#### ISO/IEC JTC1/SC2/WG2 N4513

Title: Comments on Nushu in WG2 N4484, /IEC 10646:2014/PDAM1

Status: Individual Contribution

Source: suzuki toshiya, Hiroshima University, Japan

After the interrupt from 2009 to 2012 for the standardization of Nushu, I appreciate the efforts by China NB, SEI and the experts from US, UK, Ireland, Finland and TCA. Japanese expert contribution by Endo Orie, WG2 N3705, was almost ignored in the reopened discussion, but the efforts by Deborah Anderson and Michael Everson improved some questionable parts greatly. However, Nushu charset in PDAM text (WG2 N4484) still needs more improvement, or clarification. This document lists the questionable parts left in PDAM1 code chart.

In general, following points are expected to be discussed for further improvement.

• Stability of the statistical investigation

The idea choosing a referential glyph by the statistical frequency is reasonable, but the current statistical results seem to be unstable to define the normative and uncancellable attributes (like character names). Therefore, although the charset in Nushu Duben (ISBN 978-7-5438-5282-2) could be useful to define the basic charset, it would not be appropriate to define the normative and uncancellable attributes (like character names).

Stroke counts and phonetic value in the character names

Some glyphs in PDAM text have different stroke counts from the section where they are placed. In addition, it seems that sometimes Nushu Duben prioritizes source Hanzi than the phonetic value. Considering that the separated indexing in Nushu Duben character is the most important criteria of the charset, and some Nushu characters are difficult to determine their referential phonetic value by the statistical result, the inclusion of the phonetic value in the character names should be reconsidered.

Considering 2 points in above, I ask Nushu expert group or China NB for a referential document including *mutually associated* data to review the proposed Nushu code chart; the proposed glyph shape, raw glyph shape(s), its semantics, its variants and the statistical results for them. The documents previously submitted from China NB included most of them, but they are not mutually associated. Taking WG2 N3598 as an example, the proposed glyph shapes are presented in the proposed code chart (page

1

9-20), the raw glyph shapes are presented in the comparison of 3 dictionaries (Appendix C, page 31-49), the variants and frequencies are presented in the source Hanzi table (page 72-91). They are no associated via shared identifier. As Endo Orie commented in WG2 N3705 and Deborah Anderson commented in WG2 N4442 and N4451, when a reviewer finds a discrepancy, it is difficult to guess the difference is acceptable one, mistakenly differentiated one, or anything else.

For example, Appendix F, the section 女書交際用字研究 (study of the Nushu glyph for interchange) picks the glyph variants to mean guest (客);

#### 2 个体用字差异

在我们所知的女书老人中,方块汉字水平最低的是阳焕宜,因此她用字最少,而且出现方块汉字也最少。例如:

Figure 1: Variants of "guest" glyph (WG2 N3598, Appendix section 3-2-2, p.98)

But I could find none of them in the proposed charset in PDAM, or the frequency table in WG2 N3598. I'm not sure why, they were unified with other characters? Or, their frequencies were too low to be coded in the basic charset? The changes since the last official proposal from China NB should be summarized.

# Stability of the statistical investigation submitted to WG2

As the proposals from China NB states, the proposed charset is not designed to be the complete set, but the basic representative glyph set. Anyway, the encoding architecture would be expected to be stabilized at the first basic charset, to avoid the ad-Hoc changing of the encoding architecture in future extensions. In original China NB proposal, the character identity was decided by the statistical investigation (e.g. the last China NB proposal, WG2 N4341, chapter I, section 2 "the judgment and collation of Nyushu basic and variant characters"). In PDAM1, although the character identity discussion is simplified as the indexing characters in Nushu Duben (according to WG2 N4461), the identity of the indexing character in the book is supposed to be based on similar statistical investigation. But checking the content carefully, the stability of the statistical investigation is questionable.

One reason would be that the preliminary classification of the glyphs before counting their frequencies. It would be described in following subsections.

