

# IDN Visual Security Deep Thinking

**xisigr**  
Feb, 2019



# About me



- Security researcher of Tencent's Xuanwu Lab
  - <https://xlab.tencent.com>
- Author of Web Front-End Hacker's Handbook
  - <https://www.web2hack.org/>
- xisigr@Xeye team
- [twitter.com/xisigr](https://twitter.com/xisigr) ,[weibo.com/xisigr](https://weibo.com/xisigr)



# Internationalized Domain Names (IDNs)



- IDNs
  - Long time ago, domains could only consist of the Latin letters A to Z, digits, and a few other characters from the US-ASCII coded character set.
  - IDNs allow characters from the Universal Character Set Unicode since 2003.
  - In 2018, Unicode 11.0 contains a repertoire of 137374 characters covering 146 modern and historic scripts.
  - In these more than 10 years, IDNs visual confusion security issues has never been interrupted.
  - What do IDNs Mean to you?

# ASCII Domain names vs. Internationalized Domain Names (IDNs)



**www.test.com**

Third Level  
Domain

Second Level  
Domain

Top Level  
Domain

**小明.我爱你**

IDN Second  
Level Domain

IDN Top  
Level Domain

## ASCII

Letters [a-z]  
Digits [0-9]  
Hyphen [-]  
Label length = 63

## ASCII

Letters [a-z]  
Label length = 63

## Unicode

Valid Unicode-Label:  
IDNA2008  
Valid ASCII-Label:  
Punycode, "xn--"

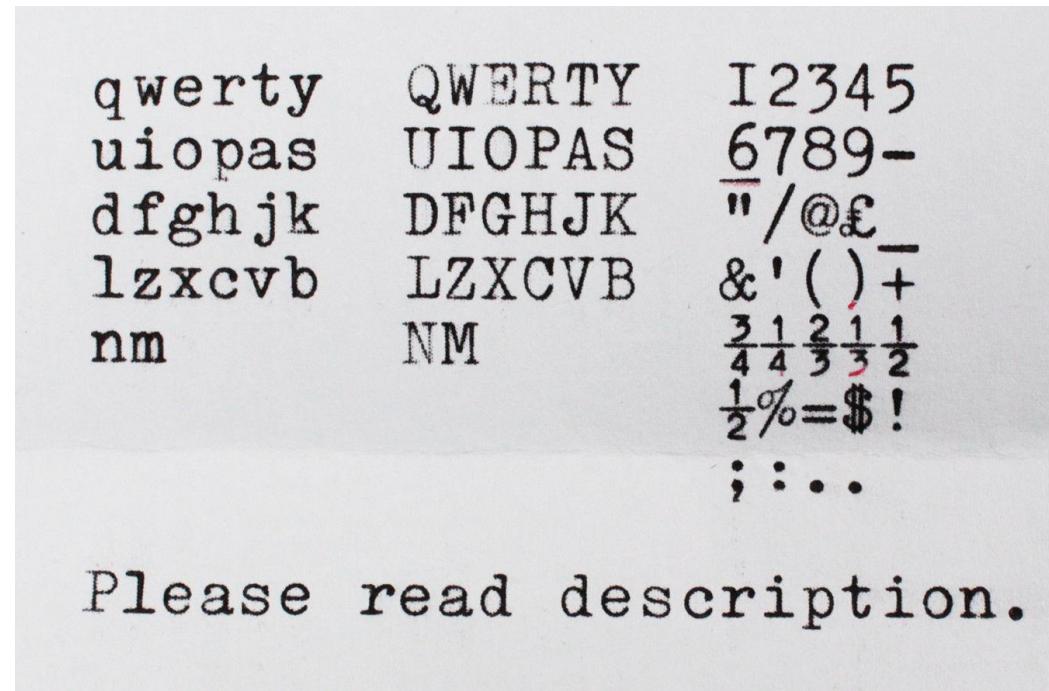
## Unicode

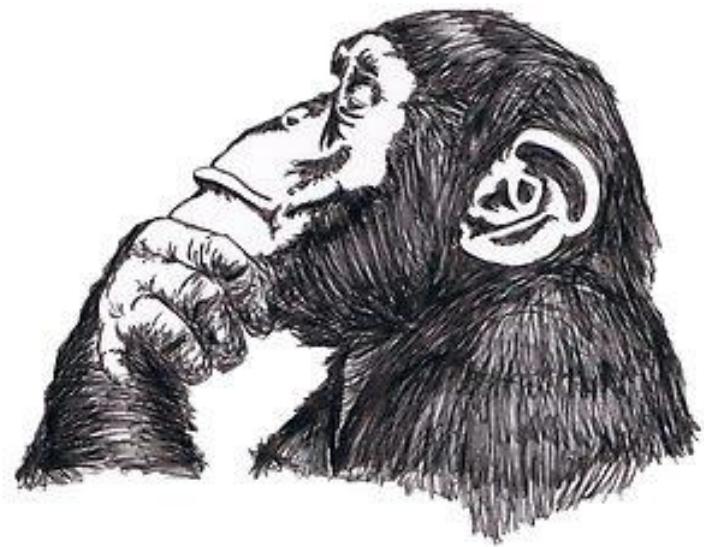
Valid Uicode-Label  
Valid ASCII-label



# Homoglyph Attack blowout

- 0 and O; 1, l and I





# Thinking

# IDN TLD

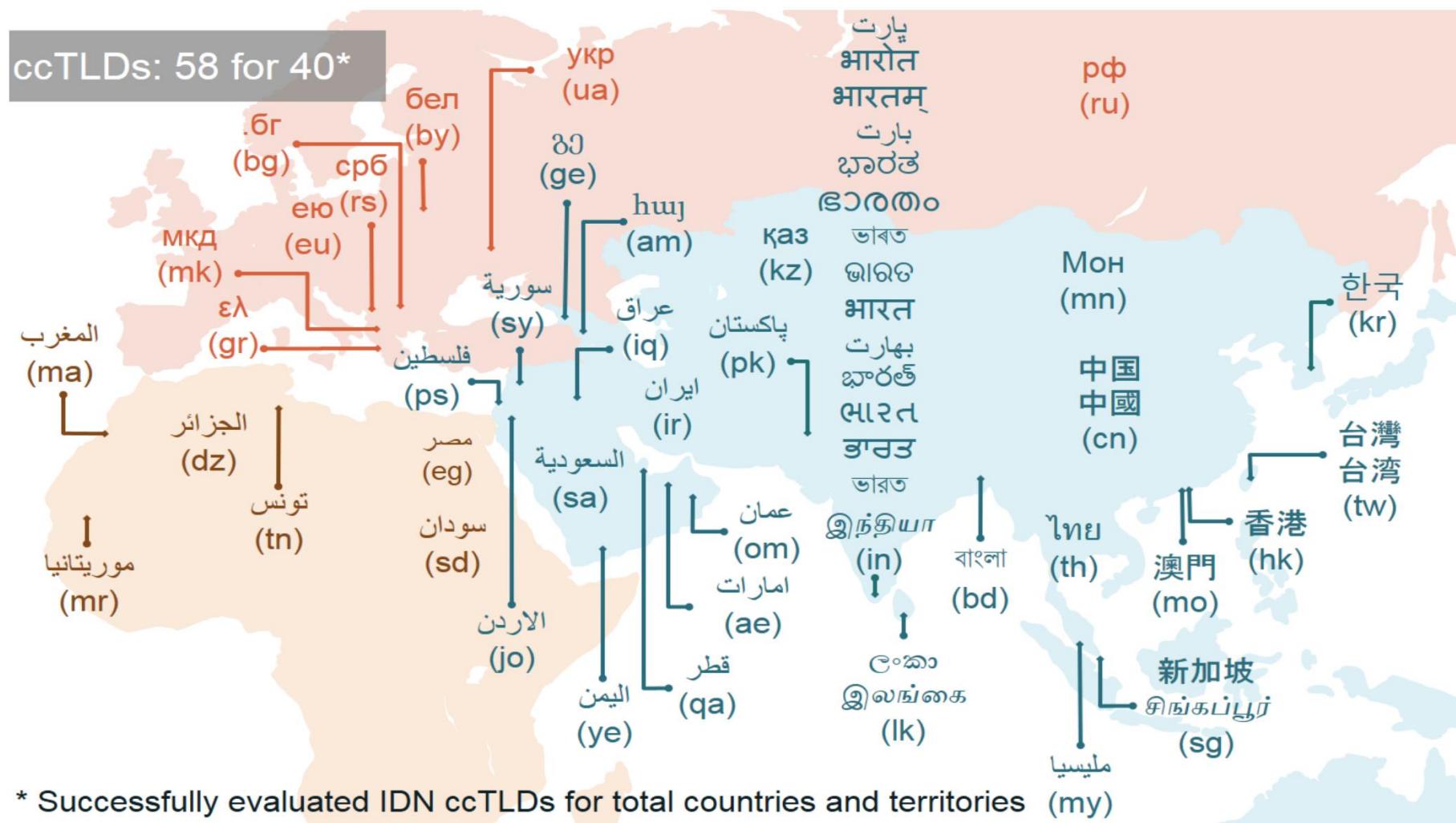


- Until late 2009, TLDs were restricted to only the Latin letters a to z . After 2009, IDN TLDs were introduced in other scripts.
  - 58 IDN ccTLDs evaluated representing 40 countries/territories
  - 56 IDN ccTLDs delegated representing 38 countries/territories
  - Requests cover 33 languages in 19 scripts
- gTLD
  - com, org, net, edu, gov, mil.....
  - 网络(网络传播), 在线(在线), 谷歌(谷歌), 游戏.....(IDN gTLD)
- ccTLD
  - cn, jp, nz, hr, be, cc.....
  - 毛里塔尼亚(毛里塔尼亚), 新加坡(新加坡), 韩国(韩国), 苏丹(苏丹).....(IDN ccTLD)

# IDN Country Code Top-Level Domains



## ccTLDs: 58 for 40\*



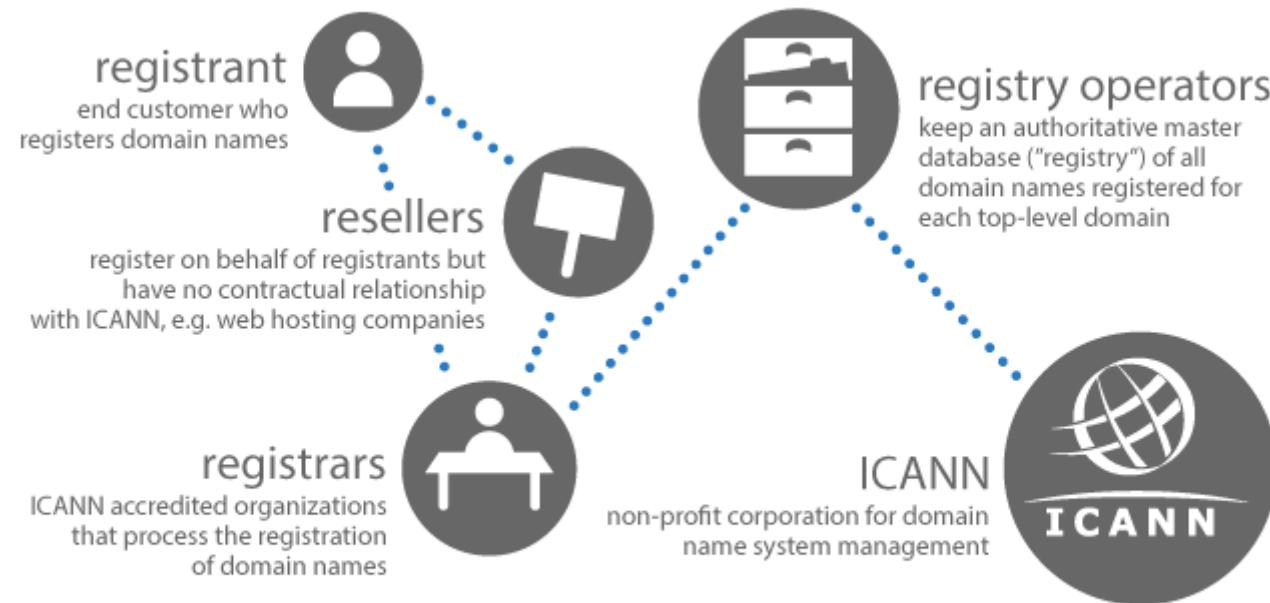
\* Successfully evaluated IDN ccTLDs for total countries and territories (my)

# Think about it



- IDN TLD has been around for ten years (2009-2019). With the emergence and increasing of IDN TLD, What changes to IDN visual security?
- Whether Emoji is allowed in IDN or in IDN TLD?

# Domain Name Registration Process



---

domain registry process

<https://whois.icann.org/en/domain-name-registration-process>

# Think about it



- In the whole domain name registration process, which links are relatively weak will help us to hunt IDN spoof.

# IDNs Registration Rules



- gTLD
  - .COM .NET
    - [https://www.verisign.com/en\\_US/channel-resources/domain-registry-products/idn/idn-policy/registration-rules/index.xhtml](https://www.verisign.com/en_US/channel-resources/domain-registry-products/idn/idn-policy/registration-rules/index.xhtml)
    - .....
    - .....
  - .....
- ccTLD
  - .السعودية (Saudi Arabia.)
    - [http://www.nic.sa/en/view/writing\\_arabic\\_idn\\_guideline](http://www.nic.sa/en/view/writing_arabic_idn_guideline)
  - Association's Hebrew IDN rules
    - [https://www.isoc.org.il/files/docs/ISOC-IL\\_Registration\\_Rules\\_v1.5\\_ENGLISH\\_-26.6.2016.pdf](https://www.isoc.org.il/files/docs/ISOC-IL_Registration_Rules_v1.5_ENGLISH_-26.6.2016.pdf)
    - .....
    - .....

# .COM IDNs Registration Rules



- 1. IETF Standards
  - Compliance with all of the RFC that comprise the IDNA2008 standard.
- 2. Restrictions on Specific Languages
  - All IDN registrations require a 3 letter Language Tag.
- 3. Restrictions On Commingling Of Scripts
  - All code points within an IDN must come from the same Unicode script
  - .....
- 4. ICANN's Restricted Unicode Points
  - .....
- 5. Special Characters
  - .....

