# SimChar:
# Building a Dataset of
# Visually Similar Characters

Sep 18 2019 @ W3C TPAC Break Out Session (AHA)

Tatsuya Mori

Waseda University

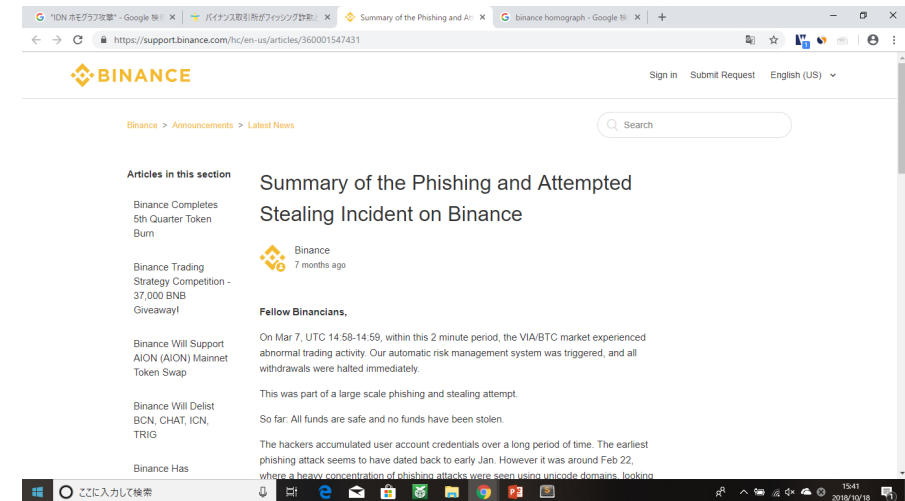# Background

IDN homograph attacks getting vital these days.

Aug 2017

Mar 2018



adobe.com targeted (adobe.com)

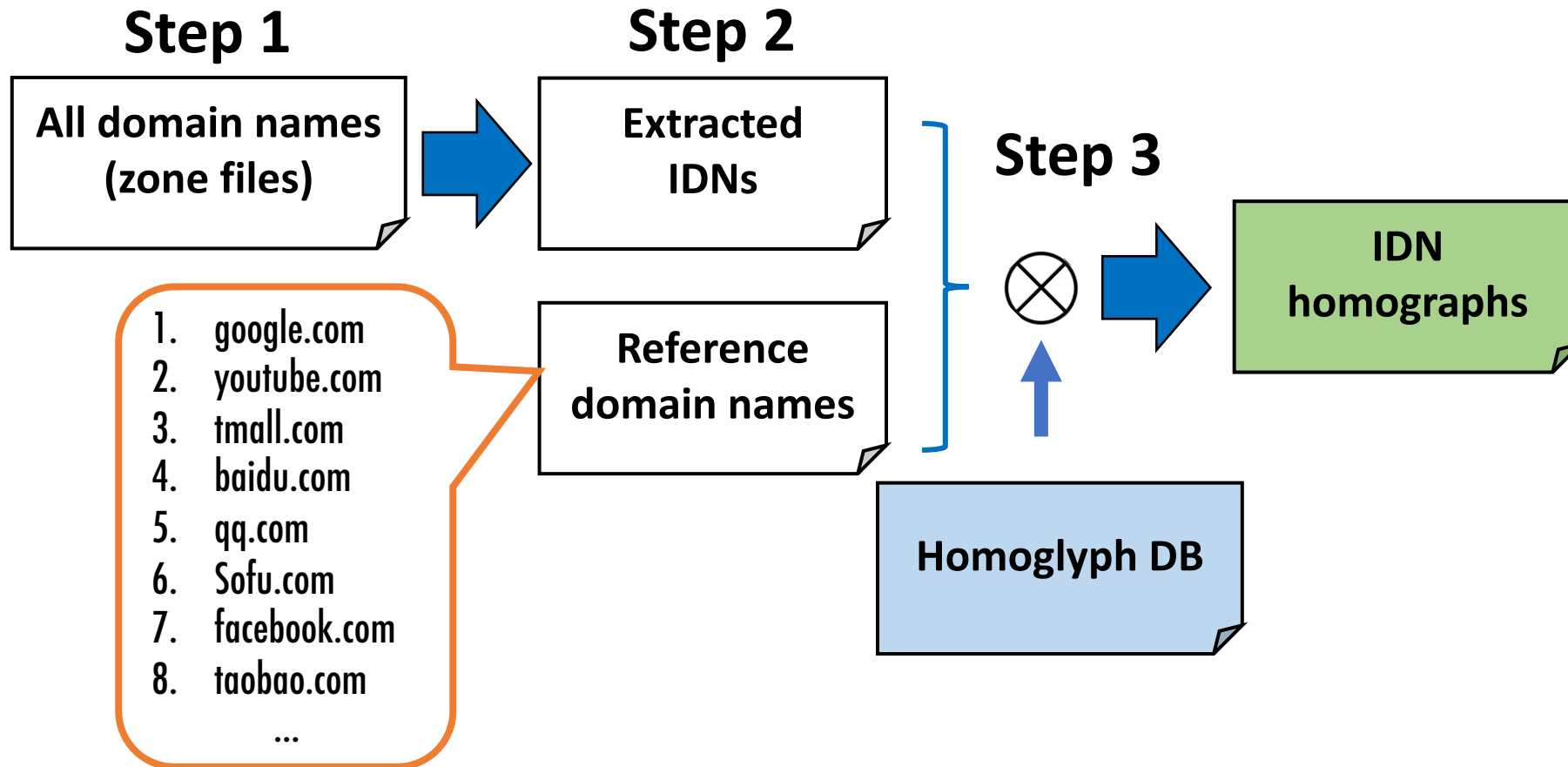Latin small letter b with dot below (U+1E05)



binance.com targeted (binance.com)

Latin small letter i with dot below (U+1ECB)

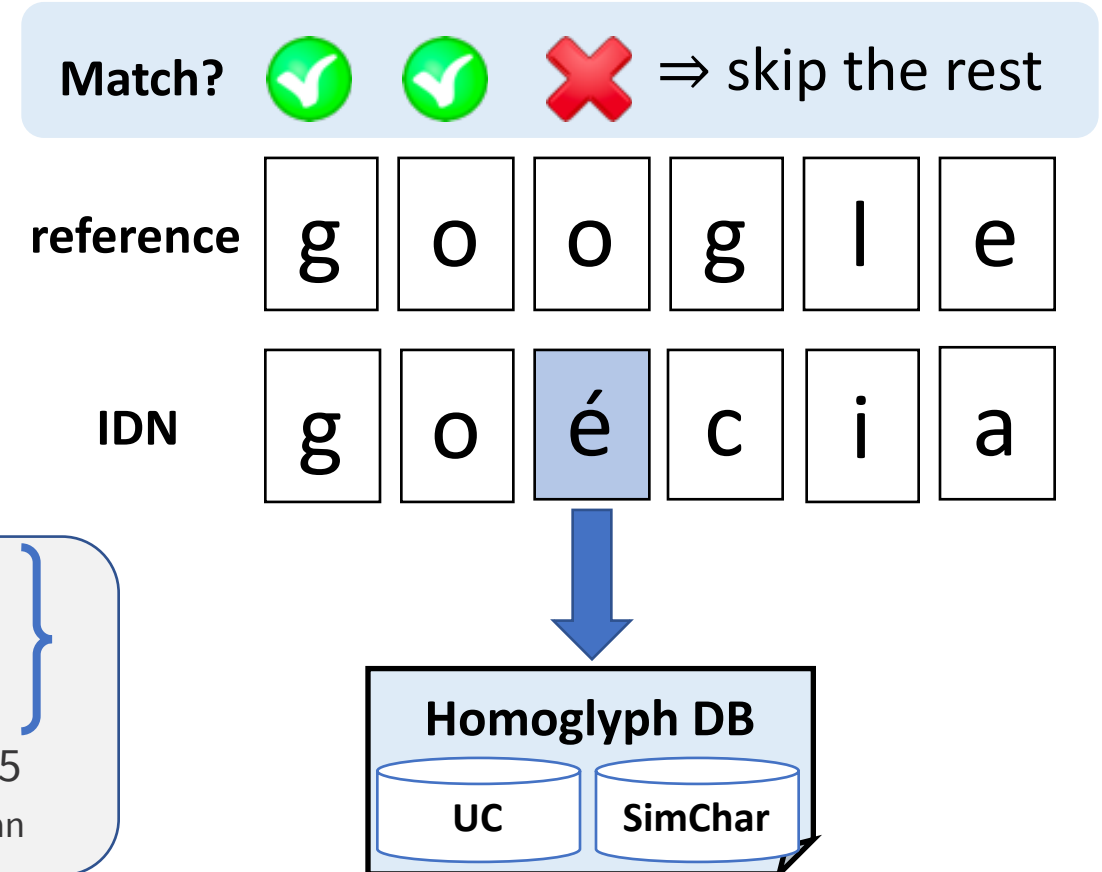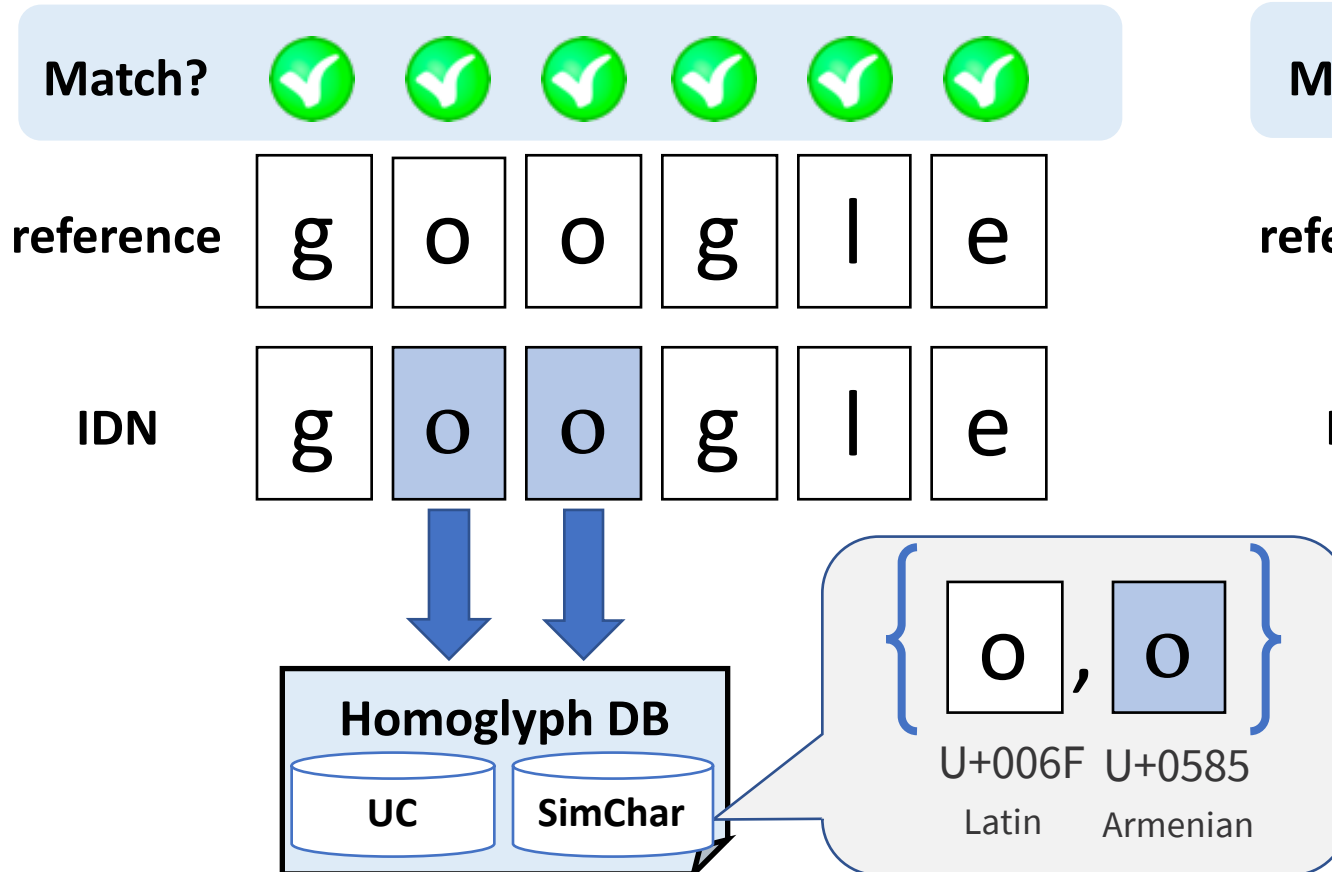Latin small letter a with dot below (U+1EA1)

# ShamFinder

- A framework we built to detect potential IDN homographs automatically.