Another reason would be that the number of the authors of the Nushu document collection is not so large; according to WG2 N3598 chapter II, section 1 "the collection and collation of Nushu characters", the materials by 4 authors (Gao Yinxian, Yi Nianhua, Yang Huanyi and He Yanxin) document is 85%. Therefore, the statistical

result may reflect individual preferences. Originally, Nushu Yongzi Bijiao presented per-author statistics. Choosing a general purpose (or most neutral) referential glyph from per-user referential glyph would be more stable approach.

# 1.1. Statistical summarization is questionable

As Endo Orie had pointed in WG2 3705 (2009), the statistical investigation result seems to be inappropriately presented, and questionable for its reliability in as-is manner. For example, the process how the representative glyph for Y42 is determined was described in WG2 N4341 section 2-B.

#### B. Prior to basic characters according to their frequency

That is to say, within the characters of the same character cell, we chose the character with the highest frequency. The character with the highest frequency is called basic representative and others as allograph. The final result of the study and the statistics can be referred to 'Nyushu basic characters and their origin verification' (published in Nyushu character-comparison, 2006).

Take the following table to explain:

y <sup>42</sup>	如余餘 儒虞娱	<b>於</b> 如105	<b>№</b> 如 131 余 10	<b>炒</b> 如 213 余 2	<b>⋄</b> 如 75 余 3	<b>於</b> 如141余3
		<b>对</b> 如10	J.	<b>炒</b> 如 3		<b>省</b> 如 11
		<b>%</b> 如 8	₩ 如 2 儒 1	<b>沙</b> 如 1	<b>沁</b> 如 11	
					<b>省</b> 如 4	

There are three main graphs for the sound of y y y y y y y y which has the same pronunciation and the meaning and has general similar structures. The first one has a rather high frequency. So y is the representative glyph.

Figure 2: Description of the frequency statistics (WG2 N4341, section 2-B)

It notes that there are 3 major glyphs frequently used (590 times), (89 times), (12 times). But how the frequency could be calculated from the table? In the table, the 1st glyph frequency might be counted as  $y_{105} + y_{131} + y_{131$ 

reliable even if it is in published material. There are some glyphs whose frequency is less than 10, using the materials including such errors as an authorized reference would be inappropriate to stabilize the encoding architecture.

# 1.2. Statistical investigation seems to be under development

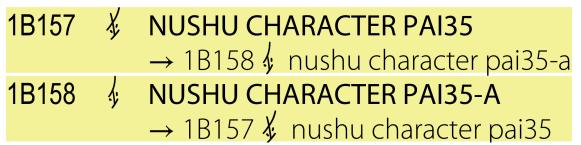


Figure 3: Variants of PAI35 in PDAM (WG2 N4484)

	WG2 3337 (200	07)	WG2 3598 (2009)		
U+1B157 glyph	<b>美</b> (本) pai <sup>35</sup> 才	<b>*</b> 176 <b>*</b> (本 <b>\$</b>	pai <sup>35</sup> 本 <sup>34</sup> )piu <sup>35</sup> 表 <sup>16</sup> t¢ioŋ <sup>35</sup> 整 <sup>5</sup> 颈 <sup>2</sup>		
U+1B158 glyph	\$ \$P\$ (正/tcyn <sup>21</sup> 转 tcion <sup>21</sup> 正 tson <sup>42</sup> 渐 pau <sup>35</sup> 宝 tcion <sup>33</sup> 整 tcion <sup>33</sup> 警 tson <sup>42</sup> 惭 tsian <sup>21</sup> 帐 tcin <sup>44</sup> 肩 p'ai <sup>35</sup> 品 cyn <sup>35</sup> 选	15 保 10 5 2 2 2	pai <sup>35</sup> 本 <sup>212-178-9</sup> pou <sup>35</sup> 打 <sup>289</sup> tçioŋ <sup>21</sup> 正 <sup>151</sup> 镜 <sup>31</sup> 政 <sup>24</sup> 敬 <sup>21</sup> 竟 <sup>1</sup> tçyn <sup>21</sup> 转 <sup>139-2-0</sup> 卷 <sup>10</sup> 眷 <sup>8</sup> tçioŋ <sup>44</sup> 正 (~月) <sup>50</sup> 京 <sup>9</sup> 惊 <sup>9</sup> tsoŋ <sup>33</sup> 渐 <sup>41</sup> pau <sup>35</sup> 宝 <sup>15</sup> 保 <sup>10</sup> tçioŋ <sup>35</sup> 整 <sup>11-40-3</sup> 颈 <sup>5-2-2</sup> 景 <sup>0-2</sup> tçie <sup>35</sup> 诊 <sup>3</sup> 拯 <sup>1</sup> tsioŋ <sup>35</sup> 井 <sup>4</sup> tsoŋ <sup>42</sup> 惭 <sup>2</sup> tçiaŋ <sup>21</sup> 帐 <sup>2</sup> p'ai <sup>35</sup> 品 <sup>1</sup> piu <sup>35</sup> 表 <sup>0-2-0</sup>		