# Unicode Scripts And Associated Code Points



Arabic	Georgian	Latin	Rejang
Armenian	Glagolitic	Lepcha	Runic
Avestan	Greek	Limbu	Samaritan
Balinese	Gujarati	Lisu	Saurashtra
Bamum	Gurmukhi	Lycian	Sinhala
Batak	Han	Lydian	Sundanese
Bengali	Hangul	Malayalam	Syloti Nagri
Bopomofo	Hanunoo	Mandaic	Syriac
Brahmi	Hebrew	Meetei Mayek	Tagalog
Buginese	Hiragana	Mongolian	Tagbanwa
Buhid	Imperial Aramaic	Myanmar	Tai Le
Canadian Aboriginal	Inscriptional Pahlavi	New Tai Lue	Tai Tham
Carian	Inscriptional Parthian	Nko	Tai Viet
Cham	Javanese	Ogham	Tamil
Cherokee	Kaithi	Ol Chiki	Telugu
Coptic	Kannada	Old Persian	Thaana
Cuneiform	Katakana	Old South Arabian	Thai
Cyrillic	Kayah Li	Old Turkic	Tibetan
Devanagari	Kharoshthi	Oriya	Tifinagh
Egyptian Hieroglyphs	Khmer	Phags Pa	Vai
Ethiopic	Lao	Phoenician	Yi

# Think about it

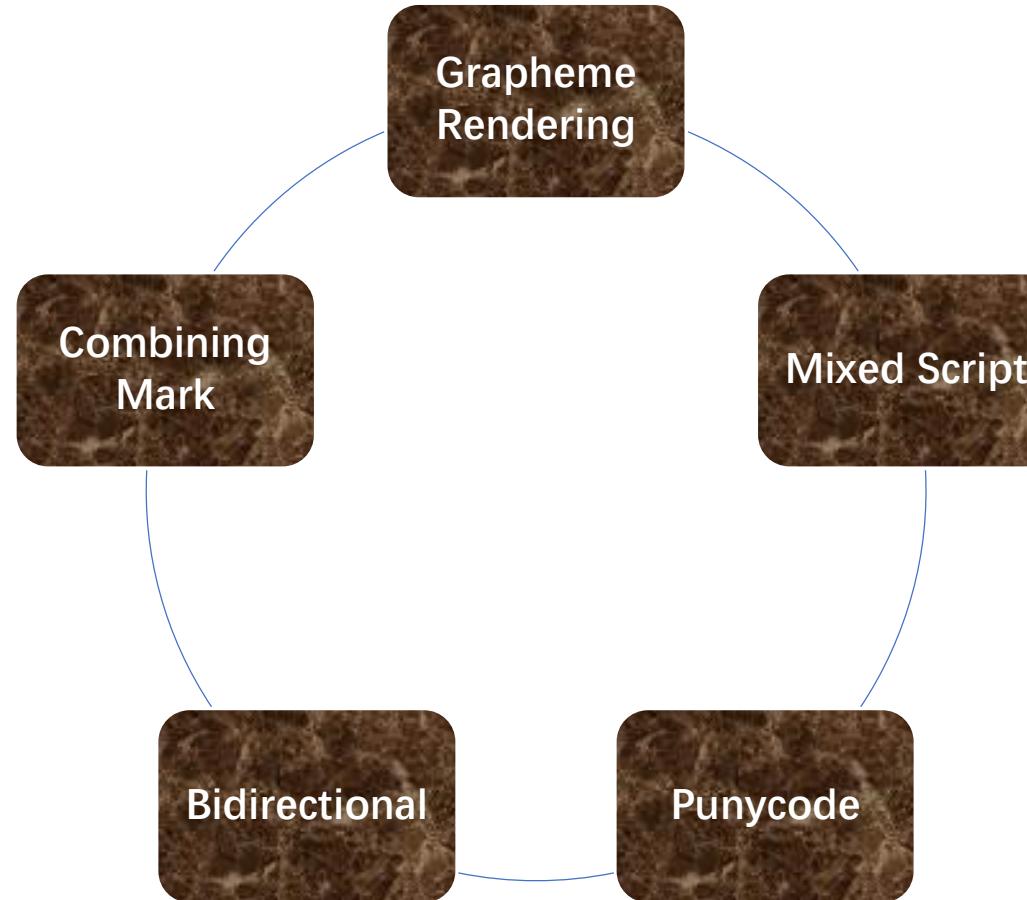


- How many IDN rules are there in the world ?
- How do we find IDN Spoof under these rules?



# Continue Thinking

# Unicode Visual Security Issues



# Grapheme Rendering



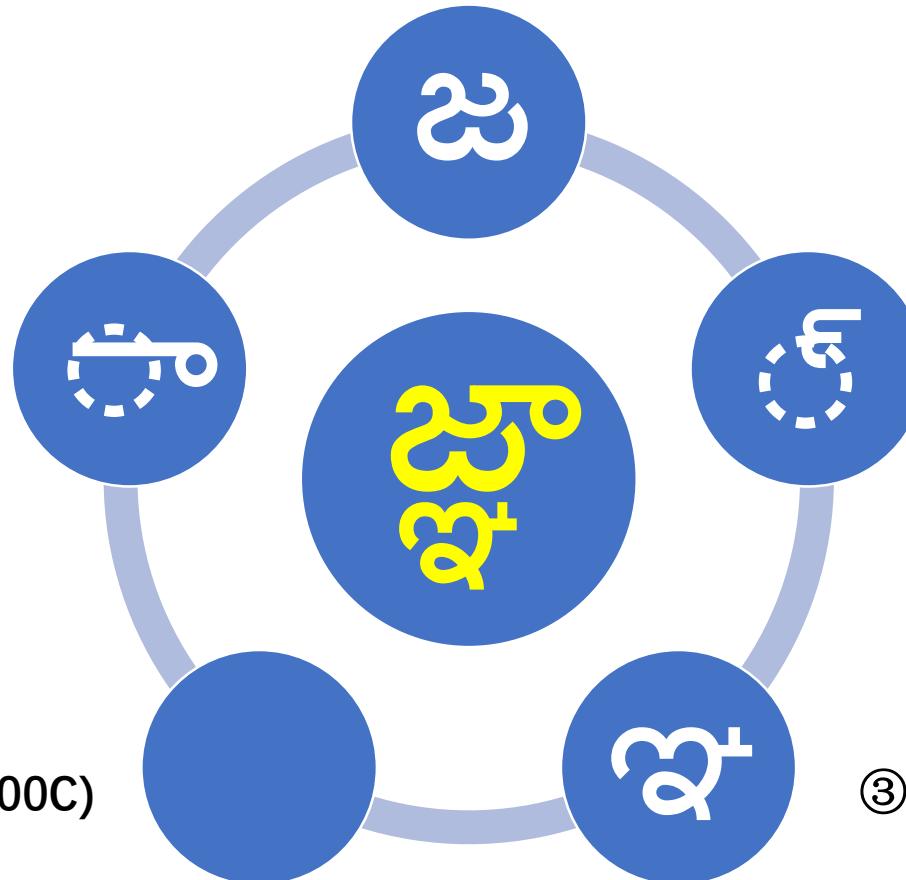
- Inadequate Rendering
  - An additional problem arises when a font or rendering engine has inadequate support for characters or sequences of characters.
- Quick Example

String	UTF-16	Punycode
el.com	<u>0065</u> 006C <u>0323</u> 002E 0063 006F 006D	xn--e-zom.com
ęl.com	<u>0065 0323</u> 006C 002E 0063 006F 006D	xn--l-ewm.com
ęl.com	<u>1EB9</u> 006C 002E 0063 006F 006D	xn--l-ewm.com

# CVE-2018-4124



① telugu sign virama (U+0C1C)



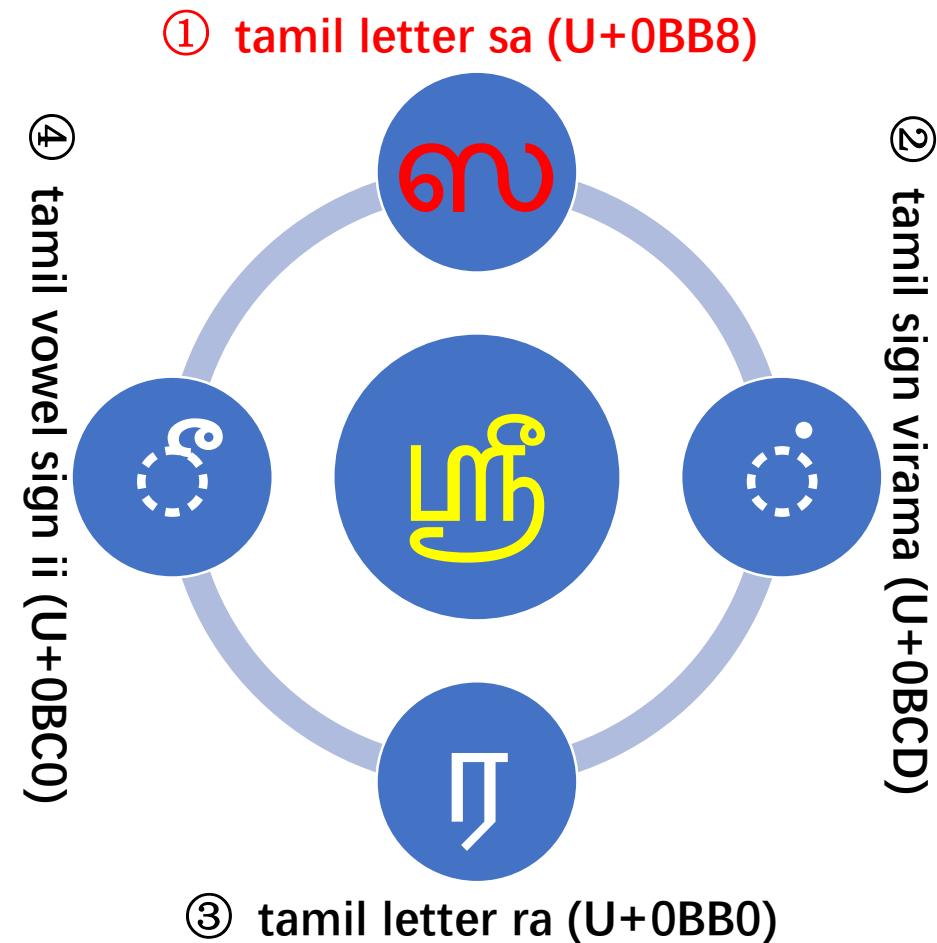
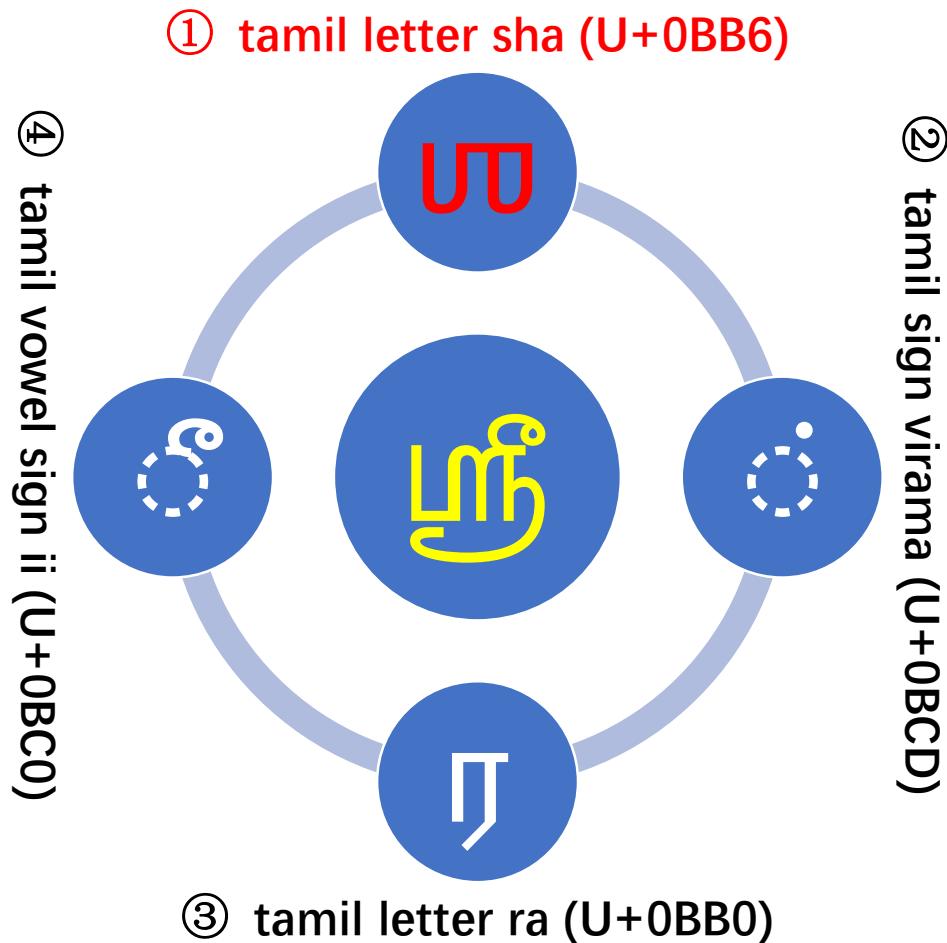
⑤ telugu vowel sign aa (U+0C3E)

② telugu sign virama (U+0C4D)

④ zero width non-joiner (U+200C)

③ telugu letter nya (U+0C1E)

# Glyphs in Complex Scripts





# Mixed-Script

- Mixed-Script Spoofing
  - The characters in some scripts, though different in meaning, are usually identical or nearly identical in appearance. However, the existence of visually confusable characters across scripts offers numerous opportunities for spoofing.
- Quick Example

String	UTF-16	Punycode
top.com	0074 <b>03BF</b> 0070 002E 0063 006F 006D	xn--tp-jbc.com
top.com	0074 <b>006F</b> 0070 002E 0063 006F 006D	top.com

# Punycode



- PunyCode Spoofing
  - Punycode is a special encoding used to convert Unicode characters to ASCII, which is used to encode internationalized domain names (IDN). The Punycode transformation is relatively dense. That means that it is fairly likely that arbitrary words after the "xn--" will result in valid labels.
- Quick Example
  - URL: <http://蘋鵝護頭.com>
  - PunyCode URL: <http://xn--google.com>

# Bidirectional Text



- Bidirectional Text Spoofing
  - When characters are mixed with left-to-right text(LTR) and right-to-left text(RTL), Unicode Bidirectional Algorithm will use a precise set of rules to determine the final visual rendering. However, presented with arbitrary sequences of text, this may lead to text sequences which may be impossible to read intelligibly, or which may be visually confusable.
- Quick Example
  - Access to <http://127.0.0.1/%D8%A7/example.org>
  - In address bar, maybe visual rendering <http://example.org/%127.0.0.1>

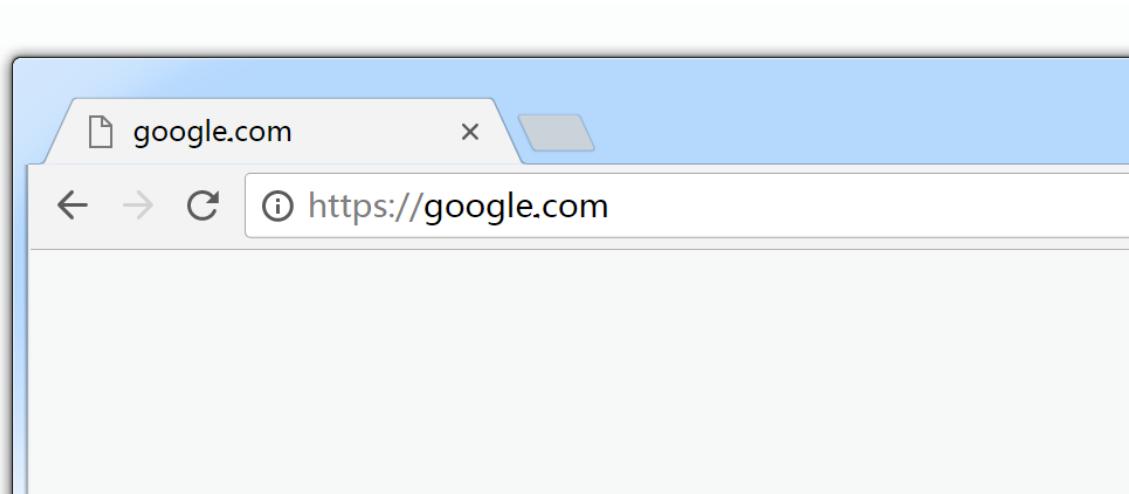
# LTR vs RTL

- URL-LTR
  - Subdomain: hi
  - Domain: google
  - TLD: com
  - Path: search
- http://hi.google.com/search
- URL-RTL(hebrew)
  - Subdomain: ה
  - Domain: ר
  - TLD: ל
  - Path: נ
- 渲染 : http://נ.ר.ל/ה



# Combining Mark

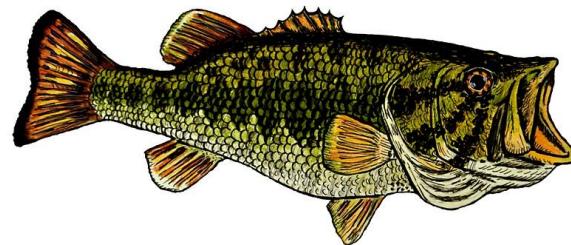
- Combing mark spoofing
  - Combining mark are characters that are intended to modify other characters. Combining character sequence maybe become visually indistinguishable with other characters .
- Quick Example
  - googleø.com
  - ø = U+18A9





**THINK  
OUTSIDE  
THE  
BOX**

# Inadequate Grapheme Rendering



**Safari Address Bar Spoof Using Latin-d  
CVE-2018-4277**

# Latin Extended-D

	A72	A73	A74	A75	A76	A77	A78	A79	A7A	A7B	A7C	A7D	A7E	A7F
0	Ł	F	K	P	W	ŋ	ł	N	G	K				
1	F	S	ќ	ප	ѡ	đ	ି	ନ	ଗ	ଳ				
2	ڙ	AA	K	P	ڙ	ڙ	ڙ	€	K	J				
3	ڙ	aa	k	p	ڙ	m	ନ	ୟ	k	X				
4	Ը	AO	K	Պ	Ծ	n	r	Ծ	Ն	B				
5	Ը	ao	ќ	պ	ծ	r	r	հ	ն	β				
6	H	AJ	L	Q	P	R	C	B	R	Ծ				
7	ହ	ାଇ	ି	କୁ	ପୁ	ତ୍ର	ତ୍ର	ବୁ	ର୍ମ	ସୁ				I
8	ବ୍ୟ	A'	L	କୁ	ି	ଧ୍ୟ	ା	ଫ୍ର	ସୁ					H
9	ତ୍ର୍ୟ	a'	ି	କୁ	ି	O	:	f	ସୁ					œ
A	Ը	A'	Թ	Ր	Յ	Ծ	=	Ը	Հ					՚
B	Ը	a'	Թ	Ր	Յ	Ծ	'	ա	Յ					՚
C	Գ	A'	Օ	Շ	Կ	Ւ	'	Ծ	Գ					՚
D	Գ	g	Օ	Շ	Կ	Ւ	Ծ	Ծ	Լ					M
E	Գ,	Ծ	Օ	Վ	Կ	Ղ	Ծ	Ծ	Լ					I
F	Գ,	Ծ	Օ	Վ	Կ	Ղ	•	Ծ	Լ					՚

d̄

latin small letter dum (U+A771)



d

In Safari address bar



# Latin: icloud.com VS Latin Extended-D: icloud.com



icloud

latin small letter i (U+0069)  
latin small letter c (U+0063)  
latin small letter l (U+006C)  
latin small letter d (U+0064)  
latin small letter u (U+0075)  
latin small letter o (U+006F)



icloud.com



icloudꝝ

latin small letter i (U+0069)  
latin small letter c (U+0063)  
latin small letter l (U+006C)  
**latin small letter dum (U+A771)**  
latin small letter u (U+0075)  
latin small letter o (U+006F)



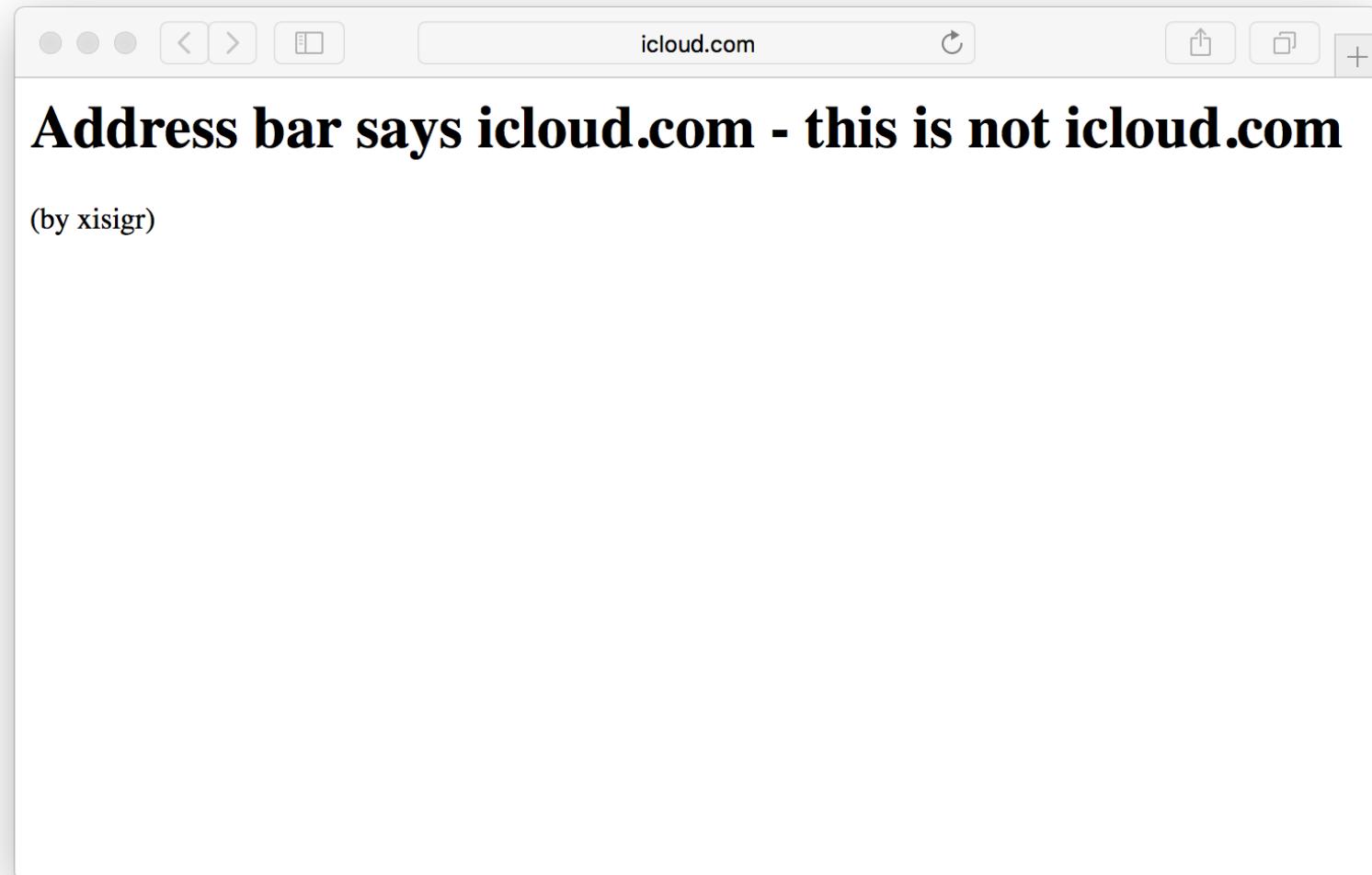
xn--icloud-rl3s.com



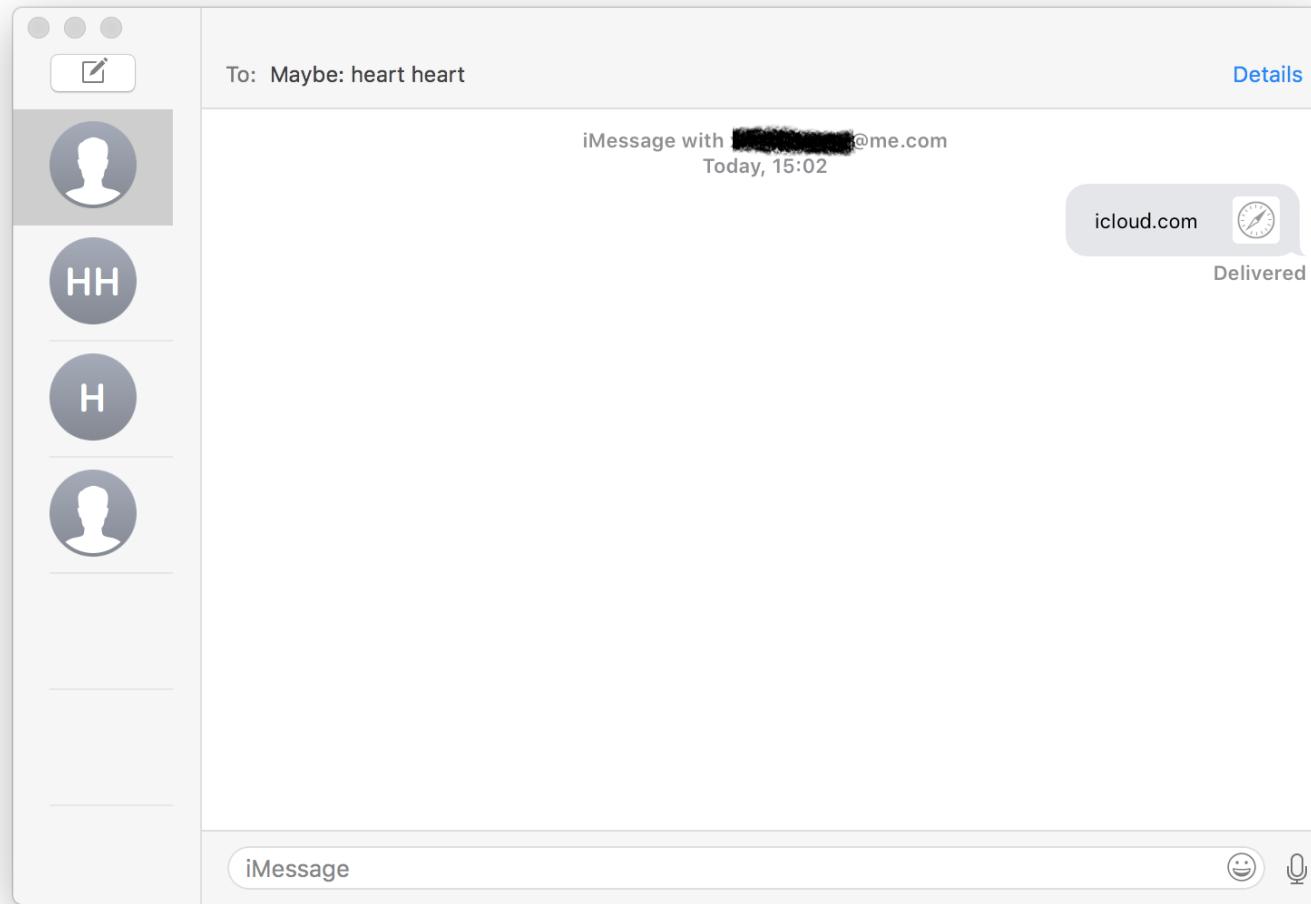
icloud.com



# xn--icloud-rl3s.com



# xn--icloud-rl3s.com





# Domain included can be spoof

- Top 10k , domain included > 25%
  - linkedin.com
  - baidu.com
  - jd.com
  - adobe.com
  - wordpress.com
  - dropbox.com
  - godaddy.com
  - reddit.com
  - .....

# Whole-Script



**Address Bar Spoof using Cyrillic  
CVE-2017-5060**

# Cyrillic

	040	041	042	043	044	045	046	047	048	049	04A	04B	04C	04D	04E	04F
0	Ѐ	Ӑ	Ӗ	Ҫ	Ӯ	Ӯ	Ӯ	Ӯ	Ӯ	Ӯ	Ӯ	Ӯ	Ӯ	Ӯ	Ӯ	Ӯ
1	Ӷ	ӷ	Ӹ	ӹ	ӻ	ӻ	ӻ	ӻ	ӻ	ӻ	ӻ	ӻ	ӻ	ӻ	ӻ	ӻ
2	Ӱ	ӱ	Ӳ	ӳ	Ӵ	ӵ	ӵ	ӵ	ӵ	ӵ	ӵ	ӵ	ӵ	ӵ	ӵ	ӵ
3	Ӱ	ӱ	Ӳ	ӳ	Ӵ	ӵ	ӵ	ӵ	ӵ	ӵ	ӵ	ӵ	ӵ	ӵ	ӵ	ӵ
4	Ӱ	ӱ	Ӳ	ӳ	Ӵ	ӵ	ӵ	ӵ	ӵ	ӵ	ӵ	ӵ	ӵ	ӵ	ӵ	ӵ
5	Ӱ	ӱ	Ӳ	ӳ	Ӵ	ӵ	ӵ	ӵ	ӵ	ӵ	ӵ	ӵ	ӵ	ӵ	ӵ	ӵ
6	Ӱ	ӱ	Ӳ	ӳ	Ӵ	ӵ	ӵ	ӵ	ӵ	ӵ	ӵ	ӵ	ӵ	ӵ	ӵ	ӵ
7	Ӱ	ӱ	Ӳ	ӳ	Ӵ	ӵ	ӵ	ӵ	ӵ	ӵ	ӵ	ӵ	ӵ	ӵ	ӵ	ӵ
8	Ӱ	ӱ	Ӳ	ӳ	Ӵ	ӵ	ӵ	ӵ	ӵ	ӵ	ӵ	ӵ	ӵ	ӵ	ӵ	ӵ
9	Ӱ	ӱ	Ӳ	ӳ	Ӵ	ӵ	ӵ	ӵ	ӵ	ӵ	ӵ	ӵ	ӵ	ӵ	ӵ	ӵ
A	Ӱ	ӱ	Ӳ	ӳ	Ӵ	ӵ	ӵ	ӵ	ӵ	ӵ	ӵ	ӵ	ӵ	ӵ	ӵ	ӵ
B	Ӱ	ӱ	Ӳ	ӳ	Ӵ	ӵ	ӵ	ӵ	ӵ	ӵ	ӵ	ӵ	ӵ	ӵ	ӵ	ӵ
C	Ӱ	ӱ	Ӳ	ӳ	Ӵ	ӵ	ӵ	ӵ	ӵ	ӵ	ӵ	ӵ	ӵ	ӵ	ӵ	ӵ
D	Ӱ	ӱ	Ӳ	ӳ	Ӵ	ӵ	ӵ	ӵ	ӵ	ӵ	ӵ	ӵ	ӵ	ӵ	ӵ	ӵ
E	Ӱ	ӱ	Ӳ	ӳ	Ӵ	ӵ	ӵ	ӵ	ӵ	ӵ	ӵ	ӵ	ӵ	ӵ	ӵ	ӵ
F	Ӱ	ӱ	Ӳ	ӳ	Ӵ	ӵ	ӵ	ӵ	ӵ	ӵ	ӵ	ӵ	ӵ	ӵ	ӵ	ӵ

a

cyrillic small letter a (U+0430)

p

cyrillic small letter er (U+0440)

1

cyrillic small letter palochka (U+04CF)

e

cyrillic small letter ie (U+0435)

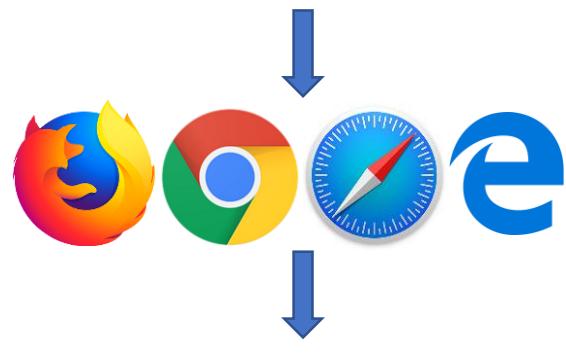


# Latin: apple.com VS Cyrillic : apple.com



apple

latin small letter a (U+0061)  
latin small letter p (U+0070)  
latin small letter e (U+0065)  
latin small letter l (U+006C)



