

Response to PRI 451

submitted by Asmus Freytag, 2022-06-14

After reviewing UTS#39 we found that there are a number of potential omissions in the `confusables.txt` data file. (For source and background see the end of this document)

Our analysis is based on ICANN's recent publications of Root Zone Label Generation Rules (RZ-LGR) and Second Level Reference LGRs for almost the complete set of "Recommended" scripts.

In these LGRs, a number of characters are considered mutually exclusive with either another character or a character sequence. This determination was made by panels of local experts and users. Where this exclusion is based primarily on appearance, we consider that an omission in the Unicode data file.

We therefore recommend that these be added to the data file before publication.

Missing Data in `Confusables.txt`

Digit zero

0030	0	0AE6	◦	Gujarati digit zero
0030	0	0CE6	○	Kannada digit zero
0030	0	0E50	◦	Thai digit zero
0030	0	0ED0	○	Lao digit zero
0030	0	1040	○	Myanmar digit zero
0030	0	17E0	◻	? Khmer digit zero
0030	0	0D20	○	Malayalam letter ttha
0030	0	101D	○	Myanmar letter wa – identical to digit zero

The case can be made that digits are ipso facto confusable semantically, as users may not keep track of which digit set is used in a label when both are available. However, when the shapes are also similar, the potential for confusion increases.

Already in `confusables.txt` are 09E6 and 0B66 which in the browser used for the screen shots at left look identical to 0AE6 or to 1040.

Digit three

0033	3	0AE9	3	Gujarati digit 3
------	---	------	---	------------------

Latin small letter c

0063	c	1004	ꠘ
------	---	------	---

Myanmar letter nga

Together with letter wa,
could be used to spoof .co

Latin small letter f

0066	f	0192	ƒ
------	---	------	---

Latin small letter f with hook

An argument has been made that 0192 is a familiar shape for this letter and users may not realize it's inappropriate for the typeface. Mutually exclusive in the Root Zone

Latin small letter i

0069	i	05D5	י
------	---	------	---

Hebrew letter vav

Latin small letter n

006E	n	0572	ղ
------	---	------	---

Armenian letter ghad

The Root Zone considers these mutually exclusive, but not 006E and 057C.

Latin small letter p

An argument can be made to consider this confusable with 01BF p Latin small letter wynn. Not only are the shapes close enough but few users know about the wynn and would take it for font idiosyncrasy. Wynn is excluded from the Root Zone for that reason.

Latin small letter s

0073	s	0D1F	ശ
------	---	------	---

Malayalam letter ttha

The Root Zone considers these mutually exclusive, partially because the font size difference disappears in whole-script labels not juxtaposed with Latin: sos. 日本

Latin letters with macron and tilde

00E3	ã	0101	ā
0067 0303	ḡ	1E21	ḡ
006E 0304	ñ	00F1	ñ

There's a generic argument that at typical type sizes macron and tilde become confusables of each other.

These pairs are mutually exclusive in the Root

00F5	õ	014D	ō
0129	ĩ	012B	ī

Latin small letter g with dot above

0121	ġ	0123	ġ
------	---	------	---

Latin small letter g with cedilla

The existing file has 0123 paired with 0127

Latin small letter eng

014B	η	03B7	η
------	---	------	---

Greek small letter eta

014B	η	0572	ղ
------	---	------	---

Armenian small letter ghad

Cyrillic small letter sha

0448	ш	0561	ա
------	---	------	---

Armenian small letter ayb

Cyrillic small letter dze

045F	ұ	1EE5	u̇
------	---	------	----

Latin small letter u with dot below

Arabic letter alef

0622	آ	0623	ا
0622	آ	0625	ا
0622	آ	0627	ا
0622	آ	0672	ا
0623	ا	0625	ا
0623	ا	0627	ا
0623	ا	0672	ا
0625	ا	0627	ا
0625	ا	0672	ا
0627	ا	0672	ا

The Arabic IDN task force (TFAIDN) concluded that the various forms of Alef should be considered mutually exclusive for domain names.

Arabic letter waw with hamza above

0624	ؤ	0648	و
------	---	------	---

Arabic letter waw

Arabic letter alef maksura (0649)

0626	ئ	0649	ى
------	---	------	---

0626	ئ	067B	پ
------	---	------	---

0626	ئ	06D0	ې
------	---	------	---



TFAIDN concluded that the existing set should be extended to include 0626, 067B and 06D0

Arabic letter the marbuta

0629	ة	0647	ه
------	---	------	---

0629	ة	06BE	هـ
------	---	------	----

0629	ة	06C0	ة
------	---	------	---

0629	ة	06C1	ه
------	---	------	---

0629	ة	06C2	ة
------	---	------	---

0629	ة	06C3	ة
------	---	------	---

0629	ة	06D5	ه
------	---	------	---



TFAIDN gives the full set as shown here. 06C3 is already listed, but not the other ones. These are mutually exclusive in the Root Zone.