Table 1: Statistical results for PDAM PAI35 variants

There are 2 characters with same phonetic values, and similar shapes are coded separately. This pair was originally coded for different pronunciation (WG2 N4341 designed U+1B157 glyph for PAI, U+1B158 glyph for POU). According to WG2 N4472R, it was decided that the phonetic values were determined by the top entry of the phonetic values in Nushu Duben, and these 2 characters are now considered as same phonetic

value characters. However, tracking the older proposals, it seems that the assignment of U+1B157 glyph to PAI35 and of U+1B158 glyph to POU35 was reasonable.

In WG2 N3337, the top frequency of the phonetic values for U+1B158 glyph was not PAI35 but POU35. The frequency U+1B157 for PAI35 (本) was originally 176 (in WG2 N3337 (2007)) but decreased to 34 (in WG2 3598 (2009) in later. If we had discussed in 2007, the phonetic value for the character name of U+1B158 should be POU35. On the other hand, some frequency was moved from U+1B158 to U+1B157; the frequency of U+1B158 glyph for PIU35 (表) is decreased from 0-2-16 to 0-2-0, and the decreased frequency seems to be moved to U+1B157 glyph. By such changes, I have a concern that the glyph identification during the statistical investigation was under development. Is it reasonable to expect that the most frequent phonetic values in Nushu Duben are already stabilized?

# 1.3. Some of rarely used character are supposed to be daily-used

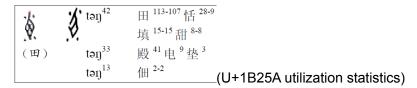
In the statistical result of WG2 N3598, there are some characters at low frequency, less than 10. Some of them are supposed to be very daily used words, like, 尿, 炭, 田, 酒, etc.

n,iu <sup>33</sup> 尿 <sup>3</sup> (尿)	<b>k</b> t'au <sup>35</sup> 讨 <sup>7</sup> (讨)	<b>½</b> huom <sup>35</sup> 反 <sup>5</sup>
<b>大</b> t'uow <sup>21</sup> 炭 <sup>3</sup> (炭)	k'aŋ³⁵ 孔 <sup>8</sup>	təŋ <sup>42</sup> 田 <sup>2</sup> təŋ <sup>33</sup> 垫 <sup>2</sup>
p'ø <sup>35</sup> 派 <sup>4</sup> (派)	ts'i <sup>5</sup> 切 <sup>5</sup> (切)	ts'ø <sup>44</sup> 差 <sup>3</sup> 猜 <sup>1</sup> (在/差)
ts'i <sup>5</sup> 妾 <sup>3</sup> (妾)	tç'iou <sup>21</sup> 臭 <sup>2</sup>	vuiə <sup>35</sup> 風 <sup>3</sup>
<b>秋</b> pui <sup>5</sup> 拨 <sup>4</sup> 北 <sup>3</sup>	pʻai <sup>35</sup> 品 <sup>4</sup>	lai <sup>35</sup>
t'uouu <sup>44</sup> 通 <sup>5</sup> (滩/炭)t 'aŋ <sup>44</sup> 滩 <sup>3</sup>	teiŋ <sup>44</sup> 沾 <sup>4</sup> teiŋ <sup>42</sup> 缠 <sup>1</sup> (缠/展/沾) teiŋ <sup>35</sup> 展 <sup>1</sup> 捡 <sup>1</sup>	ts'ie <sup>21</sup> 退 <sup>6</sup>
ti <sup>44</sup> 梯 <sup>4</sup> (梯)	<b>h</b> au <sup>33</sup> 号 <sup>4</sup>	tsiou <sup>33</sup> 酒 <sup>9</sup>