✓



apple

cyrillic small letter a (U+0430)  
cyrillic small letter er (U+0440)  
cyrillic small letter ie (U+0435)  
cyrillic small letter palochka (U+04CF)

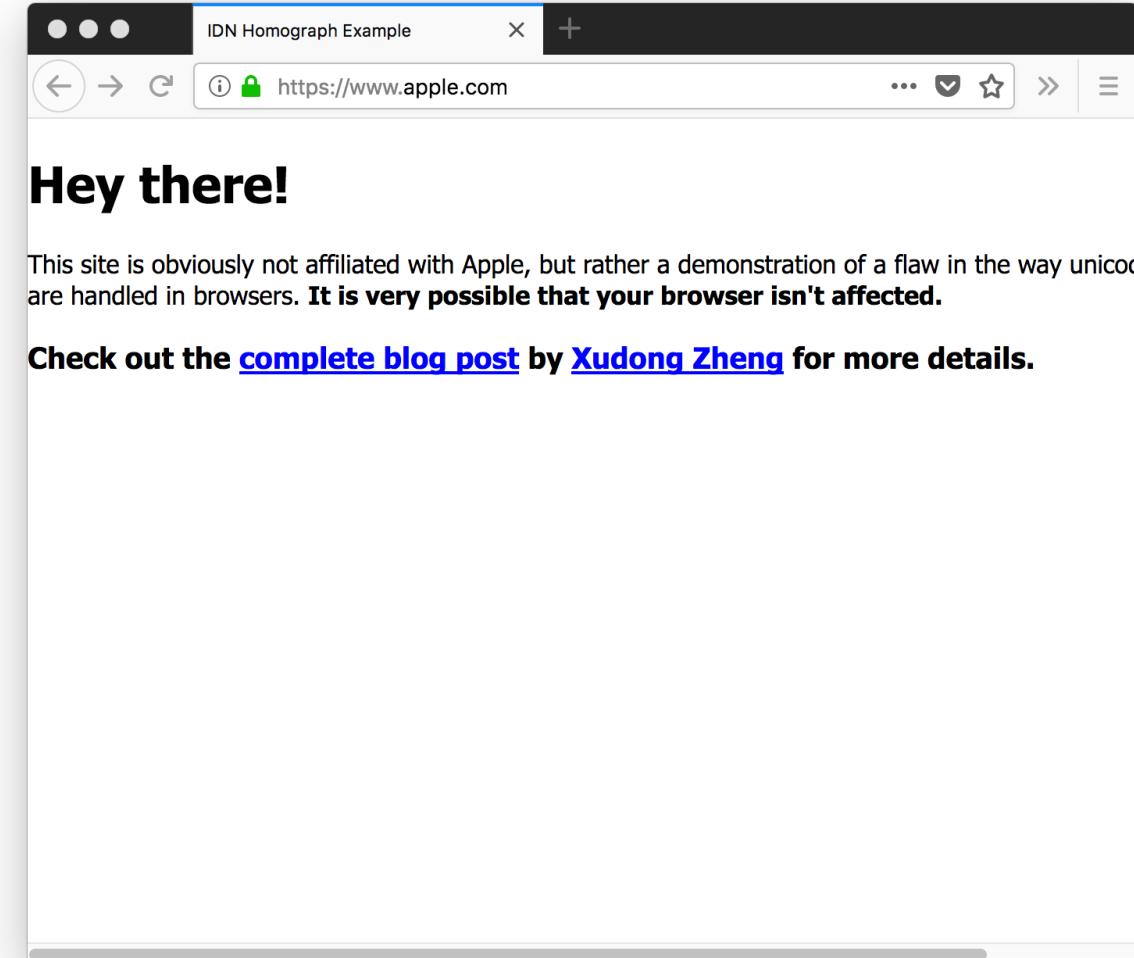
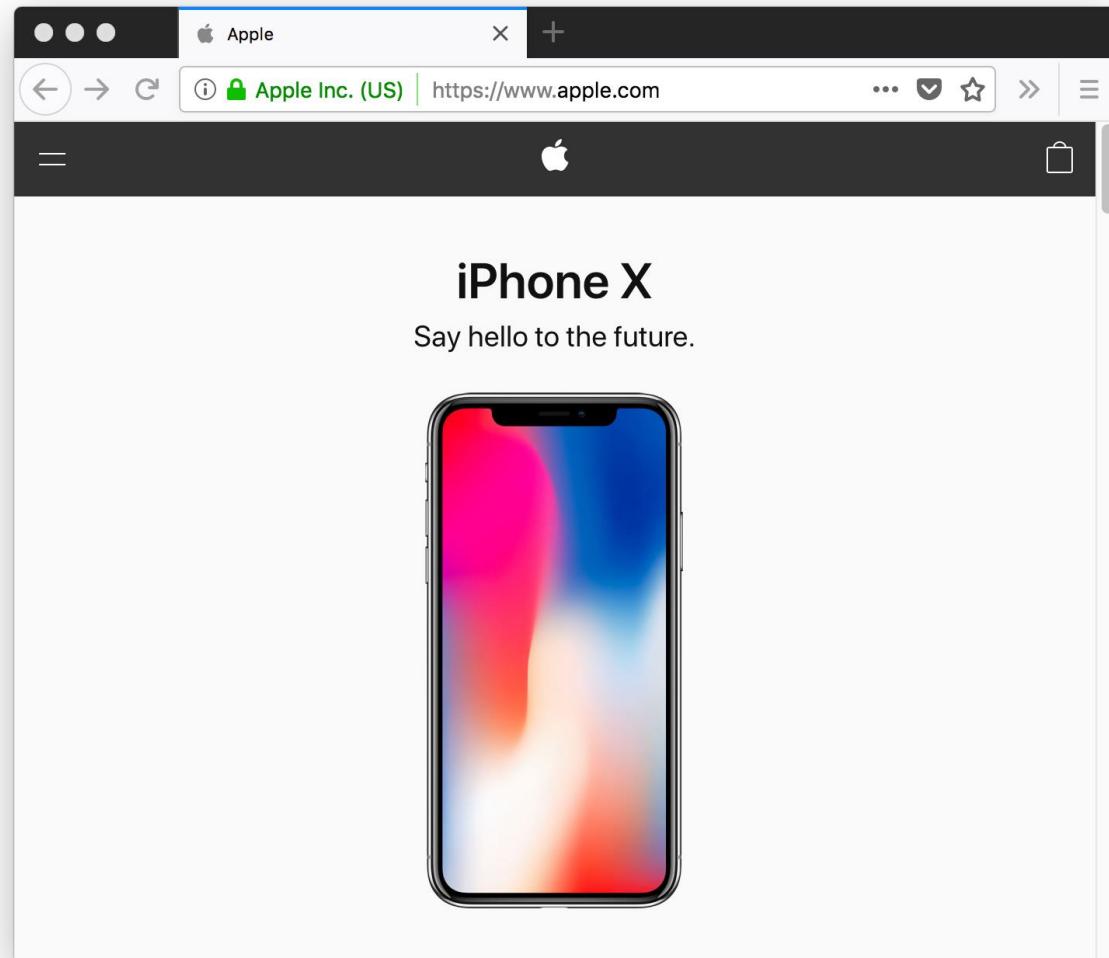