Arabic letter teh

062A	ت	067A	ت
------	---	------	---

Arabic letter tteheh

Arabic letter feh

0641	ف	0642	ق
------	---	------	---

0641	ف	06A2	با
------	---	------	----

Arabic letter qaf, Arabic letter feh with dot moved below

Arabic letter noon

0646	ن	06BA	ن
------	---	------	---

Arabic letter noon ghunna

Arabic letter peh

067E	پ	06BD	پ
067E	پ	06D1	پ
067E	پ	0752	پ

Arabic letter nyeh

0683	ن	0684	ن
------	---	------	---

Arabic letter dyeh

Arabic letter dul

068E	ذ	068F	ذ
------	---	------	---

Arabic letter dal with three dots above downwards

Devanagari sign candrabindu

0901	ँ	0945 0902	ँ
------	---	-----------	---

Devanagari candra e + anusvara

Devanagari letter a + anusvara

0905 0902	अं	0973	अं
-----------	----	------	----

Devanagari letter oe

Devanagari digit 2

0968	ર	0AB0	ર
0968	ર	0AE8	ર

Gujarati letter RA, Gujarati digit 2

Devagari nukta

0906	आ	0906 093C	आ
0913	ओ	0913 093C	ओ
093E	ा	093E 093C	ा
094B	ो	094B 093C	ो

An argument can be made that the nukta (small dot below) placed on letters where this combination is not expected will not be noticed. The Root Zone makes these mutually exclusive.

Devanagari letter aa + anusvara

0906 0902	आं	0974	आं
-----------	----	------	----

Devanagari letter ooe

Devanagari letter i

0907	इ	0A19	इ
------	---	------	---

Gurmukhi letter nga

The complete set also contains some another sequence with a nukta

Devanagari letter u

0909	उ	0A24	उ
------	---	------	---

Gurmukhi letter ta

Devanagari letter short e

090E	ऐ	0910	ऐ
------	---	------	---

Devanagari letter ai

The LGRs have context rules that prevent the "fake" decompositions that are listed in confusables.txt for these four vowels. Instead, the Neo-Brahmi panel considers them to be pairwise confusable with each other.

Devanagari letter au

0914	औ	0975	औ
------	---	------	---

Devanagari letter aw

Devanagari cross-script confusables with Gurmukhi

0917	ग	0A17	ਗ
0918	घ	0A2C	ਬ
091F	ट	0A1F	ਟ
0920	ठ	0A20	ਠ
0922	ढ	0A2B	ਢ
0924 094D 0924	त्त	0A1C	ਜ
092A	प	0A27	ਧ
092A 094D 091F 0946	ऐ	092A 094D 091F 0947	
092A 094D 091F 0946	ऐ	0A0F	

The Neo-Brahmi panel considers these confusable to the point that they are mutually excluded.

092D	भ	0A2E	भ
092E	म	09AE	म
092E	म	0A38	म
0935	व	0A15	व
0939	ह	0A35	ह
093F	ि	0A3F	ि
0948	ै	0A48	ै
0956	ु	0A41	ु
0957	ू	0A42	ू

Devanagari vowel sign ooe (with and without Anusvara or nukta.)

093B	ो	093E 0902	ो
093B	ो	093E 093C 0902	ो

Devanagari cross-script confusables with Bengali

093F	ि	09BF	ি
------	---	------	---

Devanagari in-script and cross-script additional confusables for vowel sign short e

0946	े	0947	े
0946	े	0A47	े

Devanagari vowel sign au

094C	ौ	094F	ौ
------	---	------	---

Devanagari vowel sign au

Bengali letter ra

09B0	র	09F0	ৱ
------	---	------	---

Bengali letter ra with middle diagonal

Oriya vowel sign e

0B47	େ	1031	େ
------	---	------	---

Myanmar vowel sign e

Tamil letter o + lla

0B92	0BB3	ஒள	0B94	ஒள
------	------	----	------	----

Tamil letter au

Tamil cross-script confusables with Malaylam

0BAE	ഥ	0D25	ഥ
0BB5	ഖ	0D16	ഖ
0BC6	ഌ	0D46	ഌ
0BC7	ഐ	0D47	ഐ

Telugu cross-script confuables with Kannada

0C07	ಇ	0C87	ಇ
0C10	ಐ	0C90	ಐ
0C16	ಖ	0C96	ಖ
0C17	ಗ	0C97	ಗ
0C1D	ಝ	0C9D	ಝ
0C1F	ಟ	0C9F	ಟ
0C26	ದ	0CA6	ದ
0C28	ನ	0CA8	ನ

0C30	റ	0CB0	റ
0C33	ഴ	0CB3	ഴ
0C3F	ീ	0CBF	ീ
0C41	ു	0CC1	ു
0C43	ൂ	0CC3	ൂ

Malayalam letter rra: cross-script variants

0D31	ဂ	1002	ဂ
0D31	ဂ	10D8	ဂ

Myanmar letter ga, Georgian letter in

Sinhala letter iruyanna

0D8D	යා	0D9D 0DD8	යා
0D8D	යා	0DC3 0DD8	යා

Sinhala combined sequences

Sinhala letter eyanna

0D91	ඵ	0DB5	ඵ
------	---	------	---

SINHALA LETTER MAHAAPRAANA PAYANNA

Sinhala letter eeyanna

0D92	ඵ්	0DB5 0DCA	ඵ්
------	----	-----------	----

SINHALA LETTER MAHAAPRAANA PAYANNA +
SINHALA SIGN AL-LAKUNA

Sinhala (additional cases like the prev. three)

