# PROPOSAL TO ENCODE FOUR PEGON CHARACTERS

Rikza F. Sh.

رِكْزَا ف. ص.

rikzafsh@gmail.com

22 May 2022

## 1. Introduction

Pegon is a modified Arabic alphabet used to write Modern Javanese, Sundanese, and Madurese language. The word "Pegon" originated from the Javanese word *pego*, which means "deviate". Because the Javanese language written in Arabic is considered something unusual. Pegon itself is used among Muslims, who live from religious education in *Pondok Pesantren* (Islamic boarding schools). The Pegon script is currently supported in Unicode. However, in the process of researching some documents, the author has discovered a number of characters currently not encoded in Unicode.

## 2. Proposed characters to be encoded

### 1. Arabic Vowel Sign Pepet

This vowel sign is used to represent vowel /ə/ in Pegon script, and commonly known as "pepet", which is cognate with JAVANESE VOWEL SIGN PEPET (U+A9BC). This vowel sign has similar shape with ARABIC MADDAH ABOVE (U+0653), so the users of Pegon script are commonly using it to write Pepet in the digital text. However, in the process of researching some documents, the author has discovered a contrast shape use between VOWEL SIGN PEPET and MADDAH ABOVE. The author thinks, since those character have the different shapes and functions, the VOWEL SIGN PEPET should be encoded separately and should not be unified with MADDAH ABOVE.

| Codepoint | Glyph | Name |
|-----------|-------|------|
| U+088F | | ARABIC VOWEL SIGN PEPET |
| U+0653 | | ARABIC MADDAH ABOVE |



Figure 1. Showing different shapes between VOWEL SIGN PEPET (left) and MADDAH ABOVE (right). The VOWEL SIGN PEPET is consist of three straight strokes and MADDAH ABOVE is consist of a curved stroke and a straight stroke.

## 2. Arabic Letter Dal with Two Dots Vertically Below

This letter is used in some Pegon manuscript to represent voiced retroflex stop /ɖ/. The phoneme /ɖ/ in Pegon script is usually represented with letter dal with one or three dots below. But the author has discovered some manuscript that used letter dal with two dots vertically below, which need to be encoded.

| Codepoint | Isolated | Final | Medial | Initial | Name |
|---|---|---|---|---|---|
| U+10EC2 | ڍ | ڍ | | | ARABIC LETTER DAL WITH TWO DOTS VERTICALLY BELOW |
| U+068A | ڊ | ڊ | | | ARABIC LETTER DAL WITH DOT BELOW |
| U+08AE | ڋ | ڋ | | | ARABIC LETTER DAL WITH THREE DOTS BELOW |

## 3. Arabic Letter Tah with Two Dots Vertically Below

This letter is used in some Pegon manuscript to represent voiceless retroflex stop /ʈ/. The phoneme /ʈ/ in Pegon script is usually represented with letter tah with one or three dots below. But the author has discovered some manuscript that used letter tah with two dots vertically below, which need to be encoded.

| Codepoint | Isolated | Final | Medial | Initial | Name |
|---|---|---|---|---|---|
| U+10EC3 | ط | ط | ط | ط | ARABIC LETTER TAH WITH TWO DOTS VERTICALLY BELOW |
| U+068A | ط | ط | ط | ط | ARABIC LETTER TAH WITH DOT BELOW |
| U+08AE | ط | ط | ط | ط | ARABIC LETTER TAH WITH THREE DOTS BELOW |

## 4. Arabic Letter Kaf with Two Dots Vertically Below

This letter is used in some Pegon manuscript to represent voiced velar stop /g/. The phoneme /g/ in Pegon script is usually represented with letter kaf or keheh with one, two or three dots below. The author has discovered some manuscript that used letter kaf with two dots vertically below, which need to be encoded. The LETTER KAF WITH TWO DOTS VERTICALLY BELOW has similar shape with letter KEHEH WITH TWO DOTS VERTICALLY BELOW in initial in medial form, but they have different shapes in final and isolated form.

| Codepoint | Isolated | Final | Medial | Initial | Name |
|---|---|---|---|---|---|
| U+10EC4 | كِ | لكِ | كِ | كِ | ARABIC LETTER KAF WITH TWO DOTS VERTICALLY BELOW |
| U+068A | كِ | كِ | كِ | كِ | ARABIC LETTER KEHEH WITH TWO DOTS VERTICALLY BELOW |