✓



✓

# Latin: apple.com VS Cyrillic : apple.com



# Bidirectional Text



**Firefox URL spoof using RTL  
CVE-2018-5117**

# CVE-2018-5117



TLD : شبكة

punycode=xn--ngbc5azd

Direction: RTL

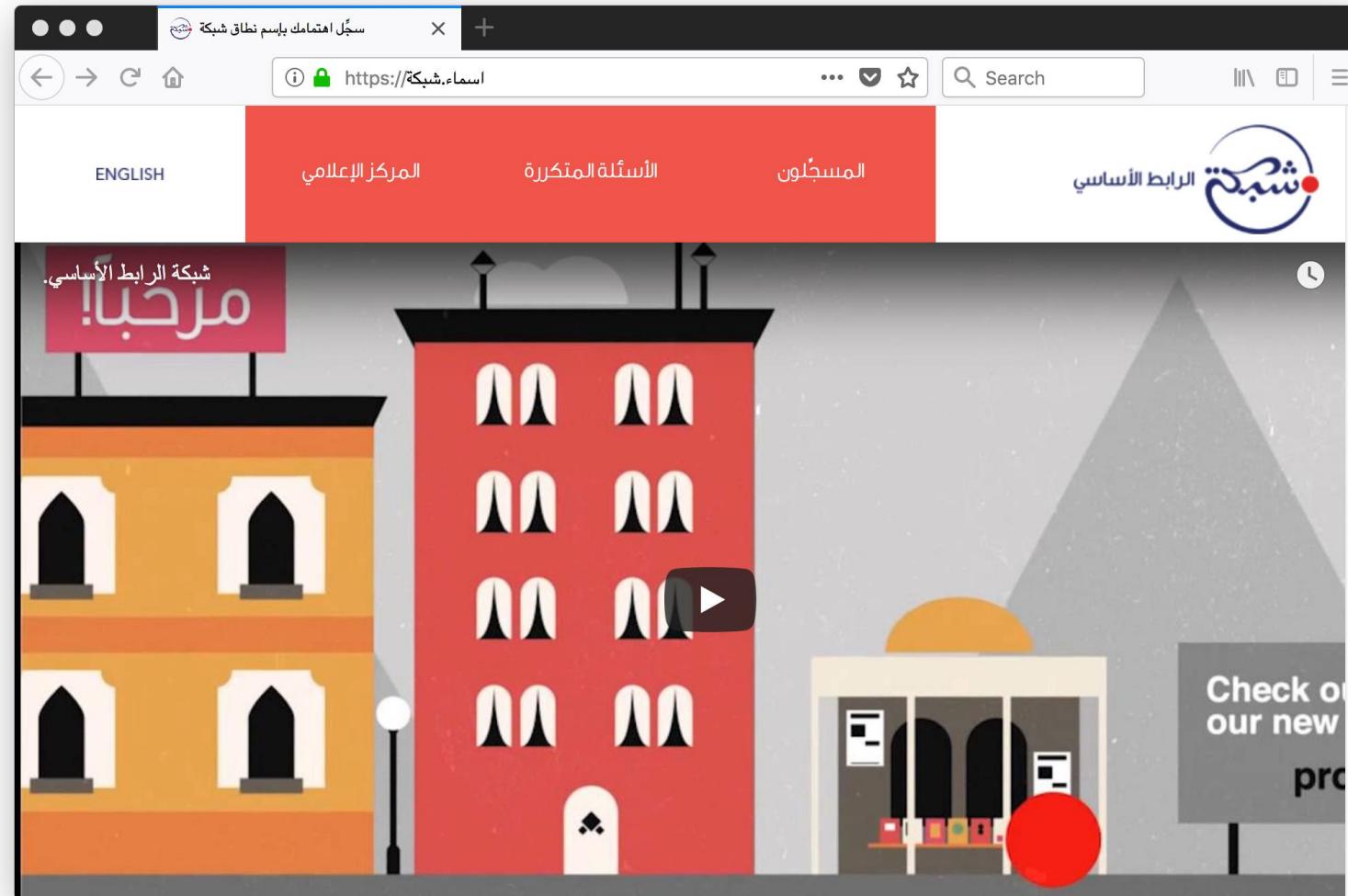
`https://xn--gbla1c4e.xn--ngbc5azd/`

SLD: اسماء

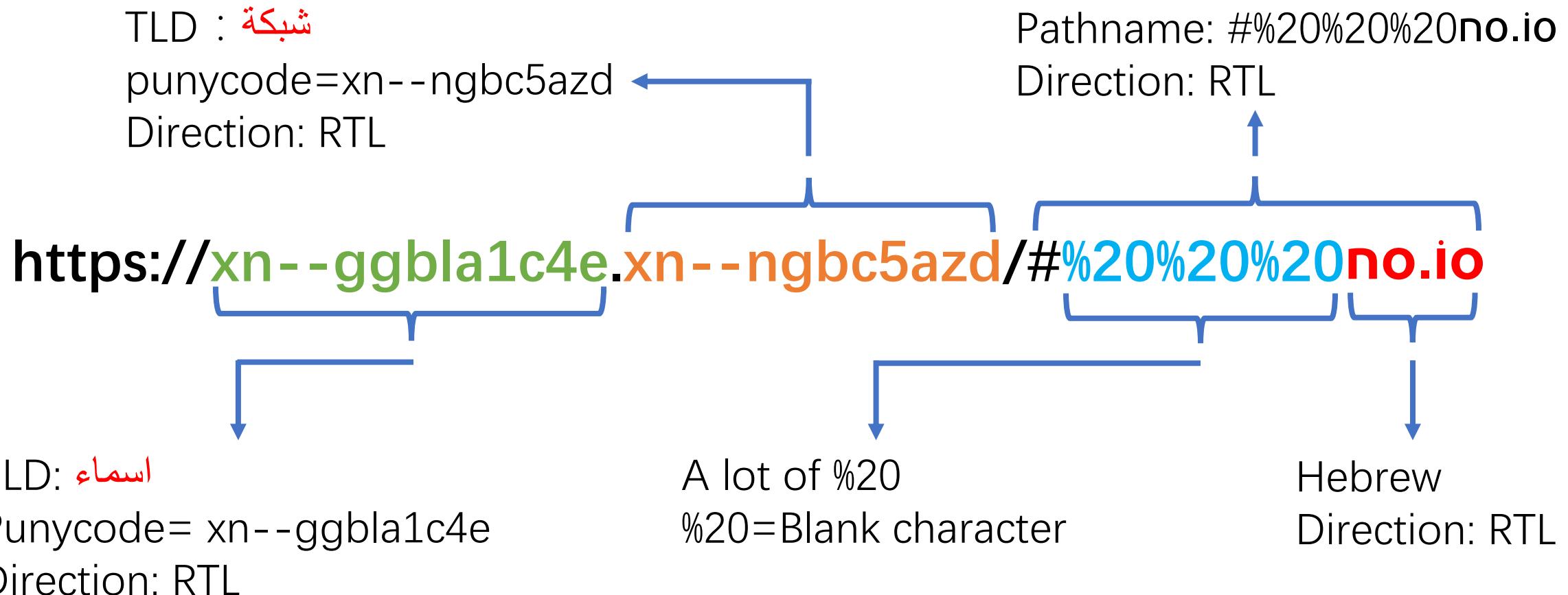
punycode= xn--gbla1c4e

Direction: RTL

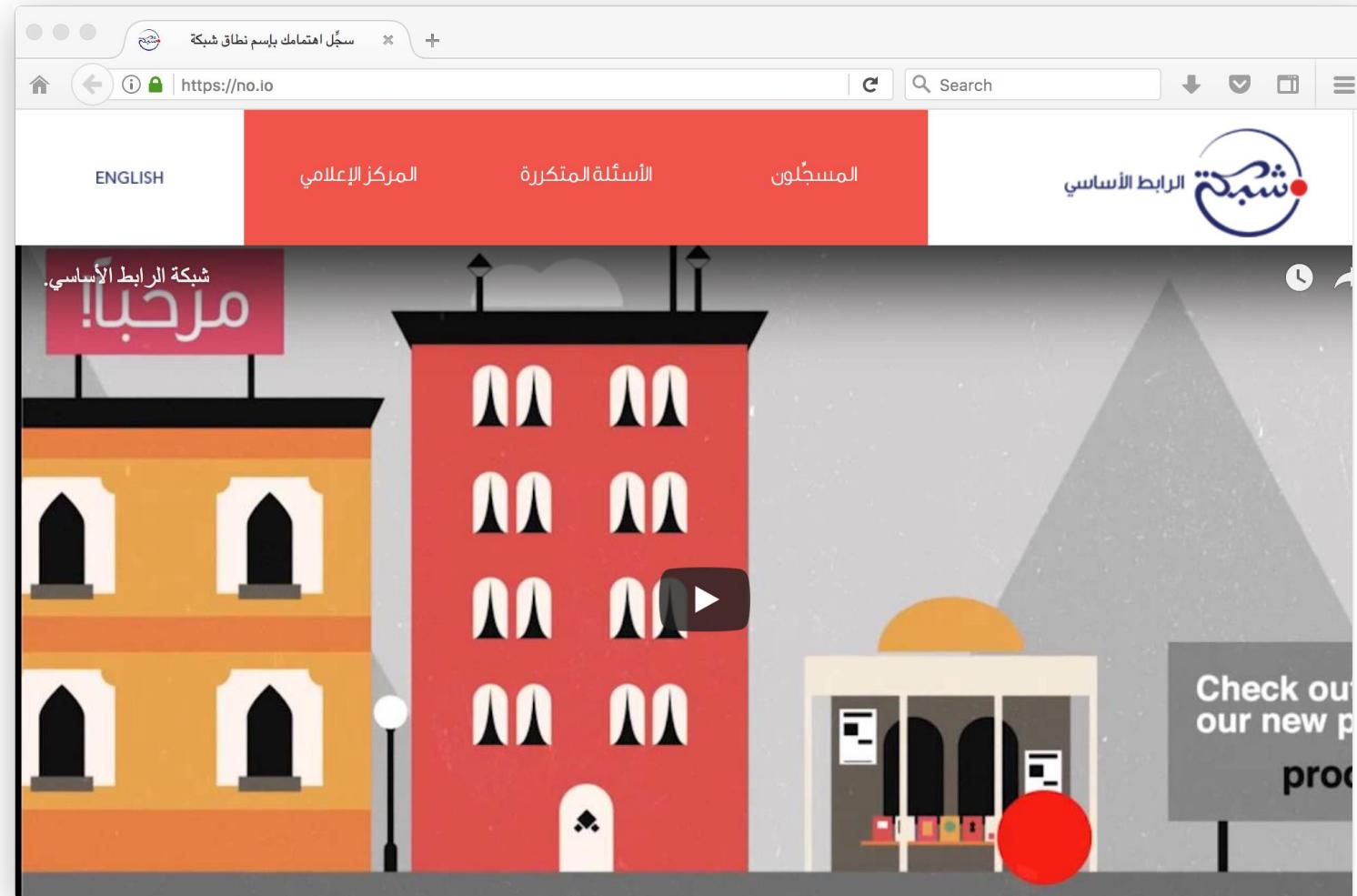
<https://xn--ggbla1c4e.xn--ngbc5azd>



# CVE-2018-5117



<https://xn--ggbla1c4e.xn--ngbc5azd/#%20%20%20ס.ה>



# Bidirectional Text



**URL Spoof Using RTL IDN TLD  
CVE-2018-4205**

<https://support.apple.com/en-us/HT208854>

# POC-1



TLD : شبكة

punycode=xn--ngbc5azd

Direction: RTL

**http://www.apple.com.xn--ggbla3j.xn--ngbc5azd/**

(345)-LD: www.apple.com

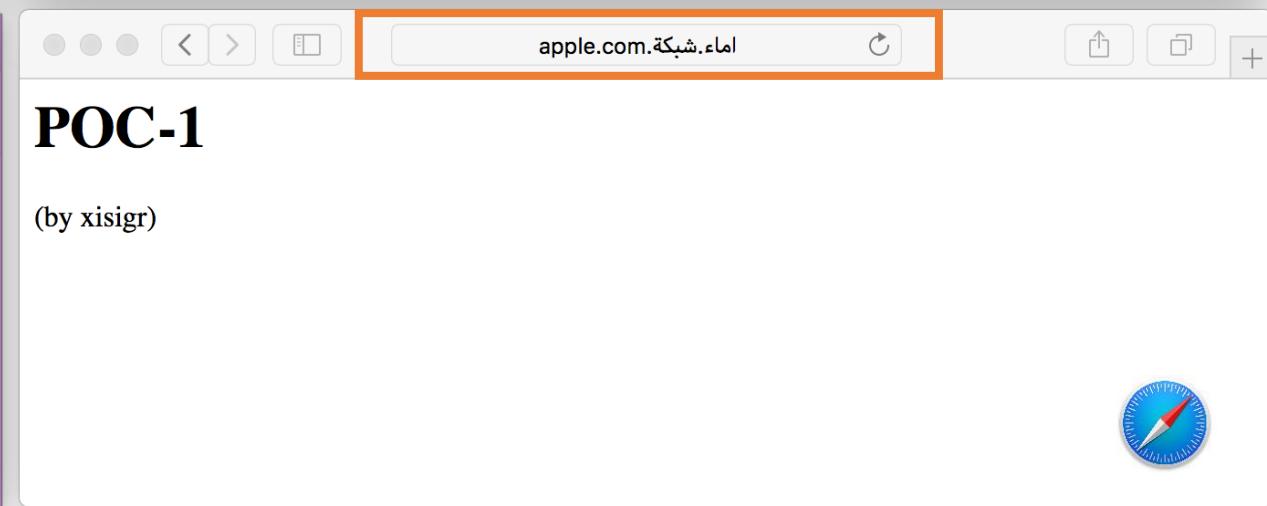
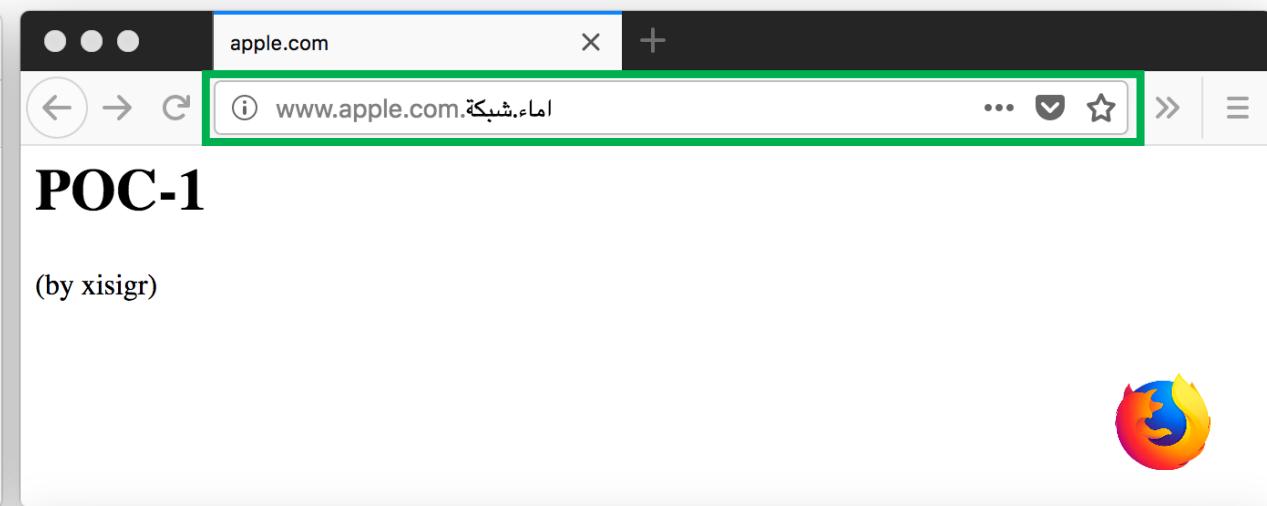
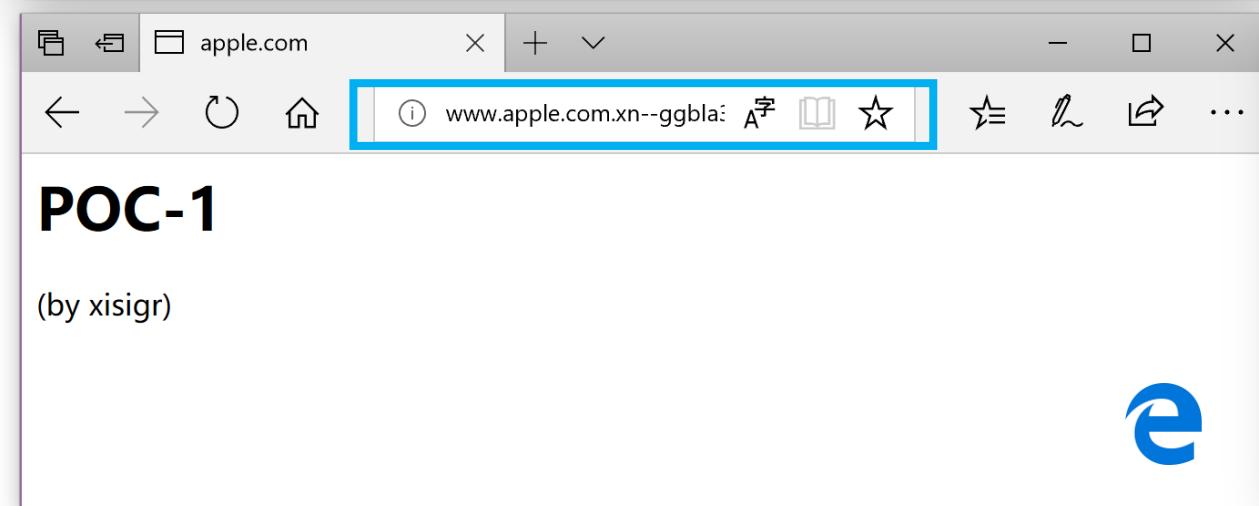
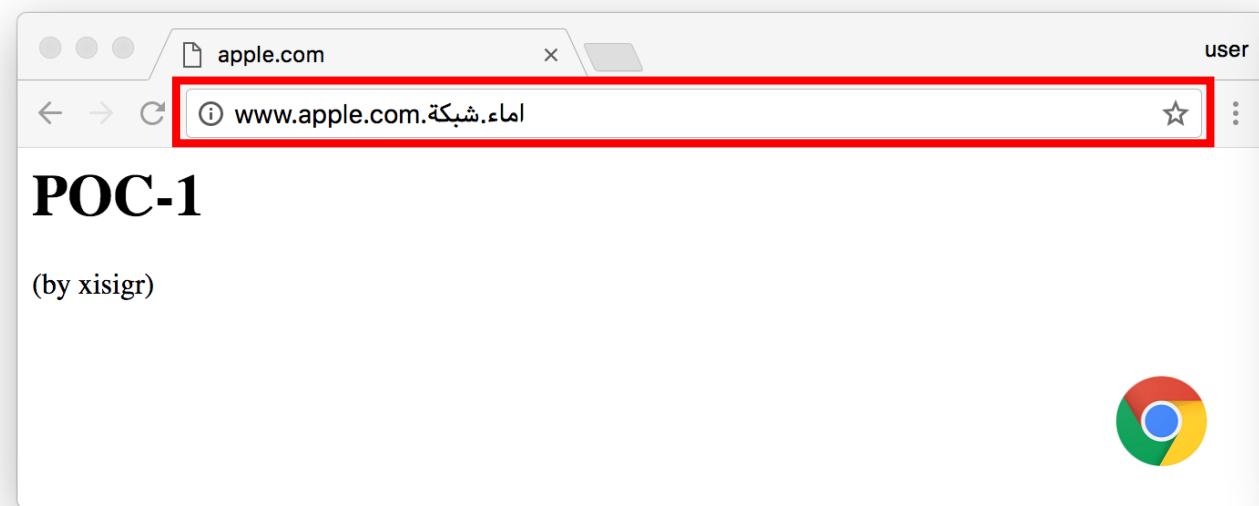
Direction: LTR

SLD : اماء

Punycode: xn--ggbla3j

Direction: RTL

# POC-1



# POC-2



TLD : شبكة

punycode=xn--ngbc5azd  
Direction: RTL

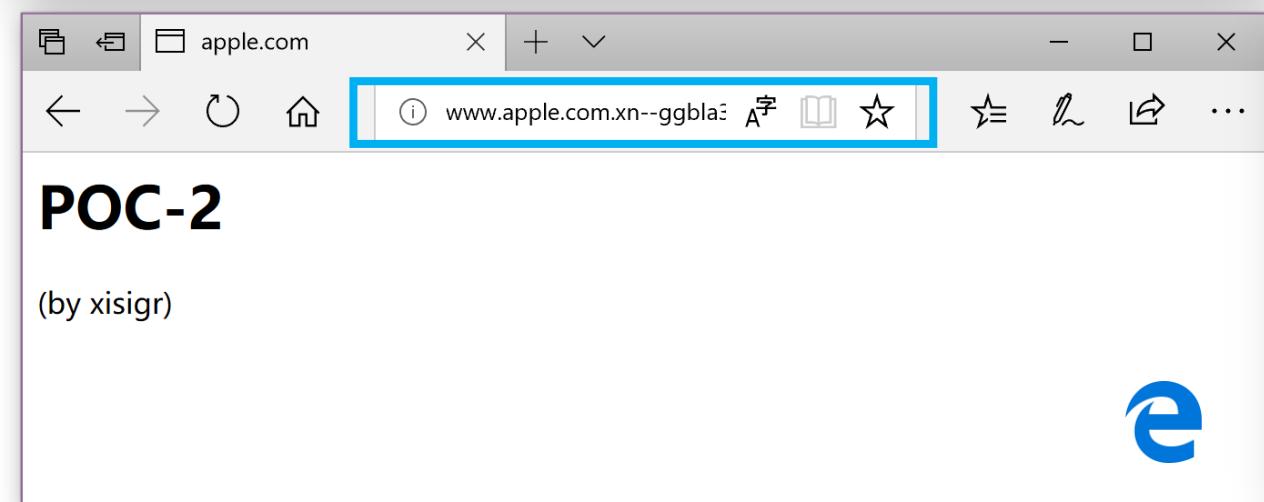
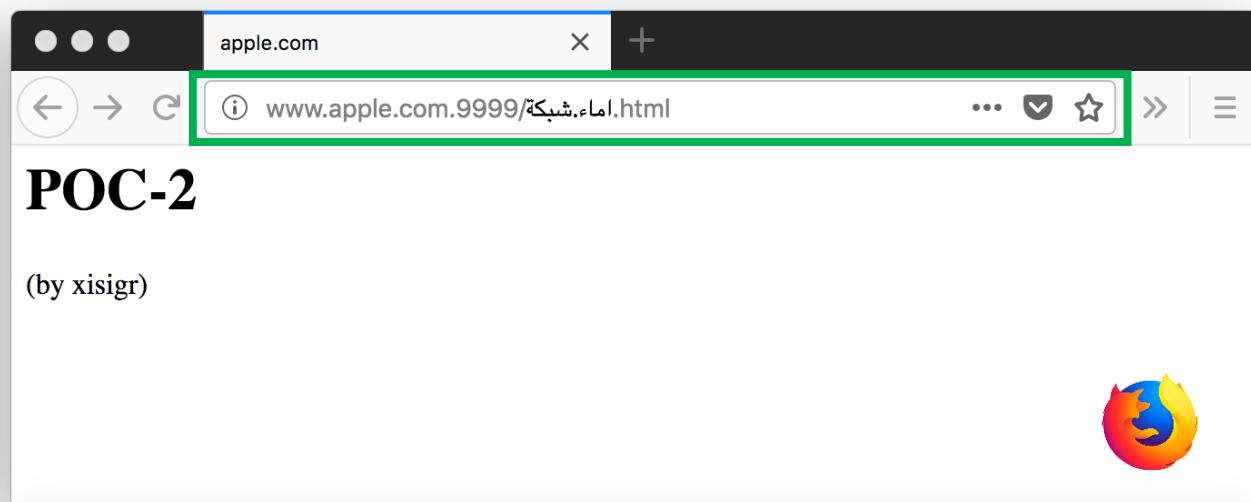
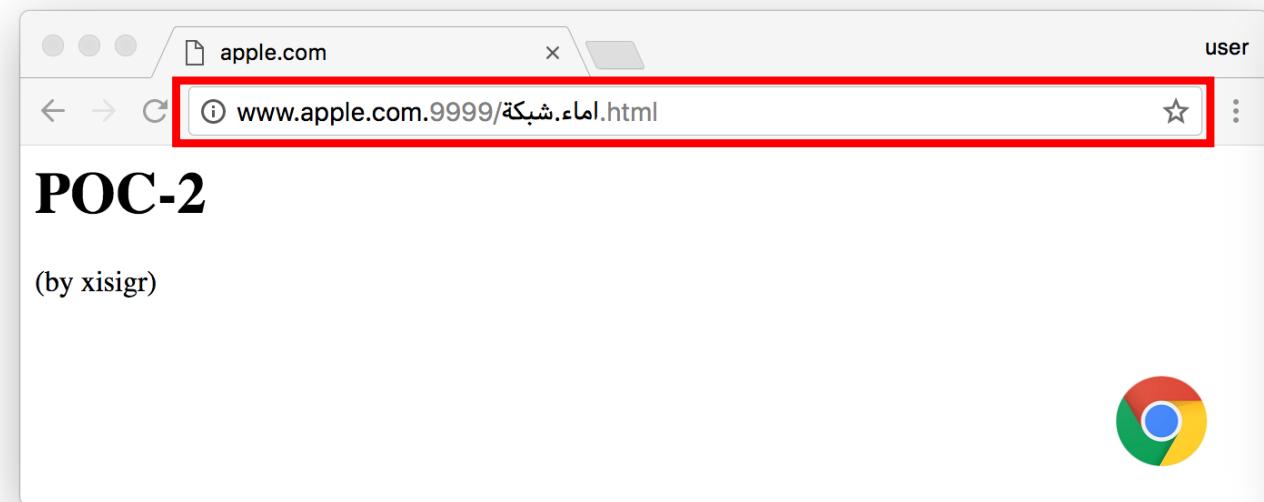
<http://www.apple.com.xn--gbla3j.xn--ngbc5azd/9999.html>

