e.g. dnstwist)

   - Machine learning-based systems



2. Detecting domain names

   - Matching domain names with the feed of created candidates

     › Certificate Transparency (CT) Logs

     › Passive DNS traffic

     › Other domain name feeds (e.g., TLD zone files)

# Generated Squatting Domains (GSDs)

- Combining multiple squatting techniques to create domain names

- Registering thousands of domain names simultaneously

- Vague similarities can be seen when GSDs are listed

**appleid.apple[.]com:**
www.appleid-signinmy.gsmserver-pro[.]com
www.applesign-in.gsmserver-pro[.]com
www.applesign-us.gsmserver-pro[.]com
www.support-appleid.gsmserver-pro[.]com

**login.microsoftonline[.]com:**
login-micro-online-doc-file-share-view.web[.]app
micro-login-drive-file-share-view-doc.web[.]app
micro-login-drive-file-share-view.web[.]app
micro-login-file-share-drive-view-doc.web[.]app

**www.amazon.co[.]jp:**
www.amaczon-co-jp.amazccn.bwyver[.]top
www.amaeozn-co-jp.amazecn.ibsmoa[.]top
www.amazcon-co-jp.amacszan.cwheoj[.]top
www.amazeon-co-jp.amazom.cvcvjj[.]top

**www.coingecko[.]com:**
www.coingeicko[.]click
www.coingjecko[.]click
www.coinsgeckko[.]click
www.coinsgqecko[.]click

**rakuten.co[.]jp:**
rakoten.co.ip.enxazgii[.]cf
rakoten.co.ip.enxazgii[.]ml
rakoten.co.ip.ciqrjrzk[.]ga
rakoten.co.ip.ciqrjrzk[.]tk

**steamcommunity[.]com:**
steamcammunnittly[.]com
steamcornmunitty[.]com
steamcoormmunity[.]ru
steamcornmunnity[.]ru

**lcloud[.]com/find:**
www.findmy.lcloud-online[.]in
www.findmy.phone-cloud-mx[.]info
www.findmy.phone-lcloud[.]info
www.findmy.phone-lcloud[.]top

**coinbase[.]com**
ccoiasbasvelog.azurewebsites[.]net
coaoiasnbaselog.azurewebsites[.]net
coiansabsabelog.azurewebsites[.]net
coinnnbaswalle.azurewebsites[.]net

**kucoin[.]com**
kuuucoinesslugincess.godaddysites[.]com
kuuucoinessslugincess.godaddysites[.]com
kuuucoinesslugincessss.godaddysites[.]com
kuuucoinessslugincesss6.godaddysites[.]com

# Generated Squatting Domains (GSDs)

Creation process of GSDs using multiple squatting techniques

www.amazon.co[.]jp (legitimate domain name)

**Typosquatting**

www.amazeon.co[.]jp

**Deceptive subdomains**

www.amazeon.co.jp.example[.]top

**Combosquatting**

www.amazeon-co-jp.example[.]top

**Deploy a phishing site**

# Challenges in Detecting GSDs



- **Evasion Tactics:** GSDs evade existing squatting detection systems

- **Rapid Evolution:** The emergence of new patterns is frequent, making manual rule creation impractical

- **False Positives:** Existing ML-based systems often generate many false positives

# PhishReplicant: Key Ideas

- Identifying domain names that resemble known phishing domain names rather than legitimate ones

- Using a fine-tuned Sentence-BERT model for measuring domain name similarity

- Detecting GSDs from newly registered or observed domain names compared using the latest phishing threat intelligence

# PhishReplicant: System Architecture

9

# Step 1: Extract Similar Domain Names

Phishing Threat Intelligence

Filtering Out

Feature Extraction

Clustering

Phishing Domain Clusters

Step 1

Domain Name Feeds

Passive D...

Matching

Gen...

Step 2

login-instagram-baitout.blogspot[.]kr

www2.amazaon.co.jp.login.lowv[.]cn

www.eki-net.con-aesceeeeesoas.ouvpkl[.]top
www.ekl-net.com-asoaeeccesaa.wbioif[.]top
www.ekl-net.conn-aocseccesas.svmm[.]top

login.miizuhabunks-co-jp[.]cyou
miizuho-bunksong-co-jp[.]cyou
login.miizuho-bunksong-co-jp[.]cyou
miizuhe-bunknet-co-jp[.]cyou

www.macesaoeod.of6jh4[.]icu
www.maceseoeod.r7302f[.]icu
www.macesarrod.mfxzq4[.]icu
www.macesaoeod.4g871x[.]icu

facebook-checkpoint-policy-violation-social-network16[.]ml

# Step 1: Extract Similar Domain Names

Extracting a set of similar domain names from known phishing domain names

**Data Sources**

- Phishing threat intelligence (TI): PhishTank, OpenPhish, CrowdCanary [ARES '23]

**Feature Extraction**

- Use Sentence-BERT for text embeddings of domain names

**Clustering Technique**

- Employ DBSCAN algorithm, utilizing cosine similarity

**Matching Rule Generation**

- Create rules based on TLDs, e2LDs, character count

# Step 2: Detect GSDs

www.macesaoeod.of6jh4[.]icu
www.maceseoeod.r7302f[.]icu
www.macesarrod.mfxzq4[.]icu
www.macesaoeod.4g871x[.]icu

...
blueriveradventures.example
sunnydaydreams.example
starrynightdelights.example
quicksilverexpress.example
mountainpeakretreats.example
**www.maceseroeod.vjhch4[.]icu**
urbanexplorerpros.example
techsavvysolutions.example
goldenharvestgardens.example
cosmicadventuresunited.example
aquamarineoasis.example
...

Domain Name Feeds

Passive DNS    CT Logs    Registered Domain Lists

Feature Extraction

Filtering Out

Calculating Similarity

Matching Rules

Generated Squatting Domains

Step 1

Step 2

# Step 2: Detect GSDs

Identifying GSDs similar to known phishing domain names

**Data Sources**

- Passive DNS traffic, CT logs, registered domain name list (e.g., zonefiles.io)

**Feature Extraction**

- (In the same way as Step 1)

**Calculating Similarity**

- Calculate cosine similarity between each cluster (Step 1) and domain names

- Extract domain name if the similarity exceeds the 0.96 threshold

**Match rules**

- If the domain name follows Step 1 cluster rules, label it as a GSD

# Evaluation of Real-time GSD Detection

**Data Sources**

- CT Logs, zonefiles.io, Passive DNS traffic: 28 days (November 2022)

**Verification of detected GSDs**

- **URL Inspection Services:** VirusTotal, URLScan, Google Safe Browsing

- **Phishing TI:** OpenPhish, PhishTank, CrowdCanary

- **Web Crawling:** Crawl websites + Phishpedia [USENIX '21])

- **Passive DNS:** Checking IP address sharing

- **Manual Validation**

# Evaluation of Real-time GSD Detection

## Results

- **92.4% Precision**: 3,498 true positives out of 3,784 detected GSDs

- **74.6% Phishing Sites**: 2,821 GSDs were used as phishing sites

- **59.0% Zero-day Phishing**: 2,233 GSDs were exclusive discoveries not identified by other services

### GSDs Detected by PhishReplicant

| Manual Validation | VirusTotal | URLScan | GSB | Phishing TI | Crawling | Passive DNS |
|---|---|---|---|---|---|---|
| 3,498 | 934 | 433 | 564 | 430 | 757 | 2,106 |

# Evaluation with Baseline Systems

**Data Sources**

- CT Logs, zonefiles.io, Passive DNS traffic: 31 days (March 2023)

**Baseline Systems**

- **dnstwist**: Generates squatting domains (rule-based)

- **Phishing Catcher**: Identifies phishing domains from certificates (rule-based)

- **StreamingPhish**: Detects phishing domains (ML-based)

- **Ctl-pipeline**: Identifies phishing domains from certificates (ML-based)

**Verification**

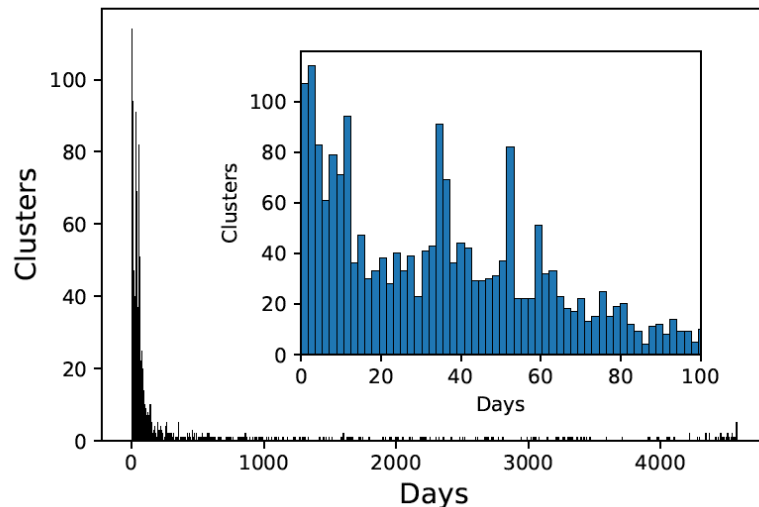- Phishing TI

- Google Safe Browsing

# Evaluation with Baseline Systems

**Results**

- PhishReplicant outperformed baseline systems with 26.1% accuracy

| System | Detected Domains | Matched with Phishing TI or GSB | Ratio |
|--------|------------------|---------------------------------|-------|
| **dnstwist** | 352,294 | 3,645 | 1.0% |
| **Phishing Catcher** | 98,326 | 705 | 0.7% |
| **StreamingPhish** | 196,677 | 3,770 | 1.9% |
| **Ctl-pipeline** | 50,441 | 201 | 0.4% |
| **PhishReplicant** | 7,358 | 1,923 | **26.1%** |

# In-depth Analysis of GSDs

Discovered 205,158 GSDs (2,842 clusters) over 150 days (since August 2022)

**Active Duration of GSD Clusters:**

- Median active period recorded at 41 days (between the first-seen and last-seen dates)

- Over 97% of clusters remained active for more than 24 hours



**IP Address Sharing:**

- 65% of clusters shared 1 or 2 IPs each

18

# In-depth Analysis of GSDs

**Phishing Targeted Brands:**

- 265 brands in 35 countries imitated by 165,643 domain names

- Credit Card category was the most targeted (70.9% of all domain names)

Top 10 categories

| Category | # of domains |
|---|---|
| Credit Card | 117,454 |
| Logistics | 17,943 |
| Telecommunications | 17,232 |
| Social Networks | 2,931 |
| Bank | 2,740 |
| Crypto | 1,842 |
| Software | 1,768 |
| E-commerce | 1,195 |
| Government | 770 |
| News | 507 |
| Other | 1,261 |

# In-depth Analysis of GSDs

**Geographical Distribution:**

- 90.3% of domain names targeted Japanese brands.

- U.S. brands were more commonly targeted (69.7%) in global phishing trends

<table>
<tr><td colspan="2" align="center"><b>Top 10 countries</b></td></tr>
<tr><td>Country</td><td># of domains</td></tr>
<tr><td>Japan</td><td>149,529</td></tr>
<tr><td>United States</td><td>12,188</td></tr>
<tr><td>France</td><td>683</td></tr>
<tr><td>United Kingdom</td><td>604</td></tr>
<tr><td>Spain</td><td>409</td></tr>
<tr><td>China</td><td>321</td></tr>
<tr><td>Turkey</td><td>305</td></tr>
<tr><td>Italy</td><td>209</td></tr>
<tr><td>Poland</td><td>197</td></tr>
<tr><td>Colombia</td><td>152</td></tr>
<tr><td>Other</td><td>1,046</td></tr>
</table>

### Top 10 brands

| Brand | Country | # of domains |
|---|---|---|
| Credit card A | Japan | 88,970 |
| Logistics A | Japan | 17,345 |
| Telecommunication A | Japan | 14,905 |
| Credit card B | Japan | 14,631 |
| Credit card C | Japan | 10,900 |
| Social networks A | United States | 2,175 |
| Credit card D | Japan | 1,526 |
| Crypto wallet A | United States | 1,454 |
| Telecommunication B | United States | 1,316 |
| E-commerce A | United States | 956 |

# Conclusion

**System Proposal:** Introducing a system to detect Generated Squatting Domains (GSDs)

github.com/tkoide398/PhishReplicant

**System Effectiveness:** PhishReplicant employs a transformer-based language model to identify domain names similar to known phishing domain names and outperformed existing systems

**Real-Time Detection:** Offering timely GSD detection using the latest phishing TI and new domain name feeds, enabling early countermeasures against newly emerged phishing sites