

Revised Proposal to Encode Alternate BA for the Bengali Language

Vinodh Rajan vinodh@virtualvinodh.com

Deepro Chakraborty deepro@ualberta.ca

17-Dec-2022

Introduction

Current mainstream Bengali orthography does not possess the equivalent of the Brahmic letter VA¹ (in Unicode parlance). Both VA and BA in Sanskrit-origin words are orthographically spelled using U+09AC ব̄ (BENGALI LETTER BA). As a result, there is no way to distinguish them in the resulting text. While this isn't an issue for writing the Bengali language (as it lacks /v/ entirely), the conflation presents a problem with Sanskrit texts, where they are distinct phonemes and require clear orthographic distinction.

(The Bangla Academy recommends [7] the orthographic differentiation of, at least, the clusters using an explicit virama. For instance, দ্ব /d̪ba/ vs দ্ব̄ /d̪[:]ɔ/ (< /*d̪ba/). Seemingly, this convention is not widely followed [4]).

This lack of a traditional letter to represent the letter VA is common to Assamese, Oriya, and Bengali. The former two eventually innovated new letters to differentiate the two phonemes by using slightly modified forms of BA to represent VA. Assamese innovated U+09F1 ব̄ (ব̄ BA with lower diagonal) and Oriya U+0B35 ବ̄ (ବ̄ BA with nukta).

Differentiating VA and BA in the Bengali orthography

In a similar vein, some Bengali publications³, particularly those printing Sanskrit and Pali (see L2/22-278) works using the Bengali orthography, also wanted to distinctly denote the two phonemes in their texts. As a result, they innovated a new letter to distinguish VA from BA in their publications. Unlike Assamese and Oriya, this was however done by using a distinct letter for BA: ব̄ (ব̄ with an inner diagonal stroke) and ব̄ was relegated to only represent VA.

From [4]:

জ্যায়সী চেৎ কর্মণস্তে মতা বুদ্ধিজর্নাদন ।
তৎ কিং কর্মণি ঘোরে মাং নিয়োজয়সি কেশব

¹ It is variously realized as /v/ or /w/ depending on the region

² Word-initially /d̪ɔ/ but word-internally /d̪[:]ɔ/

³ At least one linguistic publication [5] [8] has used ব̄ (BHA + nukta) to denote /v/. But this has not been widely adopted by publishers. One plausible reason might be that an average reader will likely ignore the nukta and will simply read it as /bha/.

জ্যায়সী চেৎ কর্মণস্তু মতা बुद्धिर्जनार्दन ।
 तं किं कर्मणि घोरे मां नियोजयसि केशव

jyāyāsī cet karmaṇaste matā buddhirjanārdana .
 tat kiṃ karmaṇi ghore māṃ niyojayasi keśava

It has been claimed such distinction can be traced back to the pre-modern/medieval era.

From [6]:

ব এবং ব।

২৩ সূত্র। আদি ভাষায় অন্ত্যস্থ বকারের আকৃতি এবং উচ্চারণ উভয়ই বর্গীয় বকার হইতে বিভিন্ন। পণ্ডিতের হস্ত লিখিত সংস্কৃত পুস্তক সমূহে বর্গীয় বকারের ব এইরূপ আকৃতি লিখিত হয়। কিন্তু বাঙালা ছাপার বর্ণ মালায় ব এবং ব উভয়ই ব সদৃশ লিখিত হয় এবং আকৃতি তুল্যতা হেতু উচ্চারণও তুল্য হইয়া গিয়াছে। এই অনিষ্ট নিবারণ জন্ত আমি ব কারের সংশোধন করিলাম। অতঃপর ব কার ইংরেজী V নামক বর্ণের স্থায় এবং ব ইংরেজী B নামক বর্ণের স্থায় উচ্চারণ করা উচিত।

ba and va

Rule: 23: In the earlier language, both the shape and pronunciation of the antyastha (sic) va were different than the vargiya ba. In Sanskrit texts handwritten by pandits the shape of a vargiya ba is written as ব̄. But in the Bengali printed alphabet, both ba and va are identical and because of their similar shape, their pronunciation has also become similar. In order to prevent this undesired thing, I have emended the letter ব̄. Afterwards, the letter ব̄ should be pronounced like the English letter V and ব̄ should be pronounced like the English letter B."

This is contentious, as there is no evidence in the manuscripts to claim the same. This usage was most likely inspired by the modern Devanagari letter forms - BA (ब) vs VA (व).

This new alternate BA shows analogous behavior to that of ব̄ in terms of vowel sign placements and the shapes of the various conjuncts as seen below. (The default behavior of the similar-looking U+09F0 is also shown for comparison)

	U+09AC	Bengali Alternate BA	U+09F0
	ব	𑂔	𑂔
Cu	বু	𑂕	𑂕
Cū	বু̄	𑂕̄	𑂕̄
rCa	ব্	𑂔̂	ব্
Cra	ব্র	𑂔̂	ব্র
Cru	ব্রু	𑂔̂	ব্রু
Cya	ব্য	ব্য়	ব্য
Cda	ব্দ	ব্দ	ব্দ
Cdha	ব্ধ	ব্ধ	ব্ধ
mCa	ম্ব	ম্ব	ম্ব
CCa	ব্ব	ব্ব	ব্ব

There are several Sanskrit publications that have used this character to print their Sanskrit texts. There are also few [blogs](#) that advocate using this distinction for the Bengali language. Clarifying the representation of the character in Unicode will facilitate the digitization of those published works. It will also allow anyone to digitally represent Sanskrit texts in the Bengali orthography without any loss of textual fidelity.

Representing the alternate BA in Unicode

Re-using U+09F0 Assamese RA

This is probably the most straight-forward solution. The Assamese RA already encoded at U+09F0 as BENGALI LETTER RA WITH MIDDLE DIAGONAL is visually similar to that of the Bengali alternate BA, even though the behavior of the two letters varies considerably as seen above. While some of the expected behavior can be achieved using joiners, others are problematic.

	Bengali Alternate BA	Sequence
bu	বু	U+09F0 ZWNJ VS-U
bū	বু̄	U+09F0 ZWNJ VS-UU
rba	ব্	U+09B0 VIRAMA U+09F0
bra	ব্র	U+09F0 ZWJ VIRAMA U+09B0
bru	ব্রু	U+09F0 ZWJ VIRAMA U+09B0 VS-U
bya	ব্য	U+09F0 ZWJ VIRAMA YA
bda	ব্দ	?
bdha	ব্ধ	?
mba	ম্ব	?
bba	ব্ব	?

Issues arise when we try to represent the conjunct forms. Unicode (or at least the fonts) treats U+09F0 as an allograph of U+09B0 Bengali letter RA. As such, U+09F0 forms a *repha* at the cluster initial positions and a *ra-phalā* at cluster-final positions, imitating the behavior of U+09B0.

ৰ্ক /rka/ - U+09B0 VIRAMA KA (OR) U+09F0 VIRAMA KA

ক্র /kra/ - KA VIRAMA U+09B0 (OR) KA VIRAMA U+09F0

As an allograph of RA, as per Unicode conventions, U+09F0 should produce a *repha* (C1-conjoining form) with the sequence U+09F0 VIRAMA ZWJ C2 and a *ra-phalā* (C2-conjoining form) with C1 ZWJ VIRAMA U+09F0. For those conjuncts marked in red, there aren't any reasonable sequences with ZWJ that will give the desired output.

It is possible that a font-level override might achieve the expected behavior, but it would violate the expected standard character behavior for U+09F0. Furthermore, even a font-level override would require additional joiners to cajole the rendering engines to treat it as a normal sequence and avoid unwanted re-ordering of the syllables (due to *reph* requiring special treatment). Below are sample renderings in MS Word of conjunct /bda/ using the alternate BA (with U+09F0 removed from the *reph* feature & /bda/ ligature defined in *akhn*)

দৰ্ দৰ্ দৰ্	U+09F0 VIRAMA DA
ব্ ব্ ব্	U+09F0 ZWJ VIRAMA DA
ব্ ব্ ব্	U+09F0 VIRAMA ZWJ DA

Also, when the appropriate font is not available it will result in an entirely erroneous and misleading rendering with no way to verify it unless one checks their underlying encoding. (e.g. শব্দ /śabda/ will be rendered as শর্দ/śarda/ in the absence of an appropriate font). In essence, it would just be a font-level hack

In essence, it would just be a font-level hack.

Encoding a new character

This would be the cleanest option of them all. The previous solution requires a very extensive use of joiners and deviation from the established behavior of U+09F0. It is already an overloaded character that requires special processing by the rendering engines. Overloading the character with joiners to create more specialized forms is counterproductive.

Given that Unicode is trying to move away from specifying complicated joiner-based solutions, a new character would be an appropriate solution after all. Even for Bengali, Khanda TA was atomically encoded to avoid using joiners to render the special form of vowelless TA.

This will, however, result in the encoding of two characters that look similar but behave differently. This won't be the first time such characters are encoded in the UCS. Similar characters already exist for Arabic where they are only differentiated in isolated/final forms.

Bengali already has ambiguous sequences due to the existence of two RAs and the encoding of the sequence is already decided based on the language of the text. This would be another language-specific addition to the Bengali block.

Bengali	ক্র/rkra/	U+09B0 U+09CD U+0995 U+09CD U+09B0
Assamese	ক্র/rkra/	U+09F0 U+09CD U+0995 U+09CD U+09F0
Bengali	ৰৰ/barba/	U+09FF ⁴ U+09B0 U+09CD U+09FF
Assamese	ৰৰ/rarra/	U+09F0 U+09F0 U+09CD U+09F0

Therefore, it is recommended that a new character BENGALI ALTERNATE BA be encoded with the character properties as shown below.

ৰ 09FF Bengali Letter Alternate BA

09FF;BENGALI LETTER ALTERNATE BA;Lo;0;L;;;;;N;;;;;

The following needs to be added to the *IndicSyllabicCategory.txt* file:

09FF ; Consonant # Lo BENGALI LETTER ALTERNATE BA

As a rarely used letter, we recommend that the letter be excluded from domain names for security purposes.

Collation

Normal Collation

As a rarely used letter it can be the last consonant in the collation following YA.

[...] < YA < ALTERNATE BA < AVAGRAHA < [...]

Tailored collation for Sanskrit

For Sanskrit, it should collate as:

[...] < PHA < ALTERNATE BA < BHA < [...] < YA < RA < LA < BA < BA WITH LOWER DIAGONAL < [...]

⁴ Proposed codepoint

Attestations

राज्ञा - वयस्य ! कः सम्पदहः। अचिन्त्या हि मणिमन्त्रोष्धीनां प्रभावः। पश्य -
 कष्टे शीपुर्बुधोत्तमस्य समरे दुष्टा मणिं शत्रुभिः
 नष्टं मन्त्रबलैर्वसन्ति वसुधामूले भुङ्क्ता हताः ।
 पूर्वं लक्षणवीरवानरभटा ये मेघनादाहताः
 पीता तेऽपि महोष्णेषुर्गनिषेर्गङ्गं पुनर्जीविताः ॥ ५ ॥

मन्त्रबलैर्वसन्ति /mantrabalairvasanti/ [1]

ङ्गसीसिचुम्बिआइं भमरेहिं सुडमारकेसरसिहहिं ।
 ओदंसअन्ति दअमाणा पमदाओ सिरीसकुसुमाइं ॥ ४ ॥
 (ङ्गदीषचुम्बितानि भमरैः सुकुमारकेशरशिखानि ।
 अवतंसयन्ति दयमानाः प्रमदाः शिरीषकुसुमानि ॥)

चुम्बिआइं /cumbiāim/ [2]

न्पित्रोर्विग्रहवग्रहं वि
 नाट्यवेदाङ्गिमालोद्य ता
 स्वात्नानाभिनयस्तं तं प्रण

नाट्यवेदाङ्गिमालोद्य /nāṭyavedāṅghimāloḍya/ [2]

গোচরো গুণঃ শব্দার্থো যস্যোঃ সা। ‘শ্রুতিঃ তে
 বিশ্বঃ। বিশ্বং জগদ্ ব্যাপ্য স্থিতা তেনাকাশঃ। অ
 শ্রুতিবিষয়গুণেত্যেতাবন্মাত্রং নোপাত্তম্। তাবতু
 চ ‘শব্দৈকগুণ আকাশঃ শব্দস্পর্শগুণো মর
 শব্দগুণত্বস্যোক্তোর্ববন্ধিতার্থালাভশ্চ। পারিশেষ
 হস্তেন বেতি সংদেহস্য দুর্নিরাসত্বাৎ সূক্তং বিশ্বমি
 যাং চ সর্বেষাং বীজানাং প্রকৃতিয়োনিরিত্যা

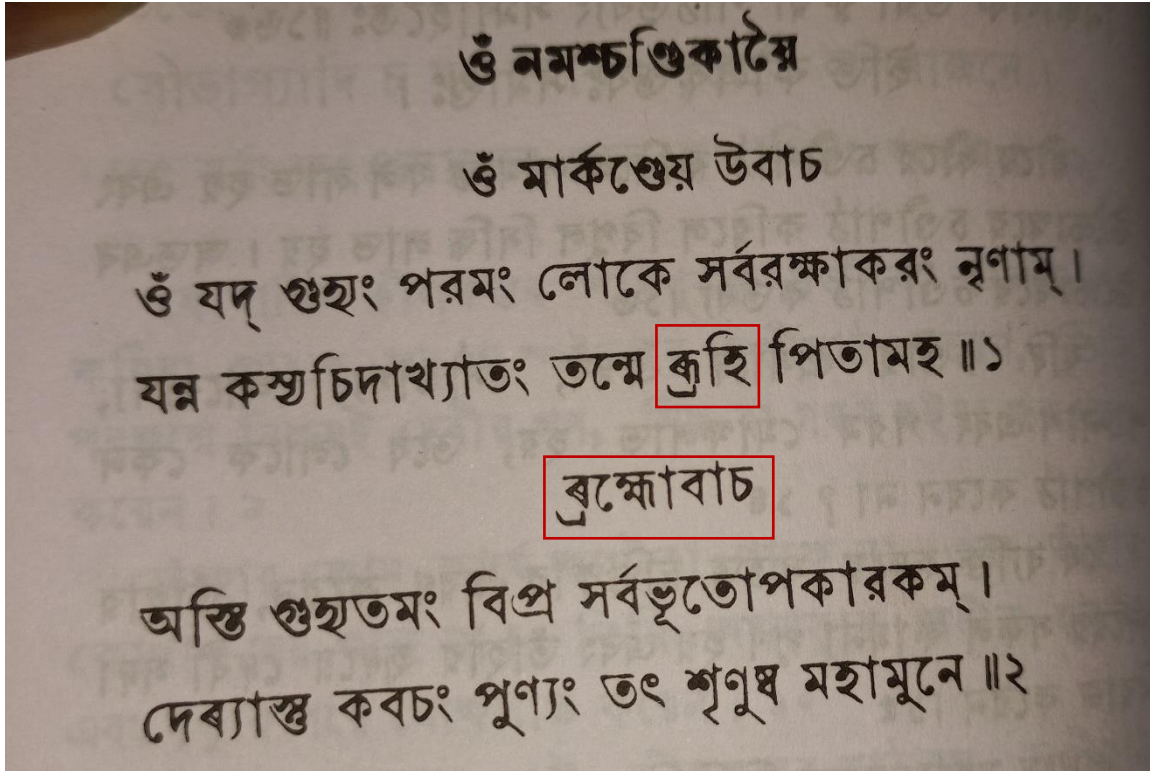
শব্দৈকগুণ আকাশঃ শব্দস্পর্শগুণো /śabdaiḥkaḡuṇa ākāśaḥ śabdaspārśagaṇo / [2]

বৈদিক গদ্যের নমুনা :
 প্রজাপতিঃ প্রজা অসৃজত তা অস্মাৎ সৃষ্টাঃ পরাচীরায়ন্তা বরুণমগচ্ছন্তা অশ্বৈত্তাঃ
 পুনরযাচত তা অশ্বৈ ন পুনরদদাৎ সোহব্রবীদ্ বরং বৃণীষ।—তৈ. সং. ২।১।২

সোহব্রবীদ্ /so'bravid/ [3]

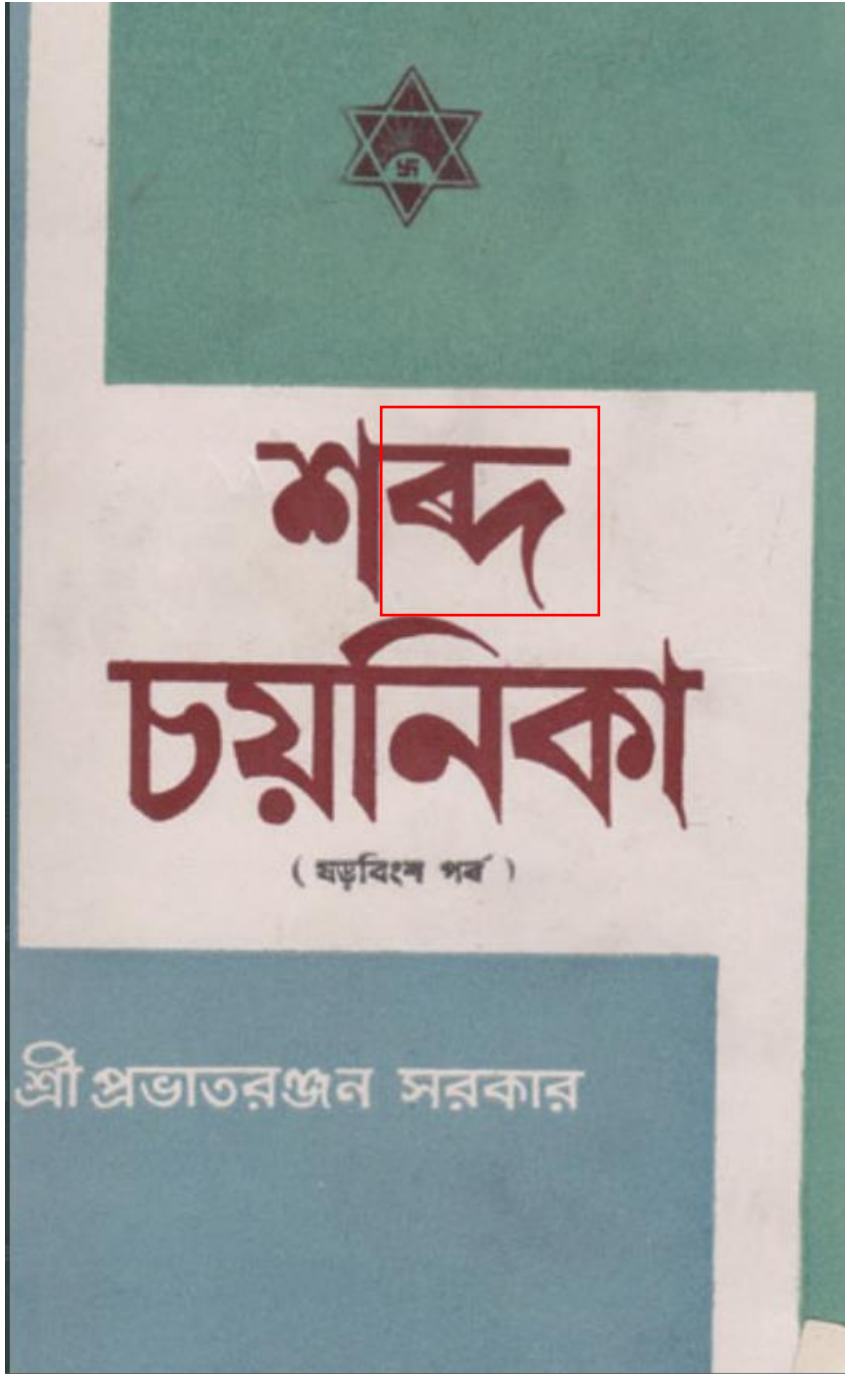
(গ) সুবন্ধুর্বাণভট্টশ্চ / কবিরাজ ইতি ত্রয়ঃ।
 বন্ধোক্তিমার্গানিপুণাশ্ / চতুর্থো বিদ্যাতে ন বা ॥ রাঘবপাণ্ডবীয় ১।৪১
 (ঘ) শব্দার্থয়োঃ সমো গুণ্যঃ / পাঞ্চালী রীতিরিষ্যতে।
 শীলাভট্টারিকাবাচি / বাণোক্তিশ্চ সা যদি ॥ শার্ঙ্গ. (Peterson) ১৭১

সুবন্ধুর্বাণভট্টশ্চ /subandhurbāṇabhṭṭaśca/ [3]



ব্রুহি /brūhi/

ব্রহ্মোবাচ /brahmovāca/ From [4]



শব্দ চয়নিকা /śabda cayanikā/ [9]

Acknowledgements

Thanks to Milind Chakraborty for providing additional attestations for the character and clarifying its usage. Publications from Ananda Marga were provided by Biswajit Mandal.

References

1. সেনগুপ্তা, আচার্য্য জ্যোতি (2016). মহাকবিশ্রীহর্ষদেবকৃত রত্নাবলী. সংস্কৃত বুক ডিপো.
Senaguptā, Ācāryya Jyoti (2016). Mahākabiśrīharṣadebakṛta Ratnābalī. Saṁskṛta Buka Dipo.
2. চক্রবর্তী, সত্যনারায়ণ (2007). মহাকবি -কালিদাস-প্রণীতম্ অভিজ্ঞান-শকুন্তলম্. Sanskrit Pustak Bhandar.
Cakrabartī, Satyanārāyaṇa (2007). Mahākabi-Kālidāsa-Praṇītam Abhijñāna-Śakuntalam. Sanskrit Pustak Bhandar.
3. বন্দ্যোপাধ্যায়, ধীরেন্দ্রনাথ (2000). সংস্কৃত সাহিত্যের ইতিহাস. পশ্চিমবঙ্গ রাজ্য পুস্তক পর্ষৎ (West Bengal State Book Board).
Bandyopādhyāya, Dhīrendranātha (2000). Saṁskṛta Sāhityera Itihāsa. Paścimabaṅga Rājya Pustaka Parṣat (West Bengal State Book Board).
4. <https://github.com/virtualvinodh/aksharamukha/issues/173>
5. <https://fr.wikipedia.org/wiki/%E0%A6%AD%E0%A6%BC>
6. সান্যাল, দুর্গাচন্দ্র (সন ১৩১৬). ভাষা বিজ্ঞান নামক বাঙ্গালা ভাষার ব্যাকরণ. হিতবাদী লাইব্রেরী.
Sānyāla, Durgācandra (sana 1316). Bhāṣā Bijñāna Nāmaka Bāṅgālā Bhāṣāra Byākaraṇa. Hitabādī Lāibrerī.
7. No author (1997). বাংলা বানানবিধি: পশ্চিমবঙ্গ বাংলা আকাদেমি গৃহীত. পশ্চিমবঙ্গ বাংলা আকাদেমি. কলকাতা
No author (1997). Bāṅlā bānānabidhi: paścimabaṅga bāṅlā ākādēmi grhīta. Paścimabaṅga Bāṅlā Ākādēmi. Kalakātā.
8. চৌধুরী, জামিল (2016). বাংলা একাডেমি আধুনিক বাংলা অভিধান. বাংলা একাডেমি. ঢাকা.
Choudhury, Jamil (2016). Bāṅlā ekāḍēmi ādhunika bāṅlā abhidhāna (Bangla Academy Modern Bangla Dictionary). Bangla Academy. Dhaka.
(https://xeroxtree.com/pdf/adhunik_bangla_ovidhan.pdf)
9. https://sarkarverse.org/wiki/Shabda_Cayanika

**ISO/IEC JTC 1/SC 2/WG 2
PROPOSAL SUMMARY FORM TO ACCOMPANY SUBMISSIONS
FOR ADDITIONS TO THE REPERTOIRE OF ISO/IEC 10646⁵.**

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

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

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

See also <http://std.dkuug.dk/JTC1/SC2/WG2/docs/roadmaps.html> for latest Roadmaps.

A. Administrative

1. Title:	<i>Revised Proposal to Encode Alternate BA for the Bengali Language</i>
2. Requester's name:	<i>Vinodh Rajan & Deepr Chakraborty</i>
3. Requester type (Member body/Liaison/Individual contribution):	<i>Individual</i>
4. Submission date:	<i>17th Dec 2022</i>
5. Requester's reference (if applicable):	
6. Choose one of the following:	
This is a complete proposal:	<i>Y</i>
(or) More information will be provided later:	

B. Technical – General

1. Choose one of the following:			
a. This proposal is for a new script (set of characters):			
Proposed name of script:			
b. The proposal is for addition of character(s) to an existing block:	<i>Y</i>		
Name of the existing block:	<i>Bengali</i>		
2. Number of characters in proposal:			
3. Proposed category (select one from below - see section 2.2 of P&P document):			
A-Contemporary	B.1-Specialized (small collection)	B.2-Specialized (large collection)	<i>Y</i>
C-Major extinct	D-Attested extinct	E-Minor extinct	
F-Archaic Hieroglyphic or Ideographic	G-Obscure or questionable usage symbols		
4. Is a repertoire including character names provided?	<i>Y</i>		
a. If YES, are the names in accordance with the "character naming guidelines" in Annex L of P&P document?	<i>Y</i>		
b. Are the character shapes attached in a legible form suitable for review?	<i>Y</i>		
5. Fonts related:			
a. Who will provide the appropriate computerized font to the Project Editor of 10646 for publishing the standard?	<i>Vinodh Rajan</i>		
b. Identify the party granting a license for use of the font by the editors (include address, e-mail, ftp-site, etc.):	<i>Vinodh Rajan vinodh@virtualvinodh.com</i>		
6. References:			
a. Are references (to other character sets, dictionaries, descriptive texts etc.) provided?	<i>Y</i>		
b. Are published examples of use (such as samples from newspapers, magazines, or other sources) of proposed characters attached?	<i>Y</i>		
7. Special encoding issues:			
Does the proposal address other aspects of character data processing (if applicable) such as input, presentation, sorting, searching, indexing, transliteration etc. (if yes please enclose information)?	<i>Y</i>		
	<i>Collation</i>		

8. Additional Information:

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 <http://www.unicode.org> for such information on other scripts. Also see Unicode Character Database (<http://www.unicode.org/reports/tr44/>) and associated Unicode Technical Reports for information needed for consideration by the Unicode Technical Committee for inclusion in the Unicode Standard.

⁵ 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

1. Has this proposal for addition of character(s) been submitted before? If YES explain	<i>This is a revised version of L2/22-268</i>	Y
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? If YES, available relevant documents:	<i>The co-author, Deepto Chakraborty, is a native user and manuscript scholar</i>	Y
3. Information on the user community for the proposed characters (for example: size, demographics, information technology use, or publishing use) is included? Reference:	<i>See proposal</i>	Y
4. The context of use for the proposed characters (type of use; common or rare) Reference:	<i>See proposal</i>	Rare
5. Are the proposed characters in current use by the user community? If YES, where? Reference:	<i>Some Sanskrit publishers using the Bengali orthography</i>	Y
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:	<i>Bengali is in the BMP</i>	Y Y
7. Should the proposed characters be kept together in a contiguous range (rather than being scattered)?		Y
8. Can any of the proposed characters be considered a presentation form of an existing character or character sequence? If YES, is a rationale for its inclusion provided? If YES, reference:	<i>See proposal</i>	Y Y
9. Can any of the proposed characters be encoded using a composed character sequence of either existing characters or other proposed characters? If YES, is a rationale for its inclusion provided? If YES, reference:		N
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:	<i>See Proposal</i>	Y Y
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: Is a list of composite sequences and their corresponding glyph images (graphic symbols) provided? If YES, reference:		N
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)		N
13. Does the proposal contain any Ideographic compatibility characters? If YES, are the equivalent corresponding unified ideographic characters identified? If YES, reference:		N