(345)-LD: www.apple.com  
Direction: LTR

SLD : اماء  
Punycode: xn--gbla3j  
Direction: RTL

Pathname: 9999.html  
Direction: LTR

# POC-2



# POC-3



## TLD : شبكة

punycode=xn--ngbc5azd

# Direction: RTL

<http://www.apple.com.xn--ggb1a3j.xn--ngbc5azd/999...html>

(345)-LD: www.apple.com  
Direction: LTR

SLD : اماء  
Punycode: xn--ggbla3j  
Direction: RTL

**POC-3**

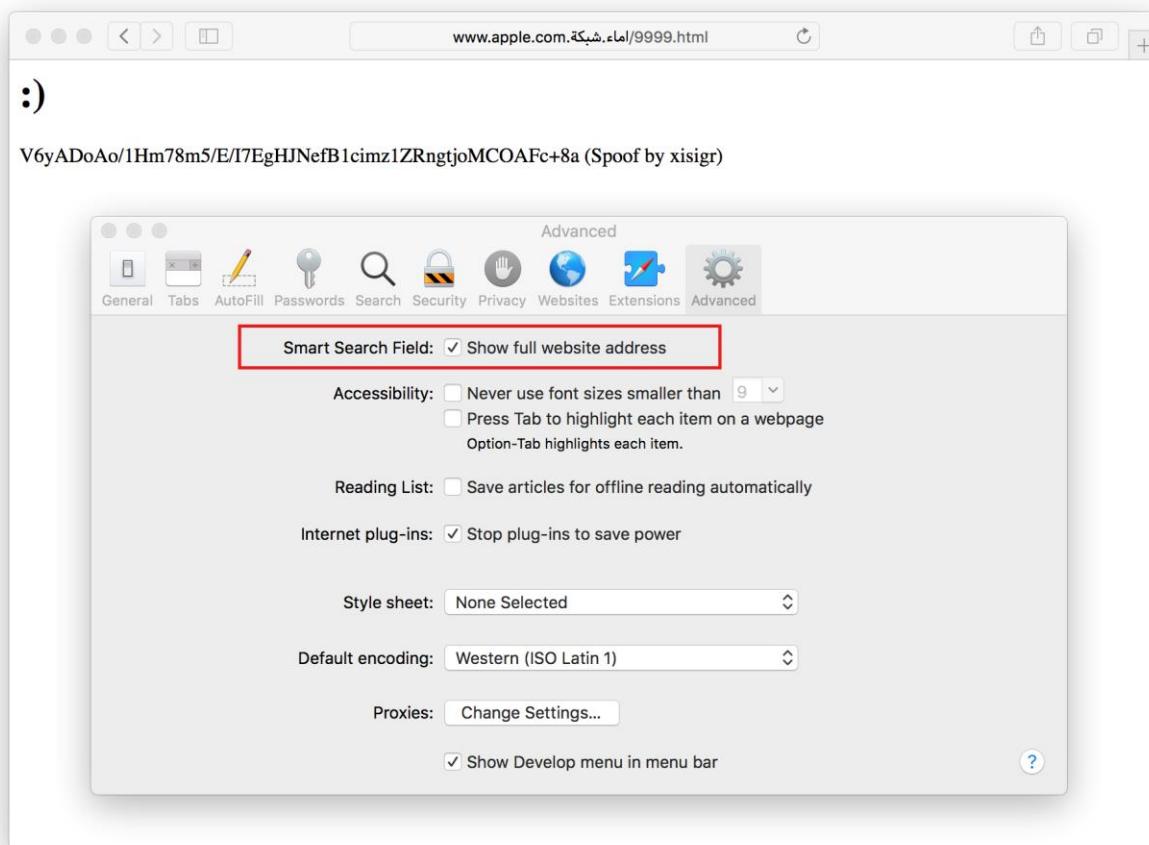


A screenshot of a web browser window. The address bar shows the URL "www.apple.com.xn--ggblaž" with a blue rectangular highlight around it. To the left of the URL is an info icon (a circle with an 'i'). To the right are icons for a magnifying glass, a star, and a document. The browser interface includes standard controls like back, forward, and search, as well as a title bar with the domain name "apple.com".

