

New Canonical decomposition and composition processes for Hangeul

Kyongsok KIM, Professor,
Division of Computer Science and Engineering,
Busan National University, Busan 609-735, South KOREA

Internet: gimgs@HANGEUL.pnu.edu
Phone: Home/ +82-51-947-6581, Office/ +82-51-510-2292
Fax: +82-51-515-2208

Abstract

In this paper, the author proposes new canonical decomposition process (NFD: Normalization Form D) and composition processes (NFC: Normalization Form C) for Hangeul, Korean script. The new processes can handle both Modern and Old Hangeul in a consistent way. Therefore, the new processes are more general than the ones described in UAX #15, which handle only Modern Hangeul syllables.

Since NFD and NFC will play an important role in sorting and searching, Internet-related applications such as Multilingual Domain names, databases, etc., we need to establish good NFD and NFC for Hangeul before they are widely used.

Keywords :

- 1) Hangeul, Hangul, Korean
- 2) canonical, decomposition, composition
- 3) normalization, form
- 4) Multilingual domain name, Internet, sort, search

1. Introduction

In Section 1, we will see a brief introduction of Hangeul, Korean script, and terms used in this paper. In Section 2, we will review previous works about Hangeul Normalization Forms, especially four Normalization Forms described in UAX #15 [UAX15].

Then, in Section 3, we will propose new canonical decomposition process (NFD: Normalization Form D) and composition processes (NFC: Normalization Form C) for Hangeul, which are better than the previous proposals.

Conclusions are described in Section 4.

Now let's see a brief introduction of Hangeul, Korean script, and terms used in this paper.

Note 1. The spelling of "Hangeul" follows ISO TS (Technical Specification) 11941, Information and Documentation -- Transliteration of Korean script into Latin characters. First edition, 1996-12-01.

Note 2. The alphabetical order of Hangeul letters follows the one adopted in South Korea.

1) Hangeul

It is the name of Korean alphabet or Korean script. 'Han' means great and 'geul' means a script.

2) Modern Hangeul simple letters

Modern Hangeul syllable-initial, syllable-peak, and syllable-final simple letters are shown below (U+xxxx indicates UCS-2 code positions [ISO10646]).

a) Syllable-initial simple letters (14):

ㄱ U+1100, ㅋ U+1102, ㆁ U+1103, ㆅ U+1105, ㆆ U+1106, ㆇ U+1107, ㆈ U+1109, ㆉ U+110B, ㆊ U+110C, ㆋ U+110E, ㆌ U+110F, ㆍ U+1110, ㆎ U+1111, ㆏ U+1112

b) Syllable-peak simple letters (10):

ㅏ U+1161, ㅑ U+1163, ㅓ U+1165, ㅕ U+1167, ㅗ U+1169, ㅛ U+116D, ㅜ U+116E, ㅠ U+1172, ㅡ U+1173, ㅣ U+1175

c) Syllable-final simple letters (14):

ㄴ U+11A8, ㄷ U+11AB, ㄹ U+11AE, ㅁ U+11AF, ㅂ U+11B7, ㅅ U+11B8, ㅈ U+11BA, ㅊ U+11BC, ㅌ U+11BD, ㅍ U+11BE, ㅎ U+11BF, ㅇ U+11C0, ㅅ U+11C1, ㅆ U+11C2

3) Modern Hangeul complex letters

Two or three simple letters can be combined to form a complex letter.

Modern Hangeul syllable-initial, syllable-peak, and syllable-final complex letters are shown below.

a) Syllable-initial complex letters (5):

ㄲ U+1101, ㄴㅇ U+1104, ㄴㆁ U+1108, ㄴㆅ U+110A, ㄴㆈ U+110D

b) Syllable-peak complex letters (11):

ㅏㅑ U+1162, ㅑㅓ U+1164, ㅓㅕ U+1166, ㅕㅗ U+1168, ㅗㅛ U+116A, ㅛㅜ U+116B, ㅜㅟ U+116C, ㅟㅠ U+116F, ㅠㅣ U+1170, ㅏㅣ U+1171, ㅑㅣ U+1174

c) Syllable-final complex letters (13):

ㄴㅇ U+11A9, ㄴㆁ U+11AA, ㄴㆅ U+11AC, ㄴㆈ U+11AD, ㄴㆉ U+11B0, ㄴㆊ U+11B1, ㄴㆋ U+11B2, ㄴㆌ U+11B3, ㄴㆍ U+11B4, ㄴㆎ U+11B5, ㄴ㆏ U+11B7, ㄴ㆐ U+11B9, ㄴ㆑ U+11BB

4) Modern syllable-initial, syllable-peak, and syllable-final letters

Modern Hangeul syllable-initial, syllable-peak, and syllable-final letters, including both simple and complex letters, are shown below in alphabetical order of South Korea:

a) Syllable-initial letters (19):

ㄱ U+1100, ㅋ U+1101, ㆁ U+1102, ㄷ U+1103, ㅌ U+1104, ㄹ U+1105, ㅁ U+1106, ㅂ U+1107, ㅃ U+1108, ㅅ U+1109, ㅆ U+110A, ㅇ U+110B, ㅈ U+110C, ㅊ U+110D, ㅊ U+110E, ㅌ U+110F, ㅍ U+1110, ㅍ U+1111, ㅎ U+1112

b) Syllable-peak letters (21):

ㅏ U+1161, ㅑ U+1162, ㅓ U+1163, ㅕ U+1164, ㅗ U+1165, ㅛ U+1166, ㅜ U+1167, ㅠ U+1168, ㅡ U+1169, ㅟ U+116A, ㅠ U+116B, ㅡ U+116C, ㅢ U+116D, ㅣ U+116E, ㅤ U+116F, ㅥ U+1170, ㅦ U+1171, ㅧ U+1172, ㅨ U+1173, ㅩ U+1174, ㅪ U+1175

c) Syllable-final letters (27):

ㄴ U+11A8, ㄷ U+11A9, ㄹ U+11AA, ㄷ U+11AB, ㄹ U+11AC, ㄷ U+11AD, ㄷ U+11AE, ㄷ U+11AF, ㅁ U+11B0, ㅂ U+11B1, ㅃ U+11B2, ㅅ U+11B3, ㅆ U+11B4, ㅈ U+11B5, ㅊ U+11B7, ㅋ U+11B7, ㆁ U+11B8, ㄷ U+11B9, ㅅ U+11BA, ㅆ U+11BB, ㅇ U+11BC, ㅈ U+11BD, ㅊ U+11BE, ㅌ U+11BF, ㅍ U+11C0, ㅍ U+11C1, ㅎ U+11C2

5) A complete Hangeul syllable: 2-letter and 3-letter syllable

In a written text, a word is a sequence of Hangeul syllables, most of which are complete syllables. Each syllable is written in a square.

There are two types of complete Hangeul syllables:

5-a) A 2-letter syllable is composed of syllable-initial and syllable-peak (i.e., vowel) letters, where each letter can be either simple or complex.

5-b) A 3-letter syllable is composed of syllable-initial, syllable-peak (i.e., vowel), and syllable-final letters, where each letter can be either simple or complex.

6) An incomplete Hangeul syllable

An 'incomplete' syllable is a syllable which is not complete. There are four types of incomplete syllables: i) a syllable-initial letter alone, ii) a syllable-peak letter alone, iii) a syllable-final letter alone, iv) syllable-peak and -final letters alone.

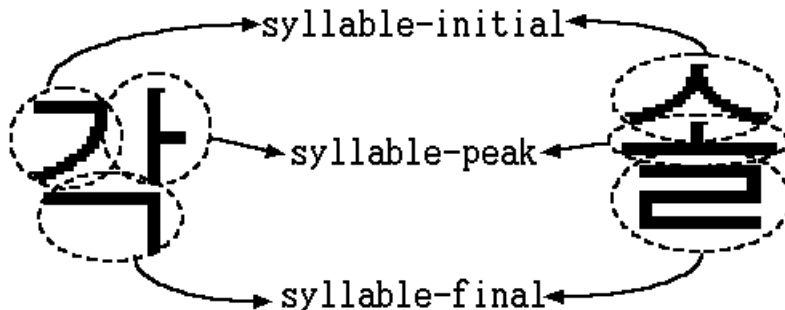
Examples: i) ㄱ (U+1100 1160), ii) ㅏ (U+115F 1161), iii) ㄴ (U+115F 11A8), iv) ㅏ (U+115F 1175 11AB)

7) Consonantal letters

In Hangeul, consonantal letters are classified into syllable-initial and -final letters depending on the physical position in a square for a syllable.

A syllable-initial letter is a consonantal letter appearing to the left and/or above a syllable-peak letter, whereas a syllable-final letter is a consonantal

letter appearing below syllable-initial and -peak letters.



8) Old Hangeul simple letters

Old Hangeul simple letters used for Old Hangeul are shown below:

- a) Old syllable-initial simple letters (3): △ U+1140, ◊ U+114C, ◡ U+1159
- b) Old syllable-peak simple letter (1): ∙ U+119E
- c) Old syllable-final simple letters (3): △ U+11EB, ◊ U+11F0, ◡ U+11F9

In addition, there are six Old syllable-initial simple letters used for foreign words.

- 8-d) U+113C, U+113E, U+114E, U+1150, U+1154, U+1155

9) Old Hangeul complex letters

There are many Old complex letters. 158 Old complex letters are included in ISO/IEC 10646. There seem fairly many (tens of ?) Old complex letters not included in ISO/IEC 10646.

10) Hangeul letters in ISO/IEC 10646-1

238 Hangeul letters (whose code positions are U+11xx) included in ISO/IEC 10646-1 are classified depending on whether they are Modern or Old, whether they are simple or complex, and whether they are syllable-intial, -peak, or -final. A summary table of letters is shown below:

Table 1. Classification of Hangeul letters included in ISO/IEC 10646-1.

	Modern			Old			Modern + Old		
	simple	complex	sub-tot.	simple	complex	sub-tot.	simple	complex	tot.
syllable-initial	14	5	19	9	62	71	23	67	90
syllable-peak	10	11	21	1	44	45	11	55	66
syllable-final	14	13	27	3	52	55	17	65	82
total	38	29	67	13	158	171	51	187	238

11) Hangeul Filler characters

Filler characters are used for representing incomplete syllables. There are two Filler characters.

syllable-initial Filler character: U+115F

syllable-peak Filler character : U+1160

Usage of Filler characters are shown below:

- ㄷ : U+1100 1160 (syllable-initial letter alone: append U+1160)
- ㅏ : U+115F 1161 (syllable-peak letter alone: prepend U+115F)
- ㅑ : U+115F 11A8 (syllable-final letter alone: prepend U+115F)

12) Old Hangeul tone marks : Bangjeom

Old Hangeul tone marks are used for representing the tones of syllables. There are two tone marks.

single dot tone mark: . U+302E

double dot tone mark: : U+302F

13) 2-complex and 3-complex letters

A 2-complex letter is a complex letter composed of two simple letters. Likewise, a 3-complex letter is a complex letter composed of three simple letters.

Examples of complex letters are shown below:

2-complex letters: ㅃ U+1101, ㅆ U+1162, ㅈㅅ U+11AD

3-complex letters: ㅅㅆ U+116B, ㅅㅆㅅ U+1170

14) IPF (Initial-Peak-Final) and Wanseong forms for Hangeul syllables

Modern Hangeul syllables can be represented either in IPF (initial-peak-final) form or in Wanseong (precomposed) form.

For example, syllable "가" can be represented as follows:

- a) in IPF form: U+1100 1161
- b) in Wanseong form: U+AC00

In contrast, Old Hangeul syllables must be represented in IPF form, since no Old Hangeul syllables in Wanseong form are included in ISO/IEC 10646-1.

For example, "Old Hangeul syllable "ㄱ" can be represented as follows:

- a) in IPF form: U+1100 119E
- b) in Wanseong form: the syllable cannot be represented in Wanseong form.

Table 2. Possibilities of representing Modern and Old Hangeul syllables in IPF or Wanseong form.

	in IPF form	in Wanseong form
Modern Hangeul	syllables composed using 67 Modern letters	11,172 syllables
Old Hangeul	syllables composed using 238 Modern and Old letters	(NO Old syllables included in ISO/IEC 10646-1)

15) 94 compatibility letters in ISO/IEC 10646-1 [U+3131 - 318E]

The 94 letters in the range of U+3131 - 318E are included for compatibility. They are from KS X 1001 (formerly, KS C 5601). Compatibility means that these letters can be used when converting Old files in KS X 1001 to files in ISO/IEC 10646. Completely new files are not supposed to be represented in compatibility letters; instead, they should be represented in U+11XX or U+AC00-D7A3.

Now let's see the usage of these letters.

- a) Each of these 93 letters (excluding the FILLER character: U+3141) can represent an independent letter.

Examples: ㄱ (U+3131), ㄱㄴ (U+3131 3134 3137), ㅏ (U+314F), ㅓ (U+3136).

b) Usage of FILLER character (U+3164): it has two different functions.

A sequence of four code positions can represent one syllable as follows:

FILLER (U+3164),

consonantal letter (U+3131-314E, U3165-3186),

vowel letter (U+314F-3163, U+3187-318E), and

consonantal (U+3131-314E, U3165-3186) letter or FILLER (U+3164).

FILLER indicates that the FILLER character and following three code positions together represent one syllable, not three independent letters.

FILLER (U+3164) before U+3131 in example c), (1) and the ones before U+3138 or 3131 in example c), (2) fall in this category.

In case of 2-letter syllable, FILLER is included at end to indicate that there is no syllable-final letter. FILLER (U+3164) after U+314F in example c), (1) below falls in this category.

c) A syllable, especially a syllable which is not one of the 2,350 syllables in KS X 1001, can be represented as follows:

(1) A 2-letter syllable can be represented in the form of FILLER (U+3164), consonantal letter, vowel letter, and FILLER (U+3164).

Example: 가 (U+3164 3131 314F 3164)

(2) A 3-letter syllable can be represented in the form of FILLER (U+3164), consonantal letter, vowel letter, and consonantal letter.

Examples: 뚝 (U+3164 3138 3157 3141). This syllable is not one of 2,350 syllables in KS X 1001.

강 (U+3164 3131 314F 3147)

d) "ㄱ ㅏ" (two independent letters) and "가" (one syllable) are represented differently as follows:

- ㄱ ㅏ: U+3131 314F (Two letters, not one syllable)

- 가: U+3164 3131 314F 3164 (One syllable. The first FILLER (U+3164 before U+3131) indicates that the FILLER and following three code positions together represent one syllable, not three letters. The last Filler (U+3164 after U+314F) indicates that this syllable is a two-letter syllable (i.e., no syllable-final letter).

2. Previous Works

In UAX #15 [UAX15], four normalization forms are described:

1) Normalization Form D (NFD): canonical decomposition

2) Normalization Form C (NFC): canonical decomposition, followed by canonical

composition

3) Normalization Form KD (NFKD): compatibility decomposition

4) Normalization Form KC (NFKC): compatibility decomposition, followed by canonical composition

Once a canonical decomposition process for Hangeul is defined, compatibility decomposition process for Hangeul can be easily defined. Therefore, in this paper, we will focus on canonical decomposition and composition processes for Hangeul, and as a result, Normalization Forms D (NFD) and C (NFC) for Hangeul.

2.1 Canonical Decomposition and Canonical Composition Processes

We may normalize Unicode-encoded text to one particular sequence. We can normalize into one of two forms [Unic30]:

1) systems that cannot handle nonspacing marks can normalize to precomposed characters;

2) In systems that can handle nonspacing marks, it may be useful to normalize so as to eliminate precomposed characters. ...

Process 1) above is NFD and process 2) above is NFC.

A character that is equivalent to a sequence of one or more other characters is called a decomposable character, a precomposed character, or a composite character.

A decomposition of a decomposable character is a sequence of one or more characters that is equivalent to a decomposable character.

2.2 Hangeul syllable decomposition process in UAX #15

The canonical decomposition process is described in [UAX15]. For a given Modern Hangeul syllable, a sequence composed of two or three of the 67 Modern Hangeul letters shown in 4), Section 1, is derived, which is a canonical decomposition. The sequence is composed of Modern Hangeul syllable-initial, syllable-peak, and optionally syllable-final letters. Whenever possible, a complex letter is used in place of two or three simple letters.

Example in [UAX15]:

- An original syllable: ㄷㅌ (U+AC03)

--> [NFD]: ㄷ + ㅌ (U+1100 1161 11AA) --> [NFC]: ㄷㅌ (U+AC03)

Note. A 2-complex letter ㅌ (U+11AA) is not decomposed any further, although it could be decomposed into two simple letters ㅌ (U+11A8) and ㅌ (U+11BA).

To help readers understand, we made up more examples:

강 (U+AC15) --> [NFD] ㄱ + ㅏ + ㅓ (U+1100 1161 11BC)
 강ㅏ (U+AC0C) --> [NFD] ㄱ + ㅏ + ㅓ (U+1100 1161 11B0)
 Note. ㅓ (U+11B0) is not decomposed any further.
 객 (U+AC1D) --> [NFD] ㄱ + ㅓ + ㅓ (U+1100 1162 11A8)
 Note. ㅓ (U+1162) is not decomposed any further.
 가 (U+AC00) --> [NFD] ㄱ + ㅏ (U+1100 1161)

2.3 Hangeul syllable composition in UAX #15

The composition process is described in [UAX15]. For a given canonical decomposition of a Modern Hangeul syllable, we can find a syllable whose code position is in the range of U+AC00 - D7A3. That syllable is a composition in [UAX15].

Note. In the Unicode 3.0 [Unic30], this composition is not described explicitly as "canonical" composition.

Example in [UAX15]:

- An original syllable: ㄱㅏ (U+AC03)
 --> [NFD] ㄱ + ㅏ + ㅓ (U+1100 1161 11AA) --> [NFC] ㄱㅏ (U+AC03)

To help readers understand, we made up more examples:

ㄱ + ㅏ + ㅓ (U+1100 1161 11BC) --> [NFC] 강 (U+AC15)
 ㄱ + ㅏ + ㅓ (U+1100 1161 11B0) --> [NFC] 강ㅏ (U+AC0C)
 ㄱ + ㅓ + ㅓ (U+1100 1162 11A8) --> [NFC] 객 (U+AC1D)
 ㄱ + ㅏ (U+1100 1161) --> [NFC] 가 (U+AC00)

3. New Canonical decomposition and composition processes for Hangeul

Although Normalization Forms KD (NFKD) and KC (NFKC) for Hangeul are not treated in this paper, based on new NFD and NFC, it should be straightforward to define new NFKD and NFKC. Therefore, in this paper, we will focus on new canonical decomposition and composition processes for Hangeul, and as a result, new Normalization Forms D (NFD) and C (NFC) for Hangeul.

3.1 The main idea of new Hangeul decomposition/composition processes

Canonical decomposition and composition processes for Hangeul, as shown in UAX #15, have limitations.

1) UAX #15 decomposition process only decomposes 11,172 Modern Hangeul. Therefore, UAX #15 decomposition process cannot decompose complex (both Modern and Old) letters.

2) UAX #15 composition process only composes Modern letters into syllables. Therefore, UAX #15 composition process cannot compose letters into Old complex letters.

Therefore, we propose new canonical decomposition and composition processes for Hangeul, Korean script. These new processes solve the limitations of UAX #15 decomposition/composition.

A new decomposition process decomposes Hangeul letters/syllables into simple letters. Then a new composition process composes simple letters into complex letters (whenever possible), but not into syllables. These are the main idea of new decomposition and composition processes.

Tables 3 and 4 compare UAX #15 decomposition/composition and new ones.

Table 3. Comparison of UAX #15 Decomposition and a new Decomposition.

Decomposition	input	output
UAX #15	11,172 Modern syllables only	1) Modern complex letters (whenever possible) 2) Modern simple letters
A New process	1) Modern syllables 2) Modern complex letters 3) Old complex letters	1) Modern simple letters 2) Old simple letters

Table 4. Comparison of UAX #15 Composition and a new Composition.

Composition	input	output
UAX #15	1) Modern complex letters 2) Modern simple letters	11,172 Modern syllables
A New process	1) Modern simple letters 2) Old simple letters	1) Modern complex letters (whenever possible) 2) Modern simple letters 3) Old complex letters (whenever possible) 4) Old simple letters

3.2 A new decomposition process for Hangeul

A new decomposition process decomposes syllables and complex (both Modern and Old) letters as follows. To decompose a syllable, start at step D1; to decompose a complex letter, start at step D2.

step D1) For each Hangeul syllable in Wanseong form, decompose a syllable into syllable-initial, syllable-peak, and possibly syllable-final letters as described in [ISO10646].

step D2) For each complex letter, decompose a complex letter into simple-letter-1, simple-letter-2, and possibly simple-letter-3. For information about decomposing a complex letter into simple letters, see Appendix 1.

step D3) The result is the canonical decomposition of the given Hangeul syllable or complex letter. Note that the canonical decomposition has only simple letters and possibly Filler characters and tone marks, but neither syllables in Wanseong form nor complex letters.

First, we will review the example in [UAX15]:

- An original syllable: ㄷᄇᆞᆫ (U+AC03)

--> [NFD] ㄷ + ㅏ + ㅓ (U+1100 1161 11AA)

As noted earlier, a 2-complex letter ㅓ (U+11AA) is not decomposed any further in UAX #15, although it could be decomposed into two simple letters ㅓ (U+11A8) and ㅓ (U+11BA).

At step D1, an original syllable (U+AC03) is decomposed into three letters:

ㄷᄇᆞᆫ (U+AC03) --> ㄷ + ㅏ + ㅓ (U+1100 1161 11AA)

Then, at step D2, a 2-complex letter (U+11AA) is further decomposed into two simple letters (U+11A8) and (U+11BA):

ㅓ (U+11AA) --> ㅓ (U+11A8) + ㅓ (U+11BA)

Therefore, the final decomposition is:

- An original syllable: ㄷᄇᆞᆫ (U+AC03)

--> [NFD] ㄷ + ㅏ + ㅓ + ㅓ (U+1100 1161 11A8 11BA)

Now let's take some more examples.

Example 3.1

Example 3.1a

⌌ (U+116B) --> [NFD] ⌌ (U+1169) + ⌌ (U+1161) + ⌌ (U+1175)

Example 3.1b

⌌ (U+116A) + ⌌ (U+1175) --> [NFD] ⌌ (U+1169) + ⌌ (U+1161) + ⌌ (U+1175)

Example 3.1c

⌌ (U+1169) + ⌌ (U+1162) --> [NFD] ⌌ (U+1169) + ⌌ (U+1161) + ⌌ (U+1175)

Example 3.1d

⌌ (U+1169) + ⌌ (U+1161) + ⌌ (U+1175) --> (Since all letters are simple, we cannot decompose any further.)

The above four original sequences are considered equal when rendered. However, UAX #15 decomposition canNOT handle any of the four examples above, since that process handles only Modern syllables. In contrast, our new decomposition process decomposes a complex letter into simple letters so that four sequences become the same in a decomposed form.

Example 3.2 Our new decomposition process decomposes the following three sequences as follows:

Example 3.2a: ⌌ (U+AC1C) --> [NFD] ⌌ (U+1100) + ⌌ (U+1161) + ⌌ (U+1175)

Example 3.2b: ⌌ (U+1100) + ⌌ (U+1162) --> [NFD] ⌌ (U+1100) + ⌌ (U+1161) + ⌌ (U+1175)

Example 3.2c: ⌌ (U+1100) + ⌌ (U+1161) + ⌌ (U+1175) --> (no further decomposition, since all letters are already simple.)

The above three original sequences are considered equal when rendered and our new decompositions of them are the same.

In contrast, the results according to the UAX #15 canonical decomposition are shown below:

Example 3.2a: ⌌ (U+AC1C) --> ⌌ (U+1100) + ⌌ (U+1162)

Example 3.2b: ⌌ (U+1100) + ⌌ (U+1162) --> (no further decomposition)

Example 3.2c: ⌌ (U+1100) + ⌌ (U+1161) + ⌌ (U+1175) --> (no further decomposition)

As you can see, in a UAX #15 decomposed form, two sequences of examples 3.2a and 3.2b are equal. However, they are different from the sequence of example 3.2c.

In contrast, our new decomposition further decomposes the complex letter into

simple letters so that the three sequences of examples 3.2a, 3.2b, and 3.2c are all equal.

The above examples 3.1 and 3.2 show why our new canonical decomposition for Hangeul is better the one in UAX #15.

Example 3.3: Now let's consider Old Hangeul syllables and Old complex letters. UAX #15 decomposition does not accept at all Old Hangeul syllables or Old complex letters. In contrast, our new decomposition process handles Old syllables and Old complex letters just like Modern syllables and Modern complex letters.

- original: ㅁㅓ (U+1122) + ㅓ (U+1161)
--> [NFD] ㅁ (U+1107 + ㅓ (U+1109) + ㅓ (U+1100) + ㅓ (U+1161)

3.3 A new canonical composition process for Hangeul

In ISO/IEC 10646-1, 11,172 Modern Hangeul syllables in Wanseong form are included; however, no Old Hangeul syllables in Wanseong form are included. Given a decomposition, we can compose letters into Modern syllables in Wanseong form (U+AC00 - D7A3) easily, but not into Old syllables in Wanseong form.

To be able to apply consistent canonical composition to both Modern and Old syllables, we suggest that letters be composed into complex letters, not into syllables.

In Unicode 3.0 Standard [Unic30] and UAX #15 [UAX15], canonical decomposition of Modern Hangeul syllables is defined; however, the author could not find an "explicit" definition of canonical composition for Hangeul.

There is a description of a composition for Hangeul in the book. Even if we assume that the composition is a canonical composition, the composition is incomplete in the sense that it handles only Modern syllables, but not Old syllables. In other words, the composition algorithm composes letters into Modern syllables; however, we cannot compose a sequence of letters corresponding to an Old Hangeul syllable into an Old syllable in Wanseong form, since there are no Old Hangeul syllables in ISO/IEC 10646.

In this paper, we propose a new canonical composition which can handle both Modern and Old syllables consistently. Given a decomposition composed of simple letters only, we can take the following steps to get a canonical composition.

Step C1) A sequence of simple letters in a canonical decomposition of Hangeul are grouped into syllable-initial, syllable-peak, and possibly syllable-final

letters.

Step C2) For each group of syllable-initial, syllable-peak or syllable-final simple letters, compose as follows:

- For a simple letter, return the simple letter;
- For a group of two simple letters composed of letter-1 and letter-2:
If there is a 2-complex letter, letter-1-2, composed of letter-1 and letter-2, then return letter-1-2;
otherwise return a sequence of letter-1 and letter-2;
- For a group of three simple letters composed of letter-1, letter-2, and letter-3:
If there is a 3-complex letter, letter-1-2-3, composed of letter-1, letter-2, and letter-3, then return letter-1-2-3;
else if there is a 2-complex letter, letter-1-2, composed of letter-1 and letter-2, then return a sequence of letter-1-2 and letter-3;
else if there is a 2-complex letter, letter-2-3, composed of letter-2 and letter-3, then return a sequence of letter-1 and letter-2-3;
otherwise return a sequence of letter-1, letter-2, and letter-3;

Note. It is assumed that a complex letter has at most three simple letters, which is true for any Hangeul syllable/letter found up to date. However, the above algorithm can be easily extended to accommodate more than three simple letters in a group.

First, we will review the example in [UAX15]:

- An original syllable: ㄷᄇᆞ (U+AC03)
--> [NFD] ㄷ + ㅏ + ㅑ (U+1100 1161 11AA) --> [NFC] ㄷᄇᆞ (U+AC03)

In contrast, our new composition composes into complex letters, not into syllables:

- an original syllable: ㄷᄇᆞ (U+AC03)
--> [NFD] ㄷ + ㅏ + ㅑ + ㅓ (U+1100 1161 11A8 11BA)
--> [NFC] ㄷ + ㅏ + ㅑ (U+1100 1161 11AA)

Now let's continue with examples 3.1, 3.2, and 3.3 above. All four original sequences in example 3.1 ended up with the same decomposition as shown below:

[NFD] ㅓ (U+1169) + