# Fighting Phishing by Introducing WikiPhish

### A new public dataset based on Wikipedia for legit URLs

**Gabriel Loiseau**

Hornet Security

HORNETSECURITY

vade

WikiPhish Paper

# WikiPhish: A Diverse Wikipedia-Based Dataset for Phishing Website Detection

Data/Toolset paper

Gabriel Loiseau
Hornet Security
Hem, France
gabriel.loiseau@hornetsecurity.com

Valentin Lefils
Hornet Security
Hem, France
valentin.lefils@hornetsecurity.com

Maxime Meyer
Hornet Security
Hem, France
maxime.meyer@hornetsecurity.com

Damien Riquet
Hornet Security
Hem, France
damien.riquet@hornetsecurity.com

https://doi.org/10.1145/3626232.3653283

# WikiPhish Team

HORNETSECURITY
vade

Gabriel Loiseau
Research Scientist

Valentin Lefils
Data Scientist

Maxime Meyer
Head of Research

Damien Riquet
Lead Research
Engineer

# Outline

- Background
- The WikiPhish Dataset
- Importance of Legitimate URL Diversity
- Training and Evaluating Classifiers
- Conclusion

# Background: phishing detection using machine learning

## 2023 CRIME TYPES

| By Complaint Count | | | |
|---|---|---|---|
| **Crime Type** | **Complaints** | **Crime Type** | **Complaints** |
| Phishing/Spoofing | 298,878 | Other | 8,808 |
| Personal Data Breach | 55,851 | Advanced Fee | 8,045 |
| Non-payment/Non-Delivery | 50,523 | Lottery/Sweepstakes/Inheritance | 4,168 |
| Extortion | 48,223 | Overpayment | 4,144 |
| Investment | 39,570 | Data Breach | 3,727 |
| Tech Support | 37,560 | Ransomware | 2,825 |
| BEC | 21,489 | Crimes Against Children | 2,361 |
| Identity Theft | 19,778 | Threats of Violence | 1,697 |
| Confidence/Romance | 17,823 | IPR/Copyright and Counterfeit | 1,498 |
| Employment | 15,443 | SIM Swap | 1,075 |
| Government Impersonation | 14,190 | Malware | 659 |
| Credit Card/Check Fraud | 13,718 | Botnet | 540 |
| Harassment/Stalking | 9,587 | | |
| Real Estate | 9,521 | | |

FBI Internet Crime Report 2023

# Background: phishing detection using machine learning

## 2023 CRIME TYPES

### By Complaint Count

| Crime Type | Complaints | Crime Type | Complaints |
|---|---|---|---|
| Phishing/Spoofing | 298,878 | Other | 8,808 |
| Personal Data Breach | 55,851 | Advanced Fee | 8,045 |
| Non-payment/Non-Delivery | 50,523 | Lottery/Sweepstakes/Inheritance | 4,168 |

"phishing detection" machine learning

Scholar    About 869 results (0.04 sec)    YEAR ▾

Since 2024

| | | | |
|---|---|---|---|
| Employment | 15,443 | SIM Swap | 1,075 |
| Government Impersonation | 14,190 | Malware | 659 |
| Credit Card/Check Fraud | 13,718 | Botnet | 540 |
| Harassment/Stalking | 9,587 | | |
| Real Estate | 9,521 | | |

FBI Internet Crime Report 2023

# The Dataset

- WikiPhish is a new open-access dataset for phishing website classification

- 110,606 URLs, HTML web pages, and screenshots

- Benign samples: **Wikipedia references**

- Phishing samples: OpenPhish, and PhishTank

# WikiPhish Advantages

**Diversity**

Wide variety of benign and phishing samples

**Up-to-date samples**

Recent phishing campaigns and benign documents

**Transparency**

Reproducible and updatable

# Data Collection Challenges

# Data Collection Challenges

- Evolving ecosystem due to concept drift
  - Phishing detection is truly adversarial

# Data Collection Challenges

- Evolving ecosystem due to concept drift
  - Phishing detection is truly adversarial

- Phishing websites are often short-lived
  - Phishers quickly remove their websites
  - URLs are quickly blocked by providers