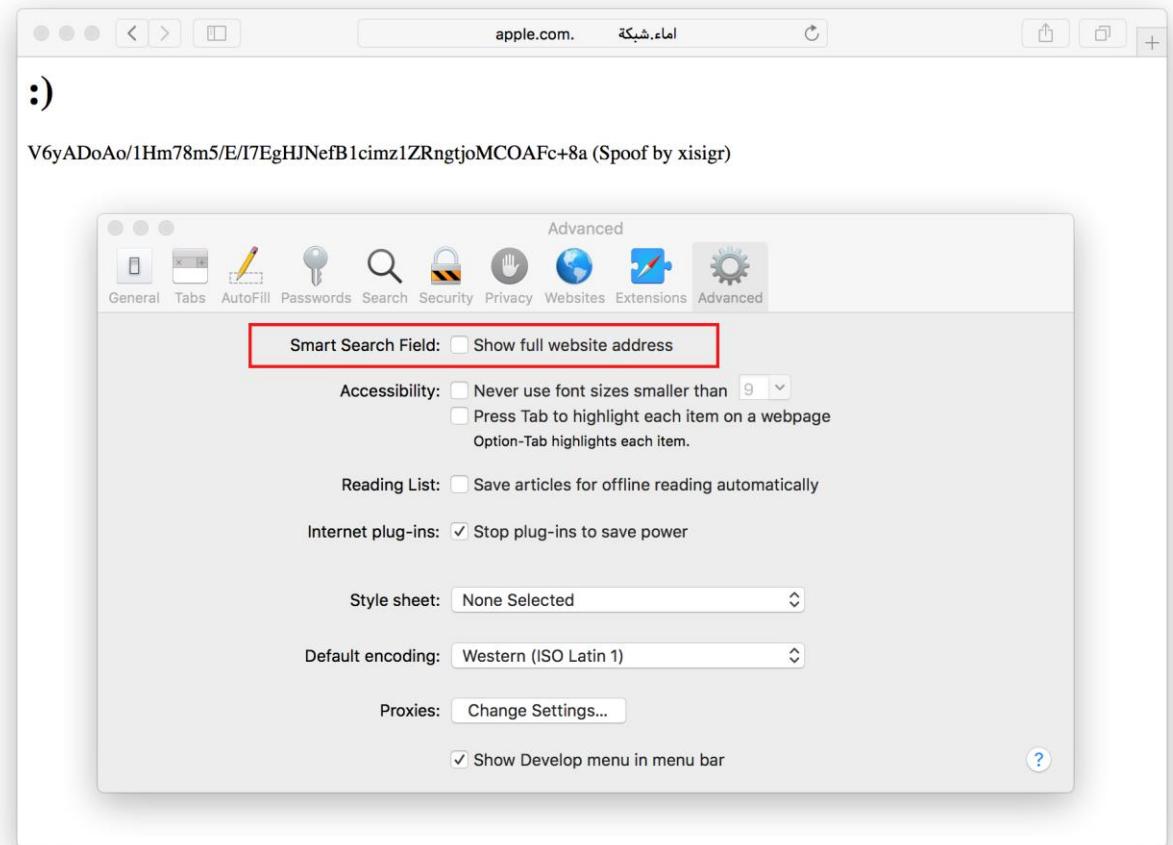
# Safari show website address



show full website address



only show domain

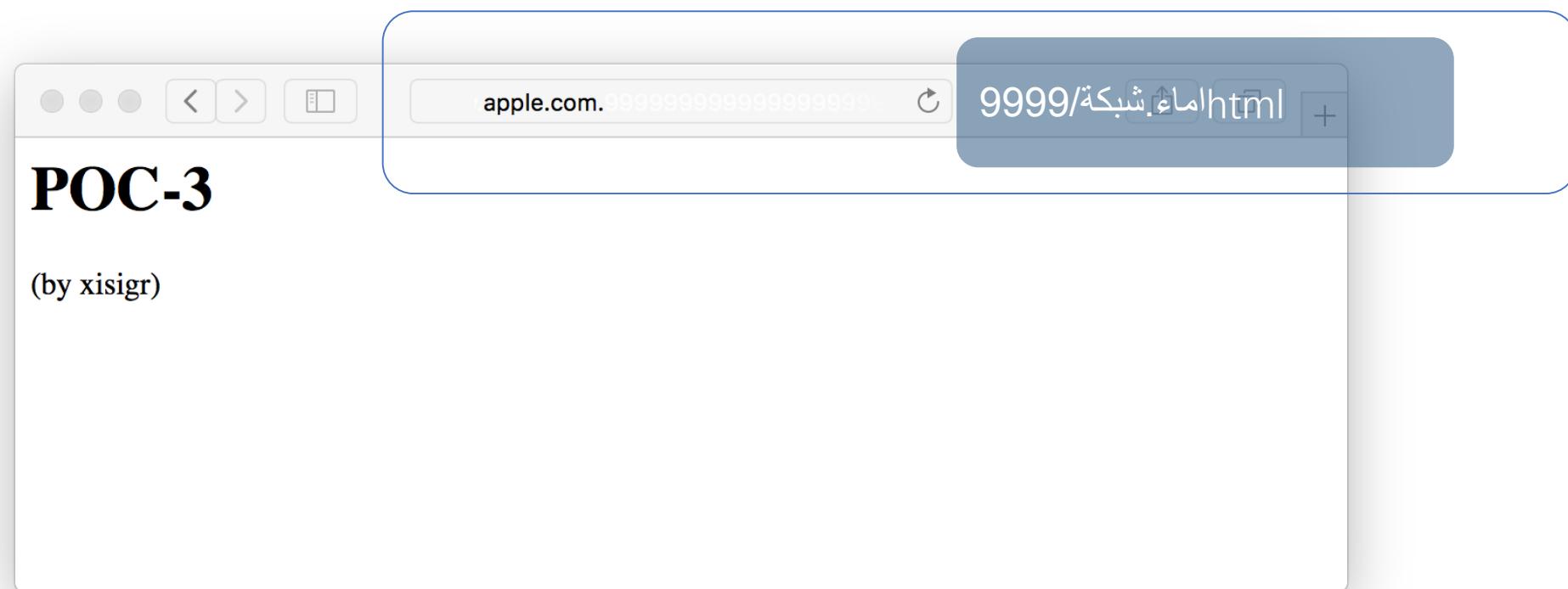


# only show domain

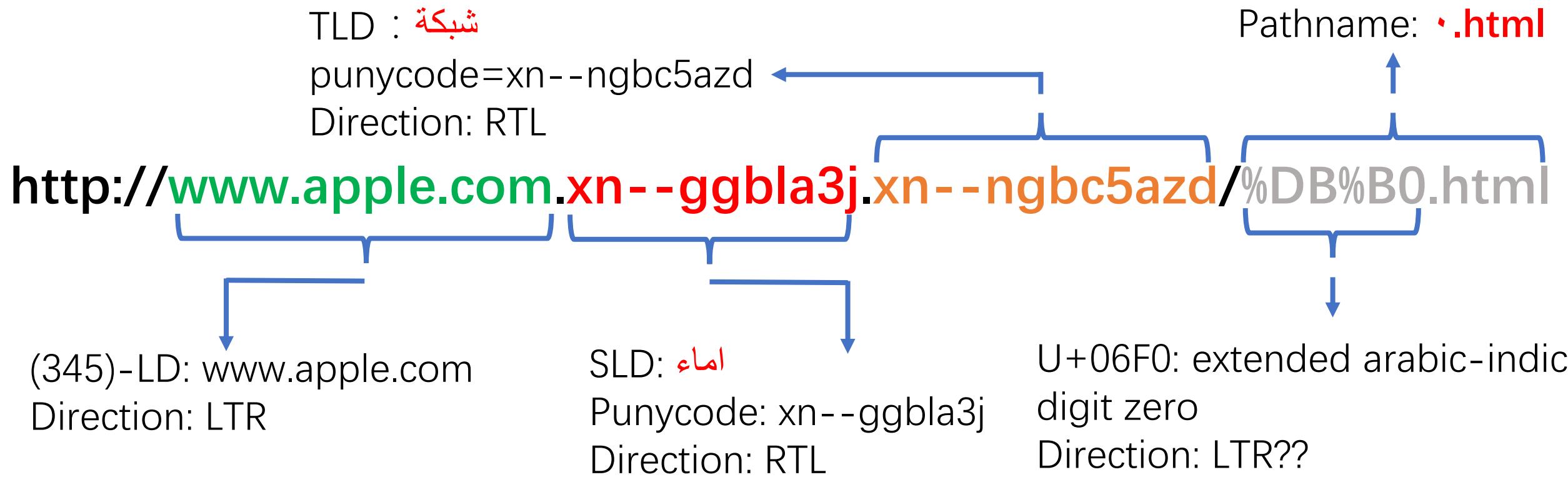


<http://www.apple.com.xn--gbla3j.xn--ngbc5azd/999...html>

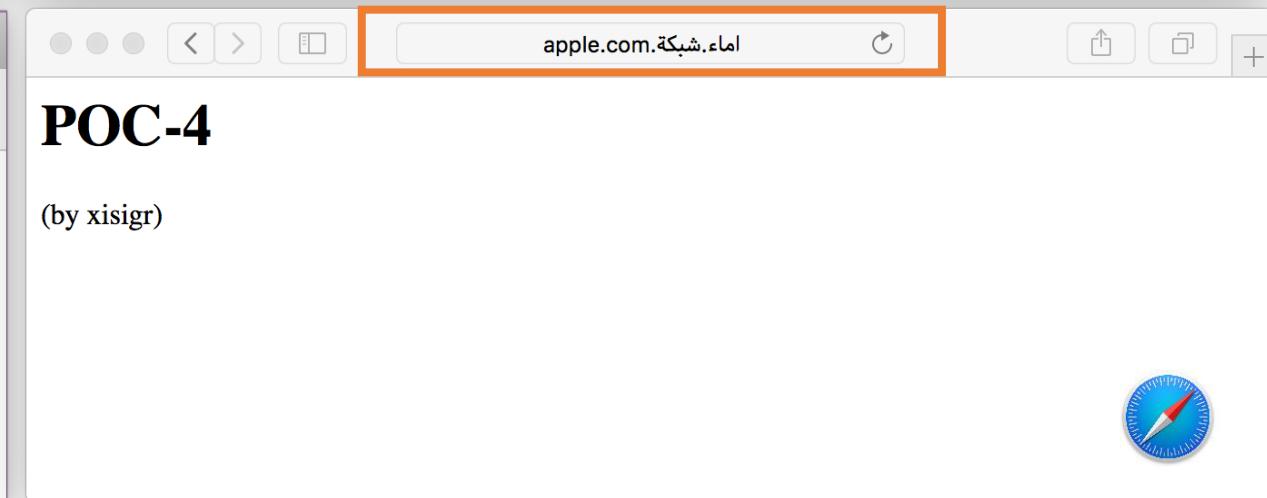
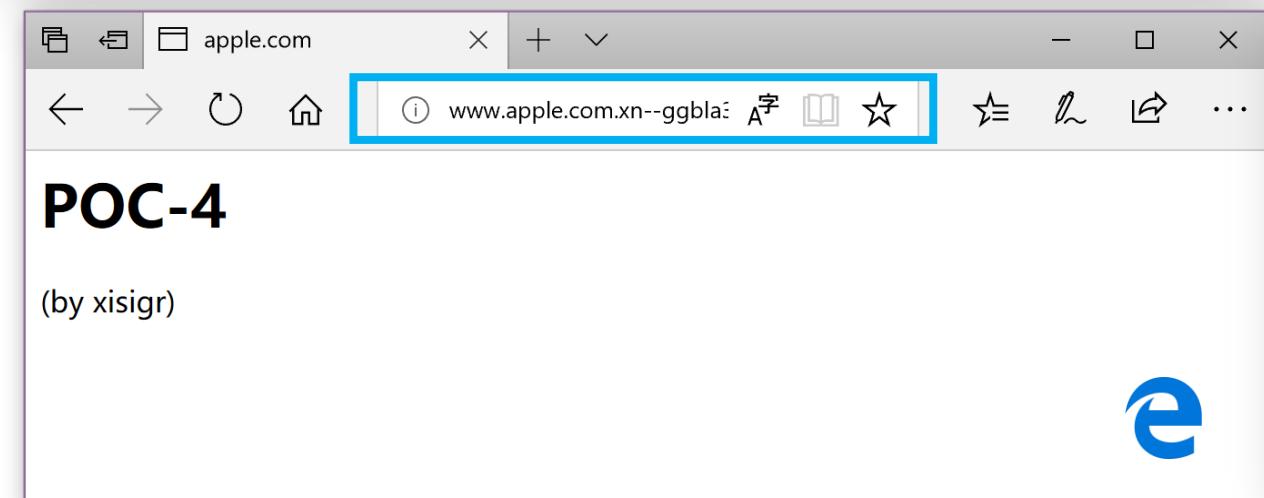
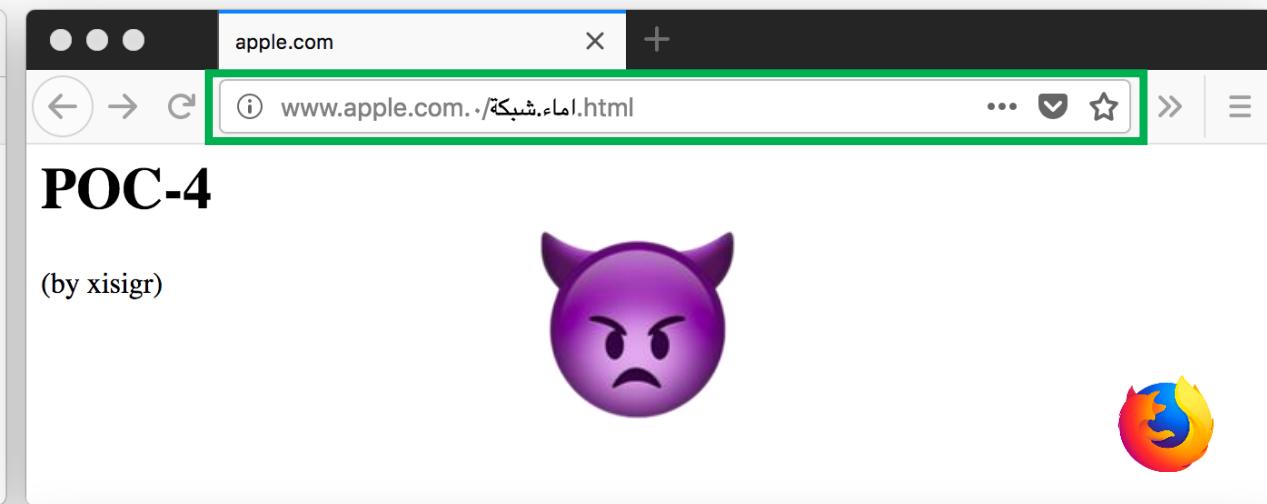
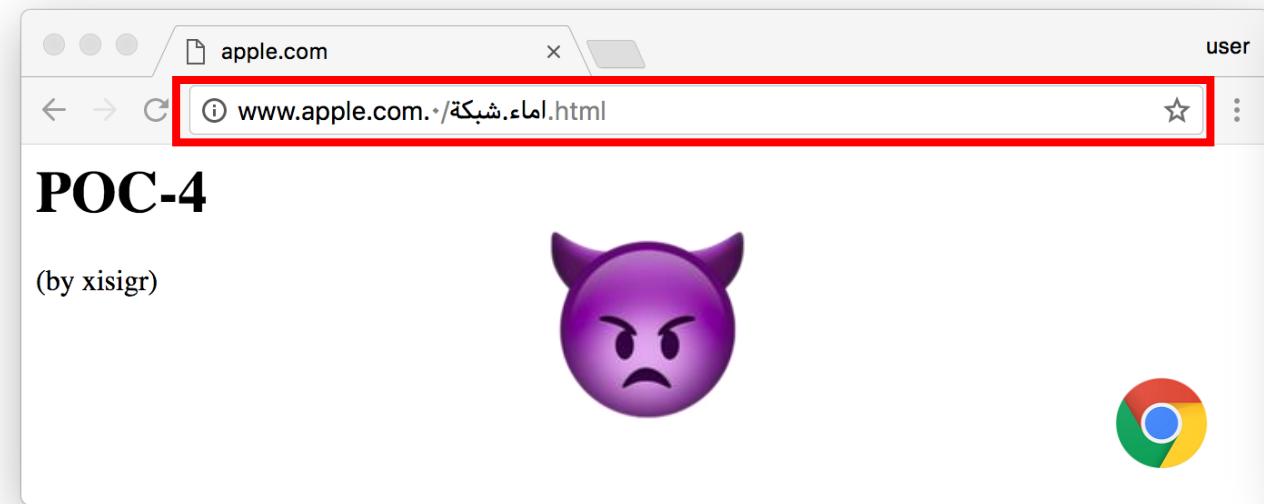
<http://www.apple.com.9999/اماء.شبكة.html>



# POC-4



# POC-4

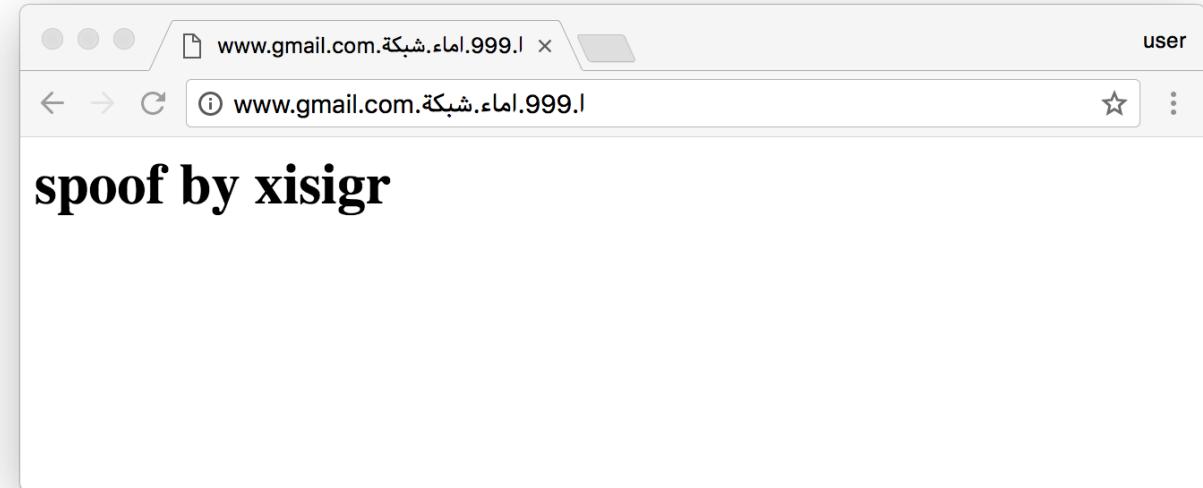
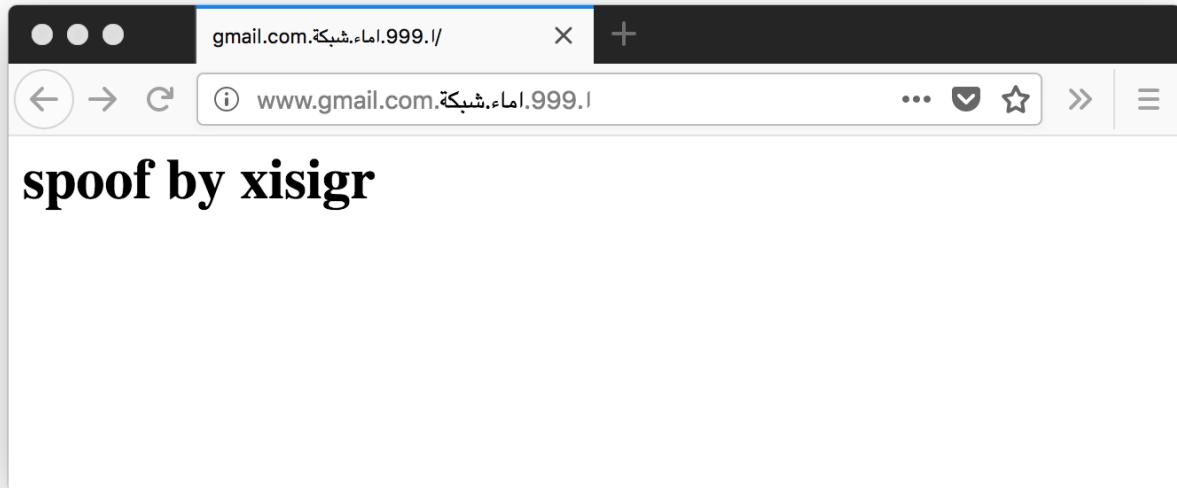


# Origin out-of-order

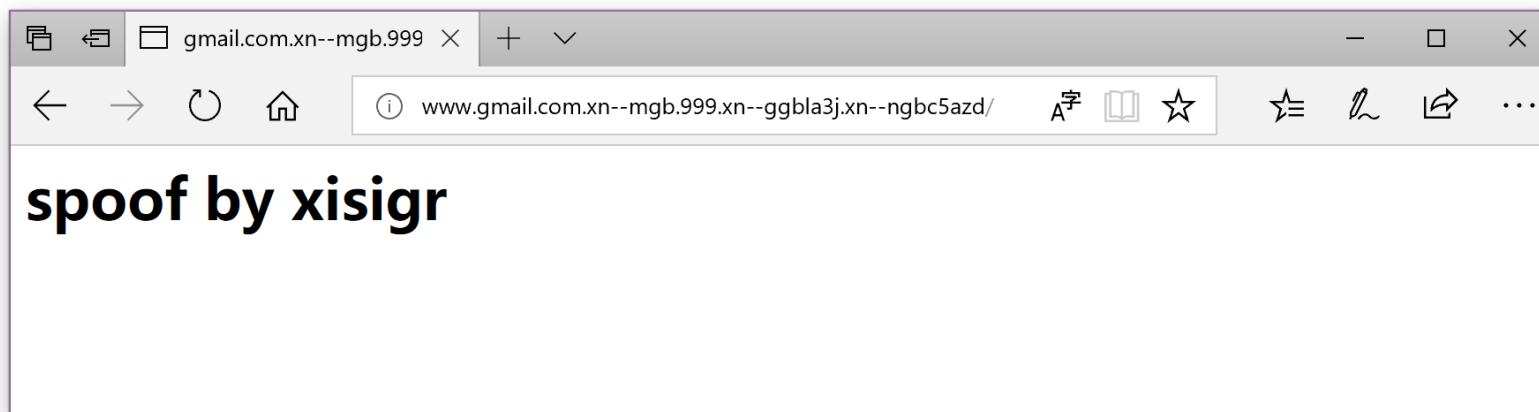
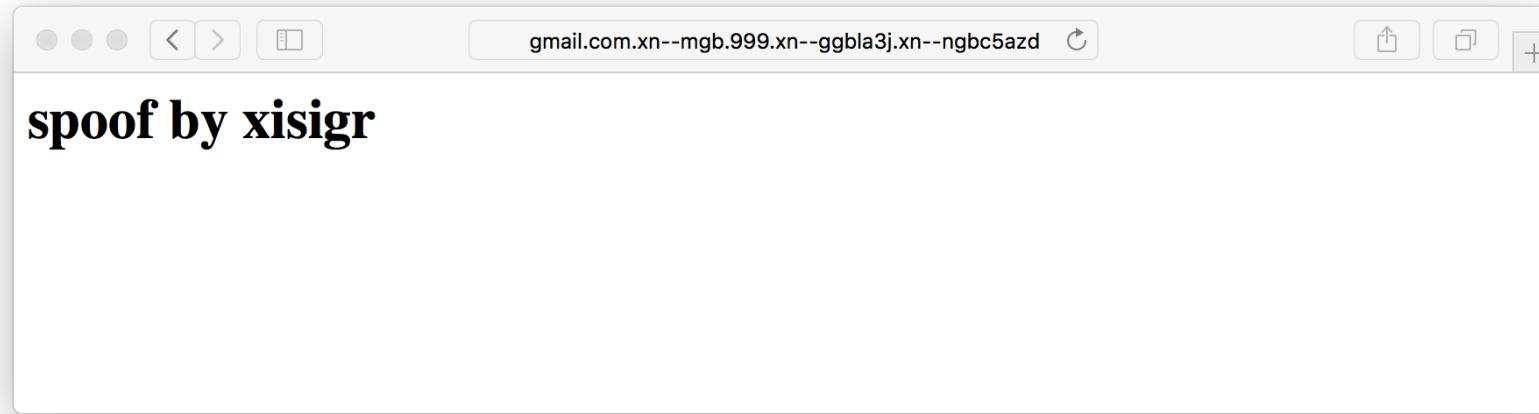