Table 2: Sample of rarely used characters in statistical result (WG2 N3598)

It is reasonable to consider a situation that some characters are low frequency because there is different character with same phonetic values. However, checking the variant for the character that are supposed for daily used words, the situation is different.

- a) U+1B150 (尿, NJIU33, urine) has no different character with same phonetic value in WG2 N3598 nor PDAM. So the word of "尿" is used only 3 times in 220000 character documents?
- b) U+1B194 (炭, THUOW21, charcoal) has no different character with same phonetic value in WG2 N3598 nor PDAM. So the word of "炭" is used only 3 times in 220000 character documents?
- c) U+1B1A0 (⊞, TENG42-A, rice field) has a character with same phonetic value (U+1B25A, TENG42), but its source Hanzi is again ⊞. What is the rational to handle U+1B1A0 separately? In addition, the phonetic value TENG42 is appropriate? The frequency of TENG33 for U+1B1A0 is same with TENG42.



d) U+1B232 (酒, TSIOU33-A, wine) has a character with same phonetic value U+1B1FB (袖/坤, TSIOU33). However, the usage of TSIOU33 does not include the usage of the glyph to mean "酒". Does it mean the investigated Nushu document collection rarely use the word "酒"?

From the cases a), b), d), it is supposed that the investigated Nushu document collection might be biased. From the case c), it is supposed that the distinction of U+1B1A0 and U+1B25A was unexpectedly introduced by the stroke number counting for unstable glyph shapes.

1.4. Some glyph shapes are not the most frequent glyph in the statistical result According to the statistical result for HANG44, the glyph shape in older China NB proposals and PDAM is not the most frequent shape in the statistics result in WG2 N3598. It seems that the proposed glyph shape is same with the first proposal on 2007; it was the most frequent glyph shape at that time (but the phonetic value was different, it should be HANG21). In later, the existence of the shape difference was recognized (in WG2 N3337, only one shape was counted, and the frequency was 32. in WG2 3598, the frequency was divided into to 2 shapes; 15 + 17 = 32), and the most frequent glyph shape and phonetic value was also changed. However, the proposed glyph shape was not changed.

# 1B238 M NUSHU CHARACTER HANG44

WG2 N3337 (2007)		WG2 N3598 (2009)		
haŋ²¹ haŋ⁴⁴ haŋ³³ haŋ³³ haŋ¹³ k'aŋ⁴⁴	汉 <sup>32</sup> 欢 <sup>15</sup> 唤 <sup>15</sup> 换 <sup>14</sup> 焕 <sup>11</sup> 汗 <sup>7</sup> 唤 <sup>6</sup> 翰 <sup>6</sup> 早 <sup>13</sup> 糠 <sup>1</sup>	(欢)	haŋ <sup>44</sup> haŋ <sup>21</sup> haŋ <sup>33</sup> haŋ <sup>13</sup> k'aŋ <sup>44</sup>	欢 <sup>305-164</sup> 荒 <sup>5-2</sup> 汉 <sup>17-15</sup> 换 <sup>14</sup> 焕 <sup>11</sup> 汗 <sup>7</sup> 唤 <sup>6</sup> 翰 <sup>6</sup> 早 <sup>13</sup> 糠 <sup>1</sup>

Table 3: Historical development of U+1B238 glyph frequency

There are 2 glyphs with different phonetic values and quite similar shapes. U+1B193 is classified in the 6-stroke characters, and U+1B1B2 is in the 7-stroke characters. In the earliest proposal, it seems that these glyphs are unified (see the statistic table in WG2 N3337). In later statistics, the shape difference might be found (thus, LEW44 glyph is classified in the 6-stroke, and OE44 glyph is in the 7-stroke) and counted the frequencies separately. OE44 glyph in WG2 N3337 was different from LEW44 glyph. However, in PDAM1 code chart, OE44 glyph is still similar to LEW44 glyph. If U+1B1B2 glyph shape is preferred, the stroke count should be reconsidered. If WG2 N3598 OE44 glyph is preferred, the glyph should be fixed.

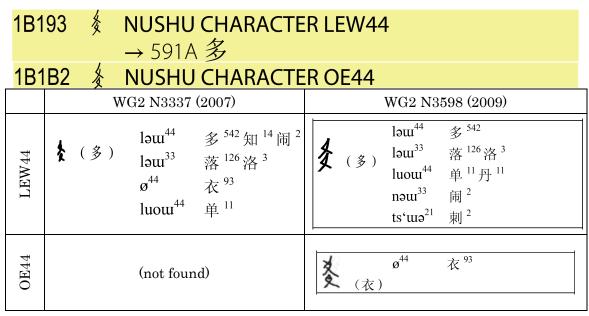


Table 4: Historical development of LEW44 and OE44 frequency

The glyph looking like U+1B105 for POE5 was not found in older China NB proposals (furthermore, it is not listed as a variant), although China NB always proposed the glyph for POE5. Also the existence of the glyph looking like U+1B105 was found in Endo Orie's He Yanxin glyph list in 中国女文字研究 (ISBN 4-625-48300-X) p.292-332. It is supposed some editorial errors exist in the statistical result table.

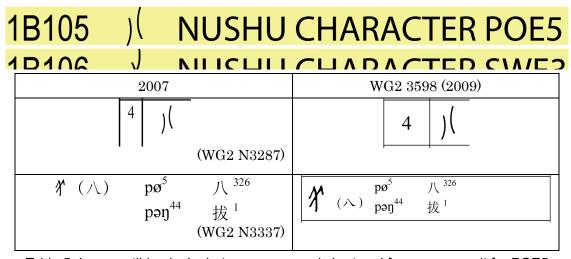


Table 5: Incompatible glyphs between proposed chart and frequency result for POE5

## 2. Stroke counts and phonetic value

The Nushu is recognized as the syllabic writing system, thus each character name includes its phonetic value. However, the identity of the character is not determined by its phonetic value, and the source Hanzi is supposed to be the main factor to consider the separated encoding.

## 2.1. How the indexing item in the statistical result is allocated?

As described in the first proposal, WG2 N3287, section 3 rule a) described, the basic rule to give the identity to Nushu glyph is the recognization of the source Hanzi. Even if the phonetic value and semantics are same, the glyphs derived from different Hanzis are coded separately. For example, U+1B101:U+1B14E are separately coded although the unification of them does not cause the semantic loss.

1B101 / NUSHU CHARACTER I5 → 1B14E ∜ nushu character i5-a → 4E00 —	<b>J</b> (-)	$i^5$ — 1547 $i^{33}$ III $i^{1}$
1B14E    NUSHU CHARACTER I5-A  → 1B101 / nushu character i5	* (壹)	i <sup>5</sup> — 1064 i <sup>21</sup> 以 <sup>3</sup> i <sup>33</sup> 叶 <sup>2</sup>

Table 6: example of separation by source Hanzi difference for the semantically equivalent pair.