## 3. Correction for U+06AE

As Unicode version 14.0, the Arabic Letter Kaf with three dots below (06AE) has annotation like this:



The author has found a mistake in second point that annotate U+06AE as alternative for U+068A in Pegon orthography. Whereas actually, U+06AE is alternative character for ARABIC LETTER KAF WITH DOT BELOW (U+08B4) in Pegon orthography to represent phoneme /g/. And U+068A has alternative character ARABIC LETTER DAL WITH THREE DOTS BELOW (U+08AE).
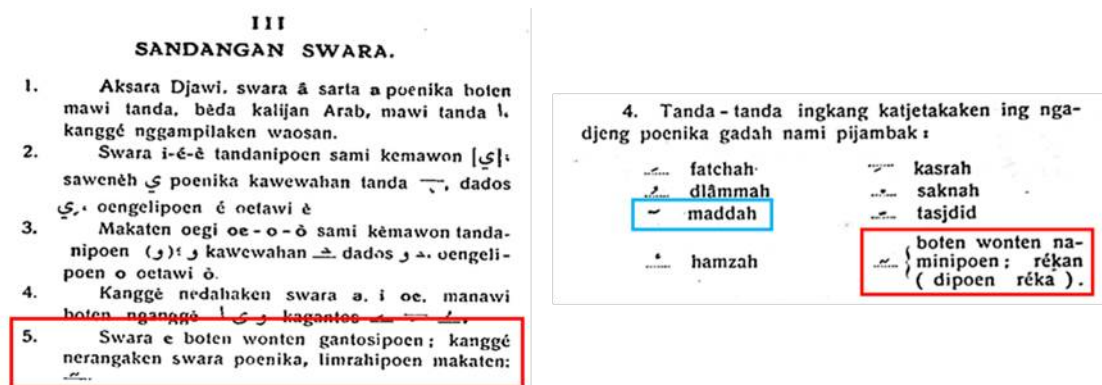
## 4. Attestation

### Arabic Vowel Sign pepet



Figure 2. Showing contrast between VOWEL SIGN PEPET and MADDAH ABOVE in "Patokanipoen Basa Djawi Kaserat Aksara 'Arab (Pegon)" (1933) p.6 and 19.
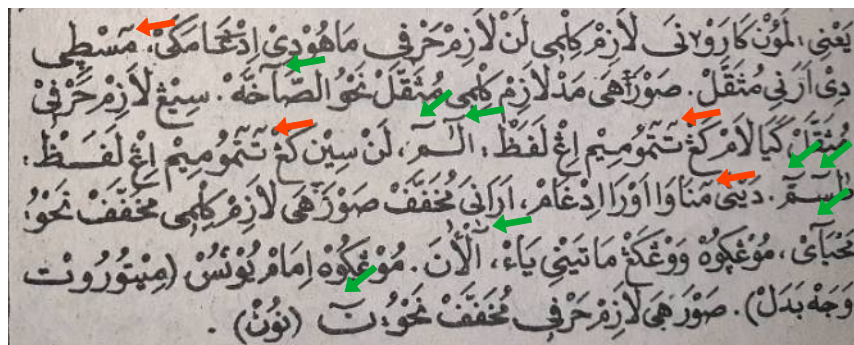


Figure 3. Showing contrast between VOWEL SIGN PEPET (marked red) and MADDAH ABOVE (marked green) in تحفة الأطفال p.23
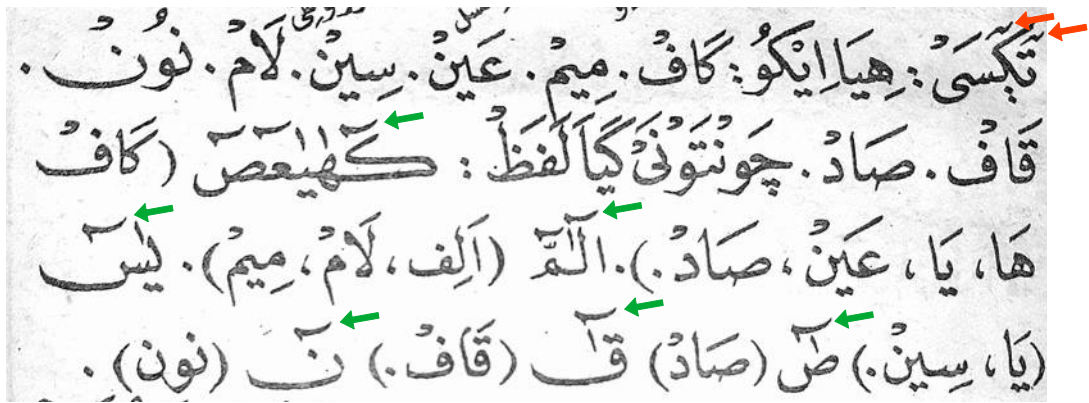
Figure 4. Showing contrast between VOWEL SIGN PEPET (marked red) and MADDA WAAJIB (U+089C) (marked green) in شفاء الجنان p.26



Figure 5. Showing table containing list of vowel sign used in Pegon, including VOWEL SIGN PEPET. From the "Al-Itqan Pêdoman Baca Tulis Arab Pegon" p.8

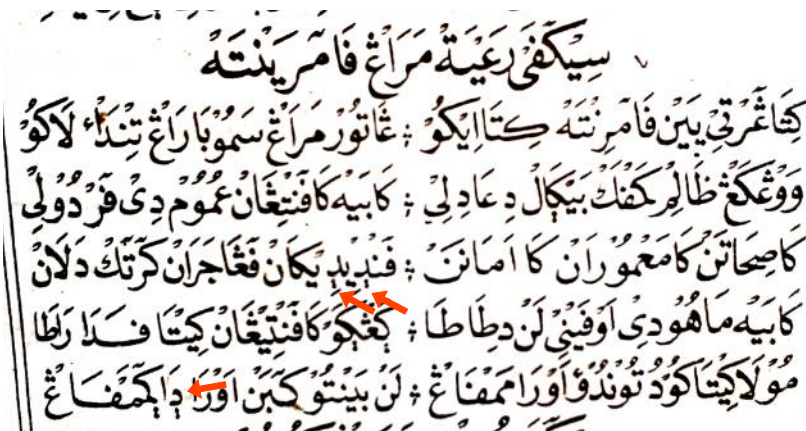**Arabic Letter Dal with Two Dots Vertically Below**



Figure 6. Showing LETTER DAL WITH TWO DOTS VERTICALLY BELOW in مترا سجاتي p.3
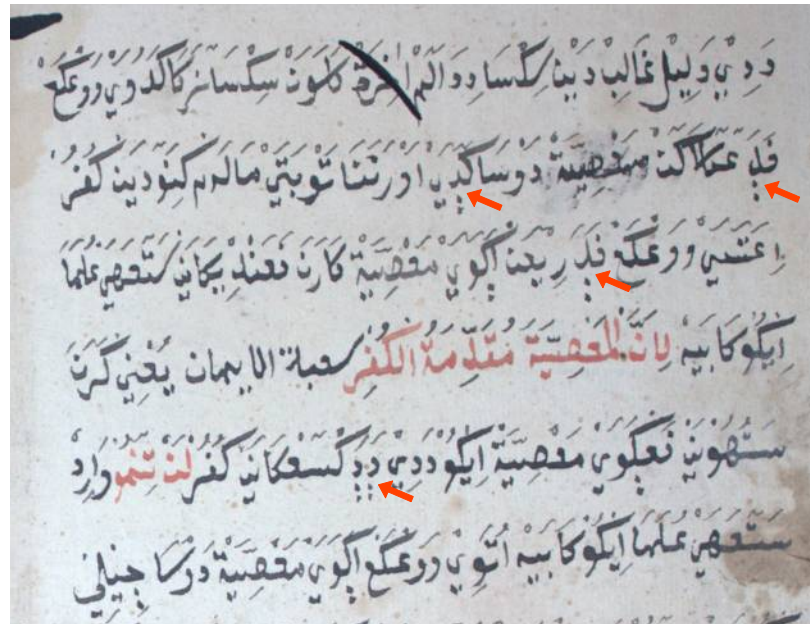
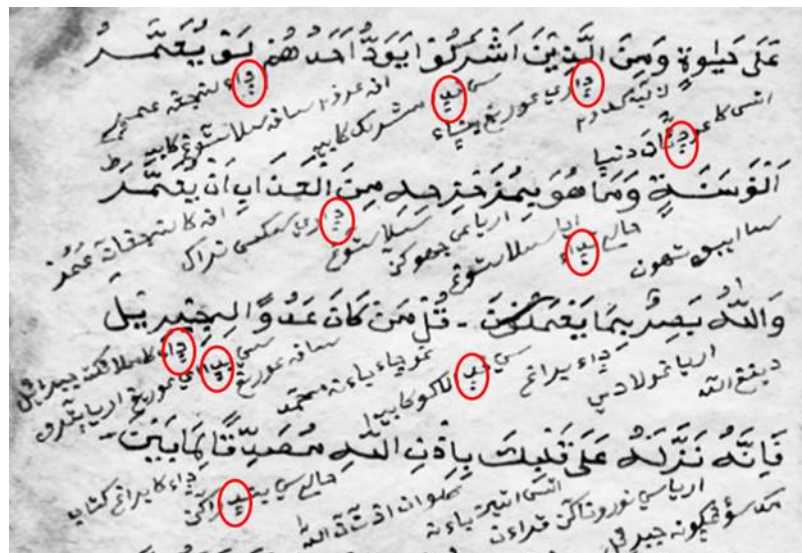Figure 7. Showing LETTER DAL WITH TWO DOTS VERTICALLY BELOW in ترهيب Tasawuf p.2



Figure 8. Mushaf SB 3 (C/TD_LK03/MM) from West Sumatra, Indonesia, showing LETTER DAL WITH TWO DOTS VERTICALLY BELOW on isolated and final form.

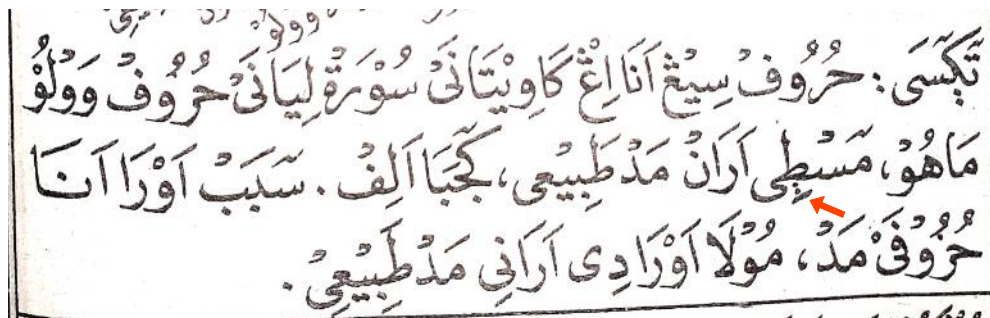**Arabic Letter Tah with Two Dots Vertically Below**



Figure 9. Showing LETTER TAH WITH TWO DOTS VERTICALLY BELOW in شفاء الجنان p.27

Figure 10. Showing LETTER TAH WITH TWO DOTS VERTICALLY BELOW in ترهيب Tasawuf p.47

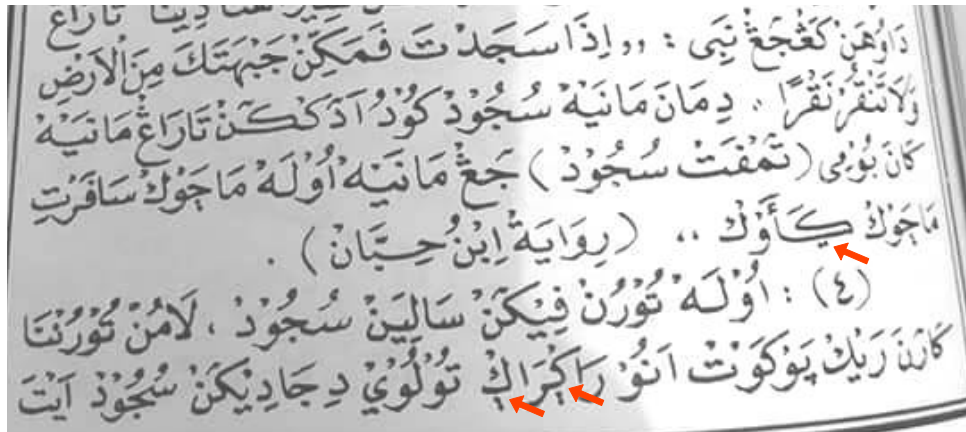**Arabic Letter Kaf with Two Dots Vertically Below**