**www.gmail.com.xn--mgb.999.xn--ggbla3j.xn--ngbc5azd**

# Firefox/Chrome



# Safari/Edge



# Think about it



- How to resolve RTL Spoof in address bar?
- Only show origin?
  - showing only the origin reduces utility, though. That's for the UX designers to decide.
  - Either way, it doesn't fully resolve these issues. You can still do RTL spoofs with just the origin (like having the labels shown out-of-order).
  - Some error logic.
    - Eg: "www" or "m" subdomain is removed
    - <https://bugs.chromium.org/p/chromium/issues/detail?id=881694>
- Spoofable RTL URLs in the UI
  - <https://bugs.chromium.org/p/chromium/issues/detail?id=351639>

# Combining character



**Safari Address Bar Spoof Using Combining character  
CVE-2018-4260**

# Hebrew

	059	05A	05B	05C	05D	05E	05F
0	05A0	05B0	05C0	05D0	05E0	05F0	וּ
1	0591	05A1	05B1	05C1	05D1	05E1	וִי
2	0592	05A2	05B2	05C2	05D2	05E2	וַיְיָ
3	0593	05A3	05B3	05C3	05D3	05E3	וְיָ
4	0594	05A4	05B4	05C4	05D4	05E4	וְיָ
5	0595	05A5	05B5	05C5	05D5	05E5	וְיָ
6	0596	05A6	05B6	05C6	05D6	05E6	וְיָ
7	0597	05A7	05B7	05C7	05D7	05E7	וְיָ
8	0598	05A8	05B8	05C7	05D8	05E8	וְיָ
9	0599	05A9	05B9	05C7	05D9	05E9	וְיָ
A	059A	05AA	05BA	05C7	05DA	05EA	וְיָ
B	059B	05AB	05BB	05C7	05DB	וְיָ	וְיָ
C	059C	05AC	05BC	05C7	וְיָ	וְיָ	וְיָ
D	059D	05AD	05BD	05C7	וְיָ	וְיָ	וְיָ
E	059E	05AE	05BE	05C7	וְיָ	וְיָ	וְיָ
F	059F	05AF	05BF	05C7	וְיָ	וְיָ	וְיָ

וְיָ

hebrew letter vav (U+05D5)

- hebrew point holam (U+05B9)

וְיָ

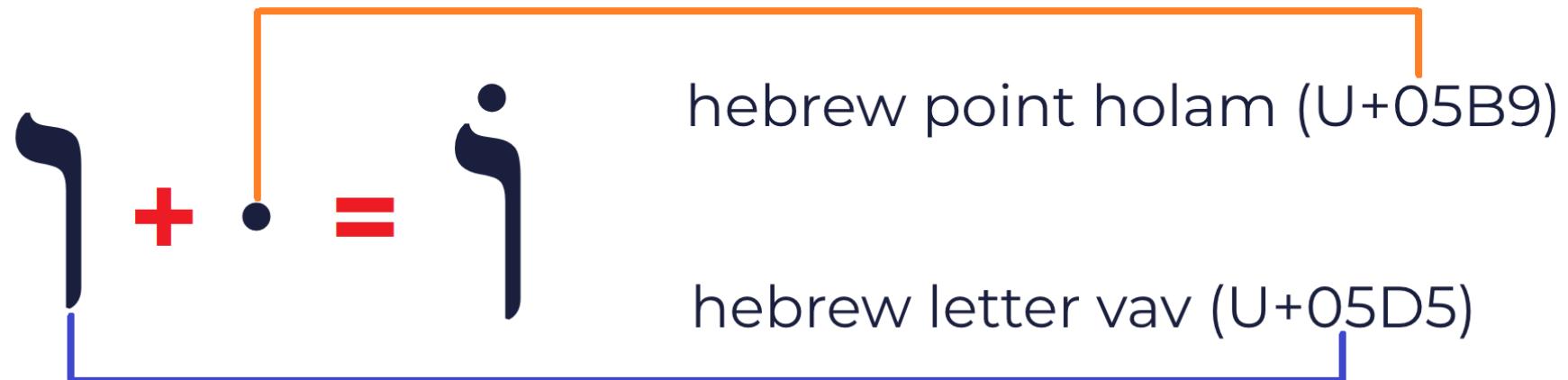
hebrew letter samekh (U+05E1)



# combining character sequence

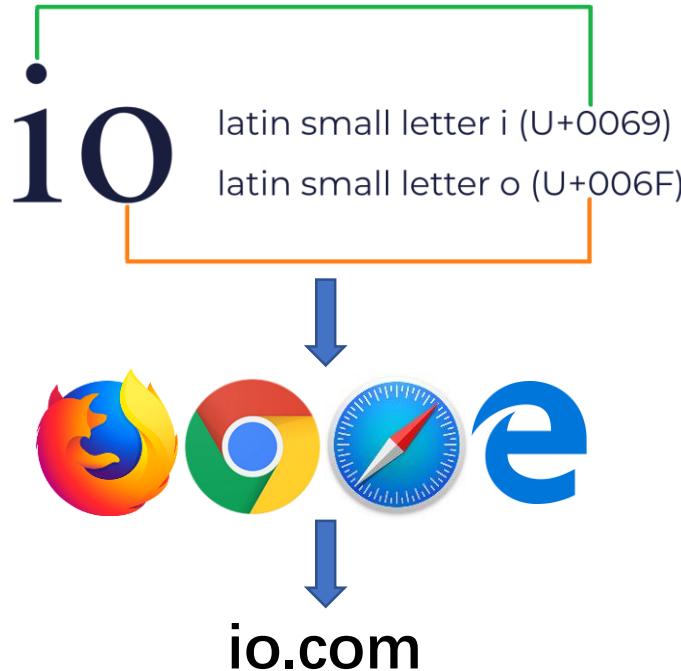


- Character + combining mark = *combining character sequence*

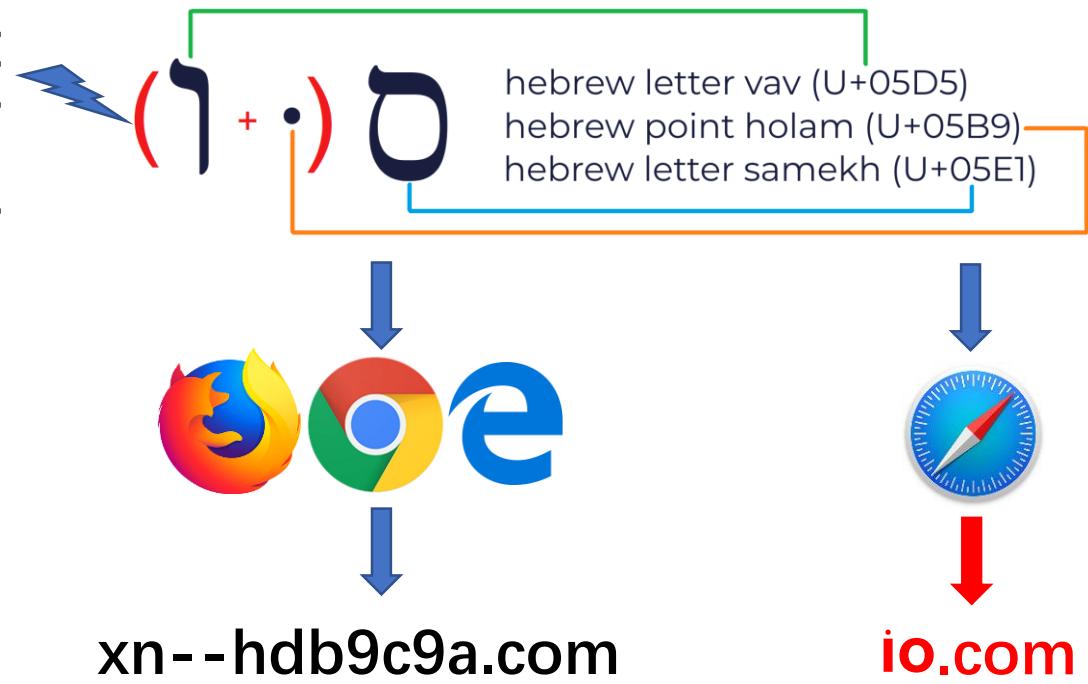




# Latin: io.com VS Hebrew: ካ.መ



*combining character sequence*



# Latin: io.com VS Hebrew: יו.com



The screenshot shows a web browser window with the URL "io.com" in the address bar. The page content features a background image of a data center with many blue and red cables. A dark banner across the middle contains the text: "Engineered for a new generation of needs and demands, IO delivers the data center as a service *@scale™*". Below this is an orange button with the text "SCHEDULE A TOUR".

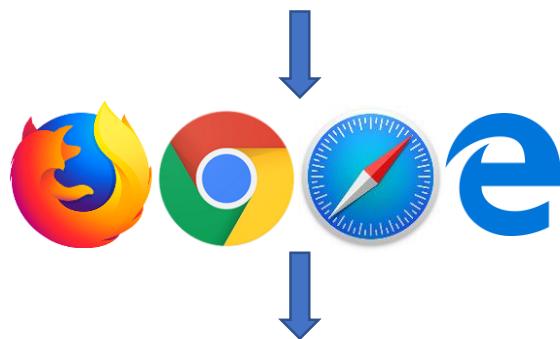
The screenshot shows a web browser window with the URL "io.com" in the address bar. The page content is a plain white page with the text "Fake io.com" in large bold letters. Below it is a long URL: "V6yADoAo/1Hm78m5/E/I7EgHJNefB1cimz1ZRngtjoMCOAFc+8a (Spoof by xisigr)".

# Hebrew alphabet: መ.መ VS Hebrew: መ.መ



መ

hebrew letter vav with holam (U+FB4B)  
hebrew letter samekh (U+05E1)



*combining character  
sequence*

( + )

መ

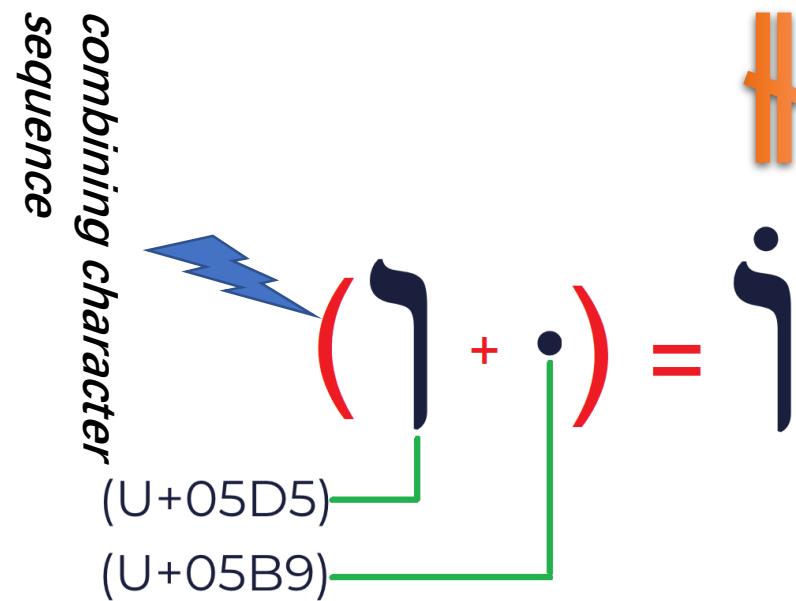
hebrew letter vav (U+05D5)  
hebrew point holam (U+05B9)  
hebrew letter samekh (U+05E1)





# grapheme is same but meaning is difference

## Alphabetic Presentation Forms



00	01	02	03
0x0	FB	4B	⚡

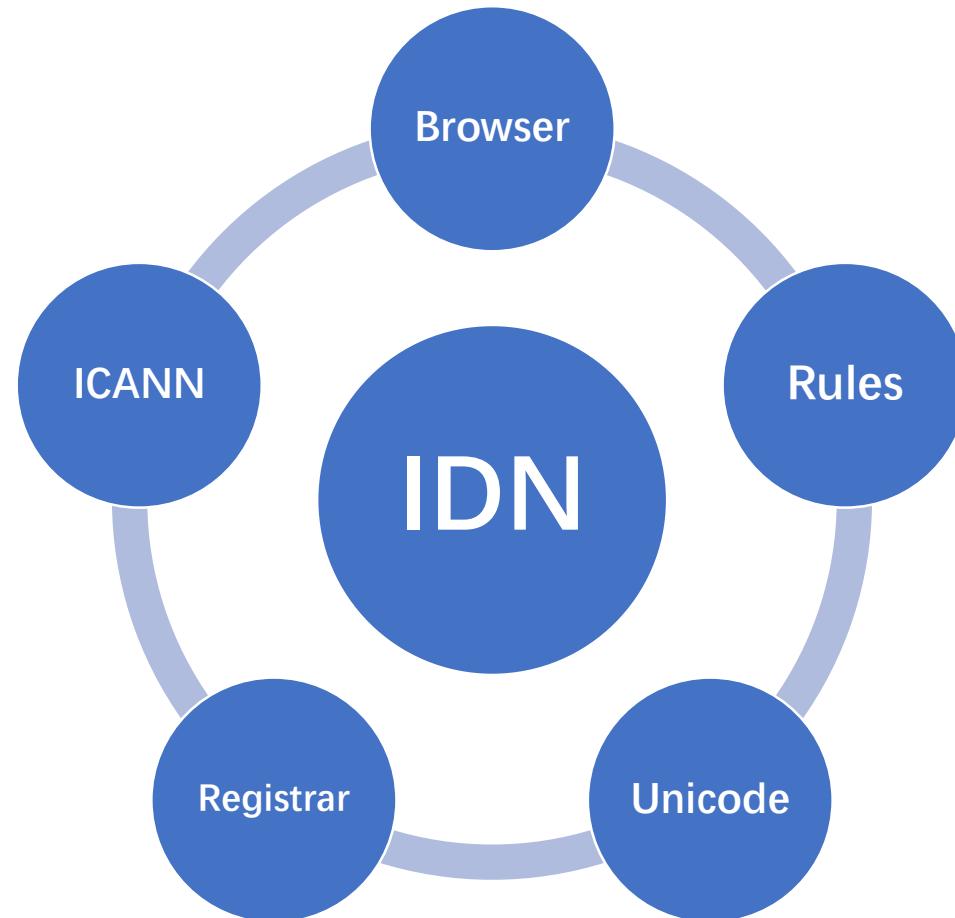
00	01	02	03
0x0	05	D5	05 B9
			⚡

# Think about it



- Should combination character sequence be supported to register?
- Is it possible to strictly check and prevent such malicious intentions at the Registrar stage?
- Think about these questions in the URL, you may find something new.
  - [http://unicode.org/faq/char\\_combmark.html](http://unicode.org/faq/char_combmark.html).

# IDN Challenging





Thinking of you