# **Proposal to Encode Telugu Sign Nukta**

Vinodh Rajan <u>vinodh@virtualvinodh.com</u>
Shriramana Sharma <u>jamadagni@gmail.com</u>
Suresh Kolichala suresh.kolichala@gmail.com

This document proposes the encoding of the Telugu Sign Nukta in the Telugu block of the UCS.

### 1. Introduction

Nukta is a consonantal diacritic used in Indic scripts to extend the native character repertoire and denote non-native phonemes. It has been traditionally used to represent Perso-Arabic phonemes, and lately, English phonemes (particularly, /f/ and /z/). Many of the encoded Indic scripts including a Nukta-like character. These include historic scripts like Grantha and Siddham, where Nukta is a modern innovation.

While Nukta is a common feature of North Indic scripts (due to significant lexical borrowing from Persian and Arabic), Kannada is the only South Indic script to have a wide-spread use of Nukta to represent /f/ and /z/. In Tamil, & U+0B83 seemingly takes a Nukta-like role to represent those phonemes by prepending itself to consonants, as in &山 /f/ and &ஜ /z/ respectively. It would appear that among the major South Indic scripts only Telugu and Malayalam do not have any characters to fulfill the role of Nukta.

### 2. Nukta Proposals in Telugu

Desikacharyulu (2000?) while elaborating the design of one of the earliest Telugu Unicode fonts Pothana2000 laments the lack of Nukta in the UCS and makes the following statements:

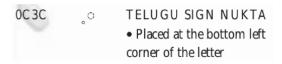
"[...] Nukta are not yet defined in the Telugu Unicode range, which they must be [...]"

"[...] it is still a good idea to have Nukta in Telugu to provide extendability [sic] of the script to foreign sounds"

He, therefore, included the Devanagari Nukta in his font as a placeholder, fully aware that it won't be rendered properly and hoping "it will be supported in the future".

One specific Nukta character has been quasi-proposed in the past for Telugu. A Government of India (GoI) publication (Rao, U.G, 2002) refers to a Nukta to extend the character set. However, there seems to be no evidence to its actual usage apart from its inclusion in the publication and a (passing) presentation to the Unicode consortium. There was no further formal follow up from the Government of India regarding the character<sup>1</sup>.

<sup>&</sup>lt;sup>1</sup> https://unicode.org/mail-arch/unicode-ml/y2009-m09/0006.html



Proposed Nukta from GoI (Rao, U.G, 2002)

## 3. Actual Nukta Usage in Telugu

Since Tamil is one of the liturgical languages of Vaishnavism in the South, Tamil texts have long been rendered in the Telugu script for the purpose of recitation in the Telugu-speaking areas of Southern India. However, Modern Telugu lacks an equivalent of the Tamil Consonant LLLA LP and couldn't faithfully represent the Tamil texts. This was overcome in several ways ranging from importing the Tamil character and treating it like a Telugu character and, quite relevant to our context, using a Nukta alongside Telugu LLA. (There is a Telugu LLLA CO U+0C34 that was encoded recently but it is an archaic character that was not known beyond epigraphic/historical circles and the common public was/is oblivious to its existence).

The below is an example from the Tamil text Tiruppāvai using both the conventions.



ழ as a Telugu letter (Link to the text)

ళ with Nukta as ళ (Kidambi, n.d-a)

Kidambi (n.d-a) clearly informs the readers on the Nukta convention used in the beginning of the document.

Attention: Please note that the letters and  $\mathfrak{B}$  and  $\mathfrak{B}$  denote  $\mathfrak{p}$  and  $\mathfrak{p}$  respectively, in Tamil. Also note that  $\mathfrak{B}$  sounds almost like  $\mathfrak{B}$ ,  $\mathfrak{B}$  like  $\mathfrak{B}$ , and so on. The consonant-cluster  $\mathfrak{F}$  is pronounced somewhere between  $\mathfrak{F}$  and  $\mathfrak{F}$ . It is, however, colloquially acceptable to pronounce the clusters  $\mathfrak{B}$  and  $\mathfrak{F}$  as  $\mathfrak{F}$  and  $\mathfrak{F}$ , respectively.

From Kidambi, S (n.d-a)

```
వాళ్లాట్ పట్టు నిస్టిరుశ్ఫీరేల్* వందు మణ్ణుం మణముం కొణ్మి € *
కూళ్లాట్ పట్టు నిస్టిర్గళై* ఎంగళ్ కుళ్లువినిల్ పుగుదలొట్టోం *
ఏళ్లాట్ కాలుం పళ్ళిప్పిలోం నాంగళ్* ఇరాక్కదర్నాళ్ల్* ఇలంగై –
పాళ్ళాళాగ పృడై పొరుదానుక్కు * పృల్లాండు కూఱుదుమే (3)
```

Tiruppallāntu (Kidambi, S, n.d-b)

வா**ழா**ட்பட்டுநின்றீருள்ளீரேல் வந்துமண்ணும் மணமும் கொண்மின் கூ**ழா**ட்பட்டுநின்றீர்களை எங்கள் கு**ழு**வினில் புகுதலொட்டோம் ஏ**ழா**ட்காலும் ப**ழி**ப்பிலோம் நாங்களிராக்கதர்வா**ழ்** இலங்கை பா**ழா**ளாக ப்படைபொருதானுக்குப் பல்லாண்டு கூறுதுமே

### Corresponding Tamil text.

It is the only case (we are aware of), where a dot-like Nukta sign is actually attested alongside a Telugu character. This form of Nukta as a dot and the Gol's Nukta as a small circle can very well be considered as glyphic variants. This is similar to the appearance of Virama in Tamil, which can appear as a dot or as a circle based on stylistic preferences. The text repository Prapatti: <a href="http://prapatti.com/slokas/slokasbyname.html">http://prapatti.com/slokas/slokasbyname.html</a> contains numerous Telugu documents using the Nukta sign (for this purpose).

Based on the actual attestation and its ability to impart extendibility to the Telugu script, we propose to include a Telugu Sign Nukta in the UCS.

# 4. Possible Objections

The attested LLA + dot below of are not glyphic variants of the already encoded Telugu LLLA cobut rather individual characters with different graphemic identities. The latter is an archaic character unknown to the majority of the populace and the former is a clear attempt to extend the alphabet and represent a non-native phoneme in a modern context. Orthographically speaking, both are distinct practices and users must be given a choice between the two, rather making it haphazardly font dependent.

It could also be said there are not any further attestations. The usage of Nukta has a consistent albeit restricted practice in at least one specific religious context/group to be considered as minority usage. Even in Siddham, the Nukta character was encoded based one modern usage from a single site committed to writing Japanese in Siddham, and in Grantha, it was entirely a modern innovation specifically requested by the user community. Also, as noted in the previous section, there is a user community that would readily consume the Nukta character if encoded.

Encoding a Telugu Nukta would especially enable the transliteration of Perso-Arabic consonants such as  $qa \ kha \ \dot{q}a \ za \ fa \ wa$  etc, the North Indic  $\underline{r}a$  and  $\underline{r}ha$ , the South Indic  $\underline{n}a$ , the East Indic  $\dot{y}a$  and so on when used in combination with appropriate existing characters. In a more modern context, it would especially allow Telugu to represent  $\underline{f}a$  and  $\underline{z}a$  unambiguously.

# 5. Suggested Transliteration Mapping

Based on the generic Indic model, we provide below suggested mappings to the existing Indic characters in Telugu using the proposed Nukta. Kannada is also shown as an existing cognate equivalent.

ISO 15919	qa	<u>kh</u> a	ġa	za	ŗa	ŗha	fa
Devanagari	क़	ख़	ग़	ज़	ड़	ढ़	फ़
Telugu	క్	ģ	ڔ۬	జ	డ	ۺ	<b>្</b> ឆ្
Kannada	호	ప్	ц	ఙ	ಡ಼	ಢ	ಫ಼

ISO 15919	у̀а
Devanagari	य
Bengali	য়
Oriya	ୟ
Telugu	య
Kannada	ಯ಼

ISO 15919	<u>l</u> a	nа
Devanagari	छ	न
Tamil	β	ன
Malayalam	φ	ഩ
Telugu	ණ ස	న
Kannada	ස	ನ:

Currently, the lack of a Telugu Nukta either leads to inaccurate transliteration or requires ad-hoc strategies to preserve one-to-one equivalence. For instance, *Aksharamukha* (<a href="http://aksharamukha.appspot.com/">http://aksharamukha.appspot.com/</a>) uses an ad-hoc spacing dot (U+00B7) as a Nukta equivalent

for Telugu, which is awkwardly placed next to the characters and does not blend well with the script. For instance, the word /zafār/ when expressed as జఫార్ is inaccurate, when using U+00B7 జ·ఫార్ looks awkward and unnatural. But the sequence జఫార్, using a Nukta looks more harmonious and native-looking. It would greatly aid transliteration if Telugu could have a native Nukta sign encoded. Some other sample words are showed below:

chattīsgaṛh छत्तीसगढ़ ఛత్తీస్గఢ్ tēṇi தேனी తేస్తి ālvār ஆழ்வார் ఆళ్వార్ (Archaic: ఆల్వార్)

malampu<u>l</u>a മലമ്പുഴ మలమ్పుళ్ల (Archaic: మలమ్పుట)

### 6. Placement of Nukta

In general, it is to position itself in the middle below a consonant. But it is recommended to move it a bit left with some consonants as shown below to avoid collision with the aspiration marker) and confusion (with existing characters). For the vowel sign /ai/, it is recommended that the second part of the sign is moved a bit lower to avoid collision with Nukta.

The below shows the Nukta sign ္ combining with క ఖ గ జ డ ఢ ఫ య ళ న alongside other vowel signs.

क्रूं क्र क्र क्षे क्षे क्ष क्षा ख़े ख़ें क्षे क्षे क्षे

क़िं क़ं क़ क़ै क़ैं क़ं क़ा क़े क़ें क़ क़ क़िं क़

ఫ్ ఫ ఫా ఫి ఫీ ఘ ఘా ఫె ఫే ఫై ఫొ ఫో ఫౌ

ర్నూ య యా రాము యా యూ ర్పూ య్లా ర్మూ ర్మూ య్హా

के के के के के के का का वे वे के के के के

న్ న నా ని నీ ను నూ నె నే నై నొ నో నౌ

### 7. Character to be encoded

It is proposed that Telugu Sign Nukta be encoded in the Telugu block of the UCS with the associated character properties.

# 0C3C ್ತ Telugu Sign Nukta

OC3C;TELUGU SIGN NUKTA;Mn;7;NSM;;;;N;;;;

## 8. Indic Syllabic Category

The following addition should be made to the IndicSyllabicCategory.txt file under:

## 9. Indic Positional Category

The following addition should be made to the IndicPositionCategory.txt file under:

### 10. Collation

As the Nukta is meant for transcribing sounds which are not native to Telugu, it is recommended that consonants with the Nukta are collated after consonants without the Nukta so as to not disturb the existing Telugu collation order.

#### References

- 1. Desikacharyulu, T.K. (2000?). Creating Unicode Compatible OpenType Telugu Fonts <a href="https://upload.wikimedia.org/wikipedia/te/c/c2/Pothanapaper.PDF">https://upload.wikimedia.org/wikipedia/te/c/c2/Pothanapaper.PDF</a>
- 2. Rao, U. G. (2002). Telugu Script. Vishwabharat@TDIL. Issue 5. <a href="http://tdil.meity.gov.in/pdf/Vishwabharat/tdil-april-2002.zip">http://tdil.meity.gov.in/pdf/Vishwabharat/tdil-april-2002.zip</a>
- 3. Kidambi, S. (n.d-a). ఆండాళ్ అరుళిచ్చెయ్ద తిరుప్పావై (Āṇḍāḷ aruḷicceyda Tiruppāvai). http://prapatti.com/slokas/telugu/naalaayiram/aandaal/tiruppaavai.pdf
- 4. Kidambi, S. (n.d-b). పెరియాళ్**వార్ అరుళిచ్చె**ర్దు తిరుప్పల్లాండు (Periyā<u>l</u>vār aruļicceyda Tiruppallāmḍu).
  - http://prapatti.com/slokas/telugu/naalaayiram/periyaazvaar/tiruppallaandu.pdf

# ISO/IEC JTC 1/SC 2/WG 2

# PROPOSAL SUMMARY FORM TO ACCOMPANY SUBMISSIONS FOR ADDITIONS TO THE REPERTOIRE OF ISO/IEC 10646.2

Please fill all the sections A, B and C below.

Please read Principles and Procedures Document (P & P) from <a href="http://std.dkuug.dk/JTC1/SC2/WG2/docs/principles.html">http://std.dkuug.dk/JTC1/SC2/WG2/docs/principles.html</a> for guidelines and details before filling this form.

Please ensure you are using the latest Form from <a href="http://std.dkuug.dk/JTC1/SC2/WG2/docs/summaryform.html">http://std.dkuug.dk/JTC1/SC2/WG2/docs/summaryform.html</a>.

See also <a href="http://std.dkuug.dk/JTC1/SC2/WG2/docs/roadmaps.html">http://std.dkuug.dk/JTC1/SC2/WG2/docs/summaryform.html</a>.

#### A. Administrative

1. Title:	Proposal to Encode Tel	ugu Sign Nukta	
2. Requester's name:	Vinodh Rajan, Shriramana S	Sharma, Suresh Kolichala	
3. Requester type (Member body/Liaison/	ndividual contribution):	Individual	
4. Submission date:		11/12/2019	
5. Requester's reference (if applicable):			
<ol><li>Choose one of the following:</li></ol>			
This is a complete proposal:			Yes
(or) More information will be provi	ded later:		
B. Technical – General			
1. Choose one of the following:			
<ul> <li>a. This proposal is for a new script (</li> </ul>	set of characters):		
Proposed name of script:			
<ul> <li>b. The proposal is for addition of cha</li> </ul>	racter(s) to an existing block:		Yes
Name of the existing block:		Telugu	
2. Number of characters in proposal:			1
3. Proposed category (select one from be	ow - see section 2.2 of P&P de	ocument):	
A-Contemporary A B.1-Specializ		B.2-Specialized (large colle	ection)
C-Major extinct D-Attested ex	tinct	E-Minor extinct	
F-Archaic Hieroglyphic or Ideographic	G-Obso	cure or questionable usage	symbols
4. Is a repertoire including character name	es provided?		
a. If YES, are the names in accorda		g guidelines"	
in Annex L of P&P document			Yes
b. Are the character shapes attache	d in a legible form suitable for	review?	Yes
5. Fonts related:	-		
a. Who will provide the appropriate of	computerized font to the Project	ct Editor of 10646 for publish	hing the
standard?	, , , , , , , , , , , , , , , , , , , ,		3
	Vinodh Rajan		
<ul> <li>b. Identify the party granting a licens</li> </ul>			il, ftp-site, etc.):
V	inodh Rajan, vinodh@virtualvii	nodh.com	<u>-</u>
6. References:			
<ul> <li>a. Are references (to other characte</li> </ul>			Yes
<ul> <li>b. Are published examples of use (s</li> </ul>	uch as samples from newspap		ources)
of proposed characters attached?		Yes	
7. Special encoding issues:			
Does the proposal address other as			
presentation, sorting, searching, ind			n)? <u>Yes</u>
	Sorting, Transliteration	)	
8. Additional Information:			
Submitters are invited to provide any addi			
that will assist in correct understanding of	and correct linguistic procession	ng of the proposed characte	er(s) or script.
Examples of such properties are: Casing i			

Submitters are invited to provide any additional information about Properties of the proposed Character(s) or Script that will assist in correct understanding of and correct linguistic processing of the proposed character(s) or script. Examples of such properties are: Casing information, Numeric information, Currency information, Display behaviour information such as line breaks, widths etc., Combining behaviour, Spacing behaviour, Directional behaviour, Default Collation behaviour, relevance in Mark Up contexts, Compatibility equivalence and other Unicode normalization related information. See the Unicode standard at <a href="http://www.unicode.org/reports/tr44/">http://www.unicode.org/reports/tr44/</a>) and associated Unicode Technical Reports for information needed for consideration by the Unicode Technical Committee for inclusion in the Unicode Standard.

<sup>&</sup>lt;sup>2</sup> Form number: N4502-F (Original 1994-10-14; Revised 1995-01, 1995-04, 1996-04, 1996-08, 1999-03, 2001-05, 2001-09, 2003-11, 2005-01, 2005-09, 2005-10, 2007-03, 2008-05, 2009-11, 2011-03, 2012-01)

## C. Technical - Justification

2. Has contact been made to members of the user community (for example: National Body,				
user groups of the script or characters, other experts, etc.)?				
If YES, with whom? Suresh Kolichala & Vinodh Rajan				
If YES, available relevant documents:				
3. Information on the user community for the proposed characters (for example:				
size, demographics, information technology use, or publishing use) is included?				
Reference: See Proposal				
4. The context of use for the proposed characters (type of use; common or rare)  Common				
Reference: See Proposal				
5. Are the proposed characters in current use by the user community?  Yes				
If YES, where? Reference: See Proposal				
6. After giving due considerations to the principles in the P&P document must the proposed characters be entirely				
in the BMP?				
If YES, is a rationale provided?				
If YES, reference: Telugu is in BMP				
7. Should the proposed characters be kept together in a contiguous range (rather than being scattered)? Yes				
8. Can any of the proposed characters be considered a presentation form of an existing				
character or character sequence?  No				
If YES, is a rationale for its inclusion provided?				
If YES, reference:				
9. Can any of the proposed characters be encoded using a composed character sequence of either				
existing characters or other proposed characters?  No				
If YES, is a rationale for its inclusion provided?				
If YES, reference:				
10. Can any of the proposed character(s) be considered to be similar (in appearance or function)				
to, or could be confused with, an existing character?				
If YES, is a rationale for its inclusion provided?				
If YES, reference:				
11. Does the proposal include use of combining characters and/or use of composite sequences?				
If YES, is a rationale for such use provided?				
If YES, reference: See Proposal				
Is a list of composite sequences and their corresponding glyph images (graphic symbols) provided?  If YES, reference:				
12. Does the proposal contain characters with any special properties such as				
control function or similar semantics?				
If YES, describe in detail (include attachment if necessary)				
13. Does the proposal contain any Ideographic compatibility characters?				
If YES, are the equivalent corresponding unified ideographic characters identified?				
If YES, reference:				