Figure 11. Showing initial and isolated forms of LETTER KAF WITH TWO DOTS VERTICALLY BELOW in ترجمة سفسنة النجا p.147. Note that its swash version may has horizontal dot alignment instead of vertical.
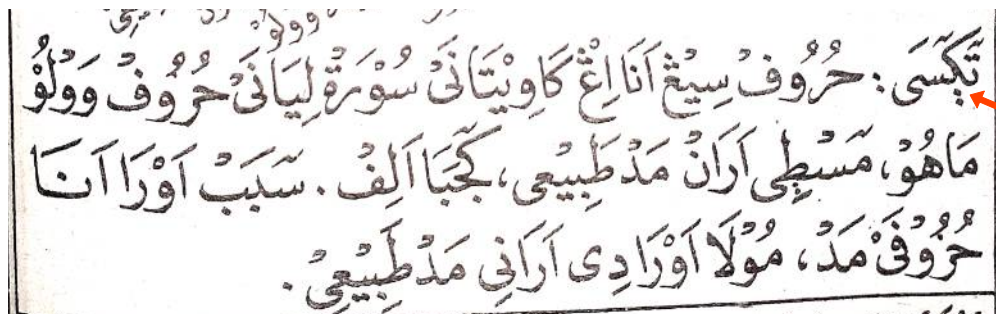


Figure 12. Showing medial form of LETTER KAF WITH TWO DOTS VERTICALLY BELOW in شفاء

الجنان p.27

6

**Arabic Letter Kaf with dot Below** and **Arabic Letter Kaf with Three Dots Below**



Figure 13. Showing ARABIC LETTER KAF WITH THREE DOTS BELOW (06AE). From "Patokanipoen Basa Djawi Kaserat Aksara `Arab (Pegon)" p.4.



Figure 14. Showing ARABIC LETTER KAF WITH DOT BELOW (08B4). From "Al-Itqan Pêdoman Baca Tulis Arab Pegon" p.7

## 5. Character Data

### Codepoints

L2/21-181 note that "… the BMP is getting filled up, … we need to reserve some space for some [potential] characters …, and put the historical and technical characters in the SMP. … This leaves 7 empty slots in the Arabic Extended-B block for important letters and combining marks."

Based on the statement above, The author decided to propose the VOWEL SIGN PEPET in codepoint U+088F in Arabic Extended-B block because that character is very commonly used in old documents and modern writing. And the letter DAL, TAH, and KAF with TWO DOTS VERTICALLY BELOW in Arabic Extended-C block because those letters are considered as variants of other letters and rarely used. Those letters are proposed in codepoint U+10EC2..10EC4 because U+10EC0..10EC1 was reserved for pending characters as mentioned in L2/21-181.

### Unicode character properties
```
088F; ARABIC VOWEL SIGN PEPET;Mn;0;NSM;;;;;N;;;;;
10EC2; ARABIC LETTER DAL WITH TWO DOTS VERTICALLY BELOW;Lo;0;AL;;;;;N;;;;;
10EC3; ARABIC LETTER TAH WITH TWO DOTS VERTICALLY BELOW;Lo;0;AL;;;;;N;;;;;
10EC4; ARABIC LETTER KAF WITH TWO DOTS VERTICALLY BELOW;Lo;0;AL;;;;;N;;;;;
```

### Joining type and group for ArabicShaping.txt
```
10EC2; ARABIC LETTER DAL WITH TWO DOTS VERTICALLY BELOW; D; DAL
10EC3; ARABIC LETTER TAH WITH TWO DOTS VERTICALLY BELOW; D; TAH
10EC4; ARABIC LETTER KAF WITH TWO DOTS VERTICALLY BELOW; D; KAF
```

7

**Annotations**

The followwing anotation should be updated in Namelist.txt

```
06AE; ARABIC LETTER KAF WITH THREE DOTS BELOW
        * Berber, early Persian
        * Pegon alternative for 08B4
```

# 5. References

Abu M. Ghithrof Danil-Barr (20??). Al-Itqan Pêdoman Baca Tulis Arab Pegon.

Nitisastro (1933). Patokanipoen Basa Djawi Kaserat Aksara `Arab (Pegon).

احمد مطهر بن عبد الرحمن المراقى. نيل الأنفال فى ترجمة **تحفة الأطفال**

احمد مطهر بن عبد الرحمن المراقى. **شفاء الجنان** فى ترجمة هداية الصبيان

كياهى بشرى مصطفى. **مترا سجاتى**

الحاج محمد عيسى بن محمد انوار. ترجمة **سفينة النجا**, دترجمهكن كان باس سوندا

**ترهيب** Tasawuf, a Cirebon manuscript

Mushaf SB 3 (C/TD_LK03/MM) from West Sumatra, Indonesia

**ISO/IEC JTC 1/SC 2/WG 2**
**PROPOSAL SUMMARY FORM TO ACCOMPANY SUBMISSIONS**
**FOR ADDITIONS TO THE REPERTOIRE OF ISO/IEC 10646**.[1]
**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:** *Proposal to Encode Four Pegon Characters*
2. Requester's name: *Rikza F. Sh.*
3. Requester type (Member body/Liaison/Individual contribution): *Individual contribution*
4. Submission date: *22 May 2022*
5. Requester's reference (if applicable):
6. Choose one of the following:
      This is a complete proposal: *YES*
      (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: *YES*
            Name of the existing block: *Arabic Extended-B and Arabic Extended-C*
2. Number of characters in proposal: *4*
3. Proposed category (select one from below - see section 2.2 of P&P document):
   A-Contemporary  *X*   B.1-Specialized (small collection) ___   B.2-Specialized (large collection) ___
   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?
      a. If YES, are the names in accordance with the "character naming guidelines"
         in Annex L of P&P document?
      b. Are the character shapes attached in a legible form suitable for review?
5. Fonts related:
      a. Who will provide the appropriate computerized font to the Project Editor of 10646 for publishing the standard?
      *Rikza F. Sh.*
      b. Identify the party granting a license for use of the font by the editors (include address, e-mail, ftp-site, etc.):
      *rikzafsh@gmail.com*
6. References:
      a. Are references (to other character sets, dictionaries, descriptive texts etc.) provided? *YES*
      b. Are published examples of use (such as samples from newspapers, magazines, or other sources)
      of proposed characters attached? *YES*
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)? *YES*
      *See proposal*

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.

---

1. 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? _ _ _NO_ _ _ _
   If YES explain _____
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.)? _ _ _YES_ _ _
   If YES, with whom? _____
   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? _ _ _YES_ _ _
   Reference: _____ *See proposal* _____
4. The context of use for the proposed characters (type of use; common or rare) _ _ *common* _ _
   Reference: _ _ _ _ *Used in many traditional Pegon manuscripts and their modern transcription* _ _
5. Are the proposed characters in current use by the user community? _ _ _YES_ _ _
   If YES, where?  Reference: _ _ _ *Modern transcription of traditional manuscript and modern writing*
6. After giving due considerations to the principles in the P&P document must the proposed characters be entirely
   in the BMP? *NO, only one character should be encoded in BMP*
   _ _ _ _ _ _ _
   If YES, is a rationale provided? _ _ _ _ _ _ _
   If YES, reference: _____
7. Should the proposed characters be kept together in a contiguous range (rather than being scattered)? _ _ NO _ _
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? _ _ YES _ _
    If YES, is a rationale for its inclusion provided? _ _ YES _ _
    If YES, reference: _____ *See proposal* _____
11. Does the proposal include use of combining characters and/or use of composite sequences? _ _ YES _ _
    If YES, is a rationale for such use provided? _ _ YES _ _
    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? _ _ NO _ _
    _ _ _ _ _ _If YES, describe in detail (include attachment if necessary)_ _ _ _ _ _ _ _ _ _ _ _ _ _
    _____
    _____
13. Does the proposal contain any Ideographic compatibility characters? _ _ _NO_ _ _ _
    If YES, are the equivalent corresponding unified ideographic characters identified? _ _ _ _ _ _ _
    If YES, reference: _____