# ShamFinder

Match? ✅ ✅ ✅ ✅ ✅ ✅

reference: g o o g l e

IDN: g o o g l e

Homoglyph DB
UC | SimChar

{ o , o }
U+006F U+0585
Latin   Armenian

Match? ✅ ✅ ❌ ⇒ skip the rest

reference: g o o g l e

IDN: g o é c i a

Homoglyph DB
UC | SimChar

# Question

- Is there a homoglyph DB out there?

# A solution

- Yes, we can make use of <span style="color:red">confusables.txt</span>

http://unicode.org/reports/tr39/#Data_Collection

**Visually Confusable Characters:**
Provides a mapping for visual confusables for use in detecting possible security problems. The usage of the file is described in *Section 4, Confusable Detection.*

# Unicode Utilities: Confusables

help | character | properties | confusables | unicode-set | compare-sets | regex | bnf-regex | breaks | transform | bidi | bidi-c | idna | langua

**Input**

paypal

Restriction [ IDNA2008 ] [ Show ]

## Confusable Characters

| p | ρ | p | Ⲣ | |
|---|---|---|---|---|
| 0070 | 03C1 | 0440 | 2CA3 | |
| LATIN SMALL LETTER P | GREEK SMALL LETTER RHO | CYRILLIC SMALL LETTER ER | COPTIC SMALL LETTER RO | |
| a | ɑ | α | а | |
| 0061 | 0251 | 03B1 | 0430 | |
| LATIN SMALL LETTER A | LATIN SMALL LETTER ALPHA | GREEK SMALL LETTER ALPHA | CYRILLIC SMALL LETTER A | |
| y | ɣ | Y | γ | у | Y |
| 0079 | 0263 | 028F | 03B3 | 0443 | 04AF |
| LATIN SMALL LETTER Y | LATIN SMALL LETTER GAMMA | LATIN LETTER SMALL CAPITAL Y | GREEK SMALL LETTER GAMMA | CYRILLIC SMALL LETTER U | CYRILLIC SMALL LETTER STRAIGHT U |
| p | ρ | p | Ⲣ | |
| 0070 | 03C1 | 0440 | 2CA3 | |
| LATIN SMALL LETTER P | GREEK SMALL LETTER RHO | CYRILLIC SMALL LETTER ER | COPTIC SMALL LETTER RO | |
| a | ɑ | α | а | |
| 0061 | 0251 | 03B1 | 0430 | |
| LATIN SMALL LETTER A | LATIN SMALL LETTER ALPHA | GREEK SMALL LETTER ALPHA | CYRILLIC SMALL LETTER A | |
| 1 | l | ǀ | ו | ן | ا |
| 0031 | 006C | 01C0 | 05D5 | 05DF | 0627 |
| DIGIT ONE | LATIN SMALL LETTER L | LATIN LETTER DENTAL CLICK | HEBREW LETTER VAV | HEBREW LETTER FINAL NUN | ARABIC LETTER ALEF |

Total raw values: 42,240

## Confusable Results

paypal paypal paypaι paypal paypal paypal paypaι paypal paypal paypaι payp
paypal paypal paypαl paypal paypaι paypal paypal paypal paypaι paypal paypal p
paypɑl paypαι paypɑl paypαι paypal paypal paypaι paypal paypal paypal paypaι p
paypal paypal paypaι paypal paypαι paypal paypal paypal paypaι paypal paypal p
paypαl paypɑl paypαι paypαl paypαι payρal payρal payρaι payρal payρal paypal
payρal payρal payρaι payρal payρaι payρal payρaι payρal payρal payρal payρal p

# Our Question

- Are there homoglyphs that are <span style="color:red">NOT</span> listed in the confusables.txt?