0D93	ඵඵ	0DB5 0DD9	ඵඵ
0D94	ඹ	0DB9	ඹ
0D9B	ඹ	0DB6	ඹ

0D9D	ઠ	0DC3	ઠ
0DA0	ઐ	0DC0	ઐ
0DB7	ઙ	0DC4	ઙ

Gujarati letter pa

0AAA	૫	0AEB	૫
------	---	------	---

Gujarati digit five

Compare also to the existing pair
04B7 ૫ 04CC ૫

Myanmar letter k + virama + ka

1000 1039 1000	က	1023	က
----------------	---	------	---

Myanmar letter i

Myanmar variant letter forms

1001	ခ	1076	ခ
1008	ဈ	105B	ဈ
101B 103E	ရှ	1061	ရှ
102B	ါ	102C	ါ
102E	ိ	1033	ိ

Considered confusable if presented
out of context or for the wrong
language

Myanmar letter nga + asat

1004 103A	င	1004 103A 1039	င
1004 103A	င	105A 103A	င
1004 103A	င	105A 103A 1039	င

The Root Zone and Reference LGRs treat these as mutually exclusive independent of how distinctions in appearance. However, note that the two sequences look identical.

Myanmar cross-script confusables with Georgian

1002	ጠ	10D8	ጠ
1010	ጡ	10D7	ጡ

Ethiopic

The Ethiopic script has a number of confusables that are based on phonetic equivalents rather than on visual similarity. The dominant language, Amharic, is commonly spelled phonetically, with apparent free alternation of homophones (for the same word). As if English had a random mixture of “lead” / “led”, “debt”/“det”, or “knight”/“night”/ “nite” and “knite”, with all forms equally acceptable in practice. And with the distinctions reduced to alternate letters, not sequences.

If this fits the Unicode definition of “confuable”, a list can be provided.

Korean Hangul confusables with Han Ideographs

4E2C	ㄱ	B258	ㄱ
723F	ㄴ	B258	ㄴ
535F	ㅁ	B9C8	ㅁ
4ECA	ㅅ	C2A5	ㅅ
5408	ㅎ	C2B4	ㅎ
4E1B	ㅆ	C4F0	ㅆ
4E15	ㅈ	C870	ㅈ
9577	長	D2BD	ㅌ

Comments on the sources for this set

These confusables were extracted mainly from Root Zone Label Generations Rules, Version 5 (RZ-LGR 5) a set of script specific repertoires for top-level IDNs that are combined with context rules (that exclude, among others, any sequences Unicode has declared as “do not use”). Any other duplicate spellings (or “close but not quite”) have been identified as “variants” and are mutually exclusive. This is equivalent to Unicode’s definition of “confusable”, except that the focus has been on cases that are either true

substitutions (users without ill intent may substitute one for the other) or those that are considered “practically indistinguishable” on visual grounds.

Some additional confusables were derived from work that ICANN is currently undertaking on Second Level Reference LGRs. These are model LGRs, often extensions of the corresponding Root Zone LGRs, that registries can use on the Second Level, with a similar attention to security.

For both sets, it is assumed that labels are restricted to a single script each, that is, no mixed-script labels are allowed (on the second level, some scripts may have ASCII add ins). However, labels of multiple scripts may coexist on a single zone, so the design includes confusables that can occur between two whole-script labels of different scripts.

There is no need to consider confusables arising from the application of combining marks out of context, because all combining marks are strictly context-limited. (And certainly cannot exist in a mixed script case). Likewise, there is no need to consider cross-script similarity for combining marks, because their allowed base characters are rarely also confusables for the same script pair.

In some cases, the proposal document for a given LGR includes a list of additional confusables that the generation panel of local experts thought did not make the cut. Some of these may well fit the slightly different criteria used in `confusables.txt`.

In determining their sets of confusables (or the actual visual variants included in the LGRs) the generation panels conducted various levels of research, from informal polling of their own members to formal research by a university. For many of the complex scripts, they also considered extensive lists of conjunct forms. The relevant details are described in their proposal documents for the Root Zone LGRs.

The easiest way to access these, is to go to <https://icann.org/idn> and look for “Root Zone LGR” and there for the list of proposals.

Because the layout of the data differs markedly between LGRs and the `confusables.txt`, the data for this current report was created by converting the latest `confusables.txt` into the LGR format and then restricting it to the characters found in the Maximal Starting Repertoire (plus digits and hyphen). The latter represents the de-facto superset of both RZ-LGR 5 and the current set of Second Level Reference LGRs.

This was manually compared to file containing the superset of variant definitions for the aforementioned LGRs, removing duplicates and variants defined for reasons unrelated to visual similarity. (As these are not identified as such, this step could not be automated).