On the other hand, it seems that the variant with significant shape difference would not be coded separately, if the source Hanzi is same. The variants of U+1B188 (☐, MAI42) have significant shape difference, but not coded separately.

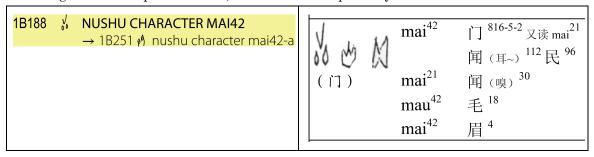


Table 7: example of unification by source Hanzi identity for the different shapes.

But there are confusing exceptions. As mentioned in above, U+1B1A0 and U+1B25A are coded separately, although their source Hanzis are same, and the shape difference is subtle than those in U+1B188 variants. If U+1BA0 shape is far different and its frequency 4 is sufficient to separate, why the second candidate glyph in U+1B25A

statistics is not coded?

U+1B1A0	U+1B25A		
təŋ <sup>42</sup> 田 <sup>2</sup> (田) təŋ <sup>33</sup> 垫 <sup>2</sup>	世報 は təŋ <sup>42</sup> 田 <sup>113-107</sup> 恬 <sup>28-9</sup> 填 <sup>15-15</sup> 甜 <sup>8-8</sup> (田) təŋ <sup>33</sup> 殿 <sup>41</sup> 电 <sup>9</sup> 垫 <sup>3</sup> təŋ <sup>13</sup> 佃 <sup>2-2</sup>		

Table 8: example of separation by shape difference for semantically equivalent and same source Hanzi.

The current status of the separation of U+1B157 and U+1B158 (see Table 2) would be similar case.

### 2.2. Stroke number is stable?

As I presented in above, I have a concern that the frequency in the statistical result is under development, and it is inappropriate to determine the uncancellable attributes of the characters. Therefore I have a concern that the number of stroke is not stabilized yet, therefore, the name of the blocks should not use the classification by the stroke number, even if the content is roughly ordered by the stroke number.

中中	mu <sup>13</sup> mu <sup>33</sup>	母 <sup>232-224</sup> 马 <sup>72-32</sup> 木 <sup>36-3</sup> 目 <sup>12</sup> 墓 <sup>4-4</sup>
(母)	mu <sup>42</sup>	麦 2-1
	mou <sup>13</sup>	亩 2-1 牡 5-2

Figure 4: Frequencies of the MU13 variants (WG2 N3598)

For example, the frequencies of 2 shapes U+1B152 MU13 (母) are almost comparable. When the most frequently used glyph is known to be the (currently) 2nd one, the number of strokes would be changed from 5 to 7.

As shown for U+1B1B2 OE 44 (衣), the representative glyph is already inconsistent with the number of strokes to be used its code point. Also the representative glyph

shape for U+1B13C CYA5 (出) is difficult to recognize its stroke number as 5 (5 is for

the 2nd frequently used glyph ??). When more documents are investigated, the most frequently used glyph could be changed. It should be noted that CYA5 glyph was not found in the code chart in WG2 N3287 and the statistical results in WG2 N3337 (also it should be noted that Endo Orie's He Yanxing glyph list published in 2002 had already included the 2nd frequently used glyph for CYA5, but the proposals in 2007 included no characters for the phonetic value CYA5).

## Summary

As presented in above, the stability of the statistical investigation of Nushu document is supposed to be still under development. Therefore, a document including the mutually associated data (code chart glyph, raw glyph, variant shapes, pronunciation, semantic Hanzi, frequency, author, etc) is needed to review the proposed code chart from the viewpoint of their stability. Also the changes from the last official proposal, WG2 N4341, should be summarized with the rationales.

Also considering that the representative glyph, its number of strokes, its phonetic value could be changed by the future investigation, the character naming convention should be reconsidered. CJK Unified Ideograph does not include any information except of its codepoint, or, Egyptian Hieroglyphs and Linear B Ideograms use the indexing number of their referential materials. Such conventions would be considerable options.

(end of document)