# Answer

- Yes, we found several homoglyphs not listed in the confusables.txt

# The process of building a homoglyph DB (SimChar)

1. Get the visual images of characters (Unicode BMP) by using GNU Unifont.

2. Compute the distance of two characters (images) with the number of different pixels.

3. If the distance is smaller than some threshold, then detect the pair as homoglyph.

# The process of building a homoglyph DB (SimChar)



$\Delta = 0$     $\Delta = 1$     $\Delta = 2$     $\Delta = 3$     $\Delta = 4$     $\Delta = 5$     $\Delta = 6$

# Stats of DBs

| Sets | # Chars | # Pairs |
|------|--------:|--------:|
| IDNA ∩ Unifont12 | 52,457 | n/a |
| UC ∩ Unifont12 | 5,080 | 3,696 |
| SimChar ∩ Unifont12[1] | 12,686 | 13,208 |

# Confusables to Latin letters

**Table 3: Number of homoglyphs of Latin letters (lowercase) contained in SimChar and UC ∩ IDNA.**

| | SimChar | | | | | |
|---|---|---|---|---|---|---|
| | # | | # | | # | |
| 'o' | 40 | 's' | 14 | 'f' | 8 | |
| 'e' | 26 | 'r' | 14 | 'm' | 8 | |
| 'n' | 24 | 'a' | 14 | 'g' | 7 | |
| 'w' | 20 | 'k' | 13 | 'j' | 7 | |
| 'c' | 19 | 't' | 13 | 'p' | 7 | |
| 'l' | 18 | 'z' | 12 | 'x' | 6 | |
| 'u' | 18 | 'd' | 10 | 'q' | 2 | |
| 'h' | 17 | 'y' | 9 | 'v' | 1 | |
| 'i' | 16 | 'b' | 8 | | | |
| Total | | | | 351 | | |

| | UC ∩ IDNA | | | | | |
|---|---|---|---|---|---|---|
| | # | | # | | # | |
| 'o' | 34 | 'c' | 4 | 'p' | 3 | |
| 'l' | 12 | 'd' | 4 | 'x' | 3 | |
| 'y' | 10 | 'g' | 4 | 'j' | 2 | |
| 'i' | 9 | 'f' | 4 | 'n' | 2 | |
| 'u' | 9 | 'a' | 3 | 'z' | 2 | |
| 'w' | 8 | 'b' | 3 | | | |
| 'v' | 6 | 'e' | 3 | | | |
| 's' | 5 | 'h' | 3 | | | |
| 'r' | 5 | 'q' | 3 | | | |
| Total | | | | 141 | | |

# Examples

è é ê ë ē
ė ǝ  ǫ ǫ ɛ  ě ᘓ ɞ ə è
θ ɵ  ɛ ǫ  ë ɑ ǫ è θ

e   ę

e_0           e_1           e_2           e_3           e_4

c o è̃ c o
o c ë̈ o o
B B   ᴓ c o ɛ̧ ǫ
ɛ ê   ẽ è̀ ᴓ

e_5           e_6

# Our Question

- Are the detected homoglyphs really confusable?

# Answer

- Yes, our human study revealed that they are <span style="color:red">more confusable</span> than those contained in confusables.txt!

Q: Are they distinct or confusing?

# Limitations & Future work

- Evaluation used GNU Unifont only
  → Need to extend the evaluation for other font families

- Participants of Human Study were English speakers
  → Need to consider linguistic/cultural spheres

  - Human perception example:
    ぬ vs. め or わ vs. ね are quite distinguishable for Japanese.

# Summary

- ShanFinder is a framework to detect IDN homographs efficiently.
- ShamFinder makes use of SimChar, which is a database of homoglyphs, and confusables.txt
- SimChar contains homoglyphs not listed in the confusables.txt.
- SimChar is available at:
  - https://github.com/shamfinder/shamfinder
  - simchar.json (47MB)
- More technical details are available at arXiv:
  - https://arxiv.org/abs/1909.07539

    The paper will appear at ACM IMC 2019
    https://conferences.sigcomm.org/imc/2019/