# Data Collection Challenges

- Evolving ecosystem due to concept drift
  - Phishing detection is truly adversarial

- Phishing websites are often short-lived
  - Phishers quickly remove their websites
  - URLs are quickly blocked by providers

- Collecting relevant legitimate documents is difficult
  - New frameworks, development practices

# Limitation of Existing Datasets

# Limitation of Existing Datasets

- Many datasets are limited to URLs, omitting HTML and visual content
  - Limited features for machine learning

# Limitation of Existing Datasets

- Many datasets are limited to URLs, omitting HTML and visual content
  - Limited features for machine learning

- Few datasets serve as benchmarks for phishing detection
  - Studies are conducted using their own dataset

# Limitation of Existing Datasets

- Many datasets are limited to URLs, omitting HTML and visual content
  - Limited features for machine learning

- Few datasets serve as benchmarks for phishing detection
  - Studies are conducted using their own dataset

- Standardization is limited, making comparison difficult
  - Dataset sizes, collection protocol, and time periods are different
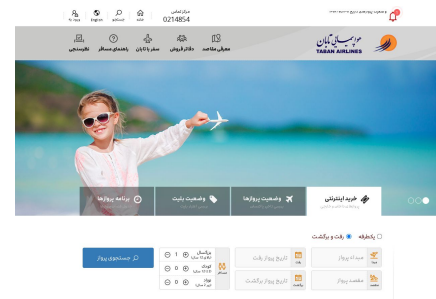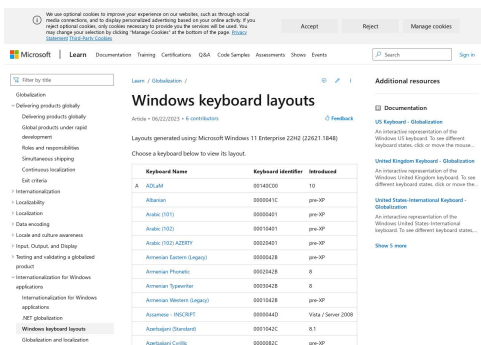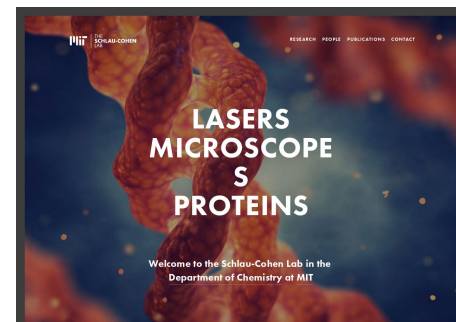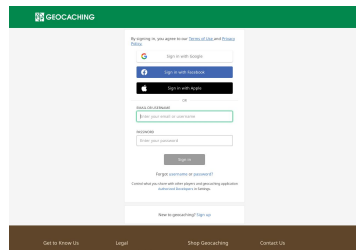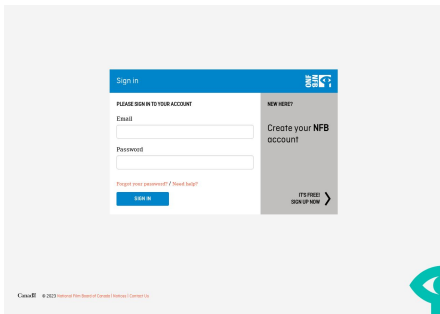  - Resulting machine learning models are hardly comparable

# WikiPhish Collection (benign)

- Extract URLs from random Wikipedia page references

- Add login page by appending "*/login*" or "*/signin*" subpath

- Includes lesser-known and "out-of-distribution" websites.

## References  [edit]

1. ^ Jansson, K.; von Solms, R. (2011-11-09). "Phishing for phishing awareness". *Behaviour & Information Technology*. 32 (6): 584–593. doi:10.1080/0144929X.2011.632650. ISSN 0144-929X. S2CID 5472217.

2. ^ Ramzan, Zulfikar (2010). "Phishing attacks and countermeasures". In Stamp, Mark; Stavroulakis, Peter (eds.). *Handbook of Information and Communication Security*. Springer. ISBN 978-3-642-04117-4.

3. ^ "Internet Crime Report 2020" (PDF). *FBI Internet Crime Complaint Centre*. U.S. Federal Bureau of Investigation. Retrieved 21 March 2021.

4. ^ Ollmann, Gunter. "The Phishing Guide: Understanding and Preventing Phishing Attacks". *Technical Info*. Archived from the original on 2011-01-31. Retrieved 2006-07-10.

5. ^ *a b c* Wright, A; Aaron, S; Bates, DW (October 2016). "The Big Phish: Cyberattacks Against U.S. Healthcare Systems". *Journal of General Internal Medicine*. 31 (10): 1115–8. doi:10.1007/s11606-016-3741-z. PMC 5023604. PMID 27177913.

6. ^ Stonebraker, Steve (January 2022). "AOL Underground". *aolunderground.com* (Podcast). Anchor.fm.

7. ^ Mitchell, Anthony (July 12, 2005). "A Leet Primer". TechNewsWorld. Archived from the original on April 17, 2019. Retrieved 2021-03-21.

8. ^ "Phishing". *Language Log, September 22, 2004*. Archived from the original on 2006-08-30. Retrieved 2021-03-21.

9. ^ Jøsang, Audun; et al. (2007). "Security Usability Principles for Vulnerability Analysis and Risk Assessment". *Proceedings of the Annual Computer Security Applications Conference 2007 (ACSAC'07)*. Archived from the original on 2021-03-21. Retrieved 2020-11-11.

10. ^ Lin, Tian; Capecci, Daniel E.; Ellis, Donovan M.; Rocha, Harold A.; Dommaraju, Sandeep; Oliveira, Daniela S.; Ebner, Natalie C. (September 2019). "Susceptibility to Spear-Phishing Emails: Effects of Internet User Demographics and Email Content". *ACM Transactions on Computer-Human Interaction*. 26 (5): 32. doi:10.1145/3336141. ISSN 1073-0516. PMC 7274040. PMID 32508486.

# WikiPhish Collection (benign screenshots)

# WikiPhish Collection (phishing)

- We monitor the public phishing databases OpenPhish and PhishTank from January 11, 2023, to October 22, 2023

- We collect every page as soon as they are listed, and only add them to the dataset if they are further validated by moderators

# WikiPhish Collection (preprocessing)

- We limit the number of items per Fully Qualified Domain Name (FQDN) to 10 to limit redundancy and enhance diversity

| Legitimate | Occurrences | Phishing | Occurrences |
|---|---|---|---|
| web.archive.org | 10,398 | storageapi.fleek.co | 402 |
| geohack.toolforge.org | 9,236 | ipfs.io | 167 |
| books.google.com | 5,309 | ipfs.fleek.co | 63 |
| scholar.google.com | 4,522 | s3.amazonaws.com | 52 |
| www.wikidata.org | 3,857 | tinyurl.com | 45 |
| doi.org | 3,177 | storageapi-stg.fleek.co | 45 |
| www.jstor.org | 2,935 | go-citien.duckdns.org | 39 |
| www.worldcat.org | 2,584 | dev.awit.ae | 37 |
| pubmed.ncbi.nlm.nih.gov | 2,234 | v.ht | 32 |
| www.imdb.com | 1,958 | s.id | 30 |
| ... | ... | ... | ... |
| Total before filtering | 268,333 | | 29,222 |
| Total after filtering | 87,563 | | 23,043 |

# Importance of Legitimate URL Diversity

- Comparison with collection strategies of two other datasets (Apruzzese et al. 2022, Ariyadasa et al. 2021)
  - Top, middle, and bottom part of **Alexa**
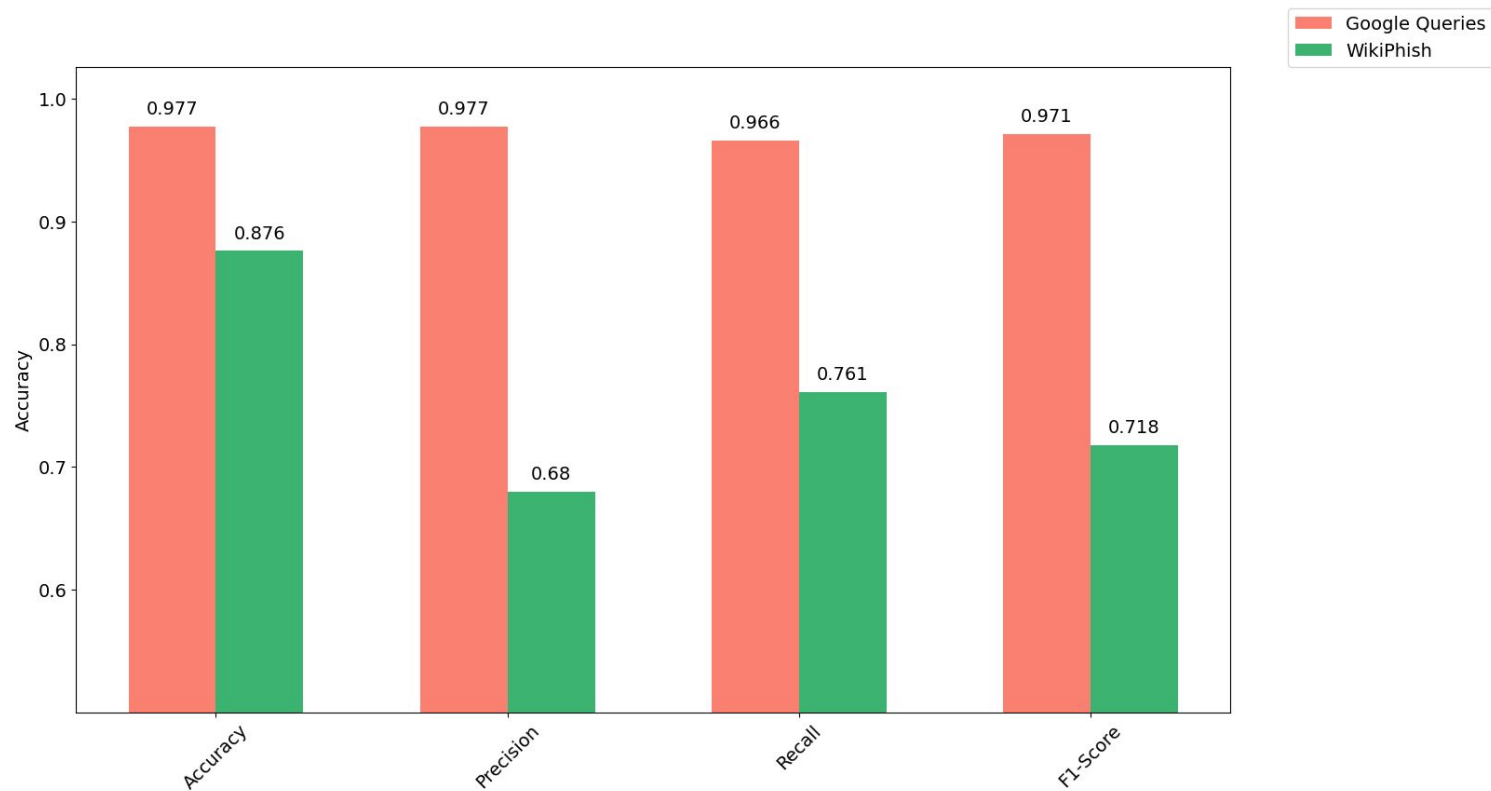  - Keywords used in Google search engine: "**Google Queries**"

| Feature Index | Description |
|---|---|
| 1 | URL subdomains count |
| 2 | URL length |
| 3 | Count of dot (.) symbols in URL |
| 4 | Count of at (@) symbols in URL |
| 5 | Count of hyphens (-) symbols in URL |
| 6 | Count of underscore (_) symbols in URL |
| 7 | Count of slash (/) symbols in URL |
| 8 | Count of www in URL |

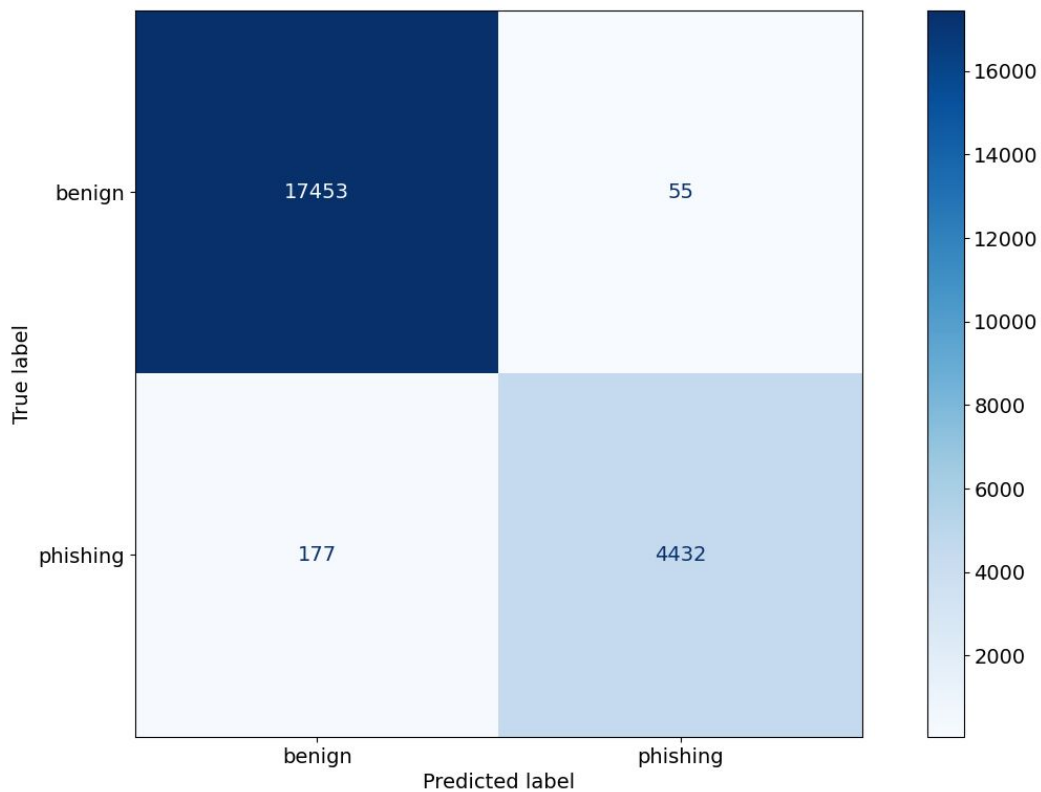| Feature Index | Alexa | Google Queries | WikiPhish |
|---|---|---|---|
| 1 | 1 | $4 \pm 1$ | $5 \pm 1$ |
| 2 | 36 | $219 \pm 5$ | $280 \pm 10$ |
| 3 | 4 | $13 \pm 3$ | $16 \pm 5$ |
| 4 | 1 | $2 \pm 1$ | $5 \pm 1$ |
| 5 | 4 | $30 \pm 2$ | $34 \pm 2$ |
| 6 | 1 | $19 \pm 2$ | $23 \pm 2$ |
| 7 | 1 | $14 \pm 1$ | $18 \pm 1$ |
| 8 | 2 | $3 \pm 1$ | $4 \pm 1$ |
| Average gain from Alexa | | $\times 7$ | $\times 10$ |

# Training and Evaluating Classifiers

- Training and evaluation of a Random Forest model
- 8 URL features, and 11 HTML features
  - 1) Train on older datasets, evaluate on WikiPhish
  - 2) Train and evaluate on WikiPhish

# Train on older datasets, evaluate on WikiPhish

# Train and evaluate on WikiPhish

# Conclusion

- WikiPhish is a new dataset for phishing website detection

- It emphasizes diversity, up-to-date samples, and transparency

- It provides a valuable benchmark for research and development of new phishing detectors

# Thank you!