Proposal to Add Data for Pairs of Confusable sequences

Proposed: Asmus Freytag 7 June 2022

This document proposes the addition a number of pairs of <u>confusable</u> sequences to the security (UTS#39) data files in analogy to the "confusables.txt" data collection for single code points.

A separate document proposes a set of eight identical sequences. These same are also absent from the collection of confusables. (Some singletons noted in comments below should also be added).

Background

ICANN recently published the Root Zone Label Generation Rules, Version 5 (see https://icann.org/idn for contents and overview). These Rules, collectively known and RZ-LGR-5, present a dedicated analysis of characters and sequences that users can be expected to freely substitute in the context of domain names. While some of them have strictly identical appearance, a larger set are either not readily distinguished or otherwise seen as alternations that users will substitute. The latter are the source for the current proposal.

The scripts covered in RZ-LGR-5 nearly exhaust the set of "Recommended" scripts from UA#31, while the actual repertoire is limited to characters and sequences that are in common, wide-spread and everyday use. They represent high-priority targets for security exposure but also mitigation efforts. Their absence from the Unicode confusable data represents thus a serious omission that should be remedied.

RZ-LGR-5 represents the culmination of a decade of effort by local volunteer panels composed of native speakers, linguists and technologists; ICANN staff; and a panel of technical experts combining linguistic, Unicode and IDN related expertise. At multiple stages throughout the process, public comments were solicited in a manner similar to the Unicode PRI process.

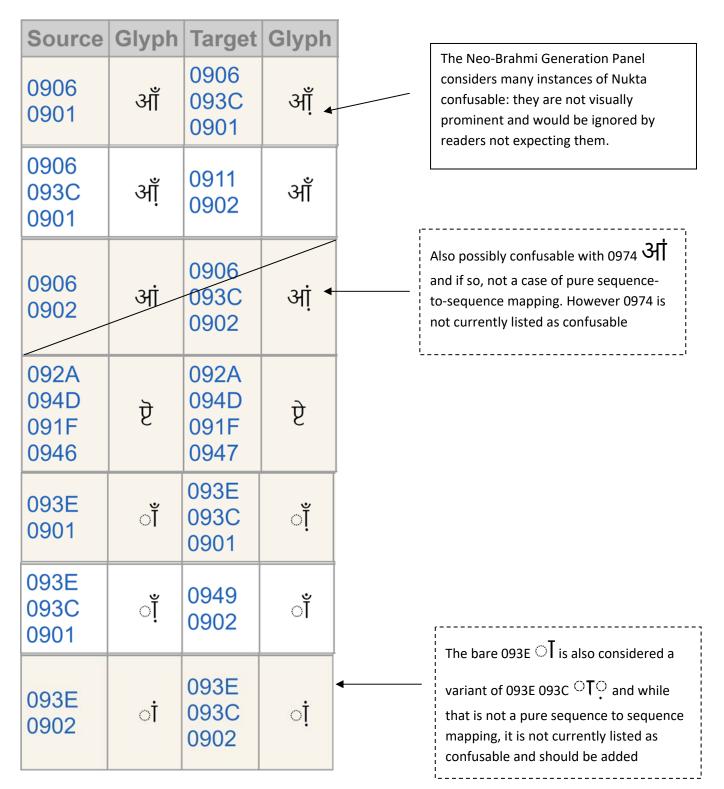
The rationale for all technical decisions is extensively documented, with each script LGR referencing the underlying proposal document, and, as the case may be, supportive data files.

There are other mappings (singleton to sequence, or singleton to singleton) that also have confusable appearance. This document focuses on the need to document property information, like confusable appearance for cases where it cannot be expressed as a character property, but must truly be considered a <u>property of strings</u>.

None of these mappings have any single code point (in the RZ-LGR-5 repertoire) that would be confusable with the sequence. No data file can be constructed that takes a single code point as the left-hand value. This makes them distinct from the larger set of confusable mappings that can be captured as a property of characters (and of which many, but not all, are already in the security (UTS#39) data files).

The following set of screen captures shows the sequences and their appearance at high resolution

Pairs of Confusable Sequences from RZ-LGR-5



09A8 09CD 09A5	ऋ	09A8 09CD 09B9	ন্হ	For the rationale for this and other similar confusables, please see https://icann.org/idn and look up the Root Zone LGR proposal document for the script.		
09B8 09CD 09A5	इ	09B8 09CD 09B9	স্হ			
0D28 0D4D	м	0D7B 0D4D	ൻ്	Also possibly confusable with 0D7B 00 and if so, not a case of pure sequence-to-sequence mapping. However 0974 is not currently listed as confusable		
0D31 0D31	ററ	0D31 0D4D 0D31	Ŋ	<u>i</u>		
0D31 0D31 0D4D 0D31	ററ്റ	0D31 0D4D 0D31 0D31	გი			
0D33 0D33	<u> </u>	0D33 0D4D 0D33	<u> </u>	The slight kerning caused by the virama was seen as not sufficiently distinct for IDNs.		
0D33 0D33 0D4D 0D33	<u> </u>	0D33 0D4D 0D33 0D33	<u> </u>			

		_	
0D9D 0DD8	සීබ	ODC3 ODD8	සෲ
1004 103A	ξ	1004 103A 1039	ీ
1004 103A	င်	105A 103A	Ć.
1004 103A	ć	105A 103A 1039	్
1004 103A 1039	ీ	105A 103A	CC.
105A 103A	CO _∞	105A 103A 1039	ఀ

Given that 0D9D and 0DC3 are already variants of each other, it's not clear why these sequence needed to be listed.

However, neither 0D9D nor 0DC3 are currently listed as confusable, and therefore should be added.

Note that 1004 C and 105A C take on the same appearance if followed by ASAT (see document on identical sequences), and may take on a rather not so distinct appearance if followed by some other diacritic (such as 103D). Shown here bracketed by 101D C on either side for context.

000000

The precise details depend on the renderer making any distinction fragile. The Root Zone LGR considers these two code points variants in any context where they are not followed by a consonant (or digits, if digits were allowed in the Root Zone). confusables.txt does not allow such context-based rules, which would therefore argue for adding all sequences of the form 1004+combining mark and 105A+combining mark, where the range of Myanmar combining marks starts at 102B and ends at 103E, but also the ranges 105E to 1060. 1062 to 1064. and 1081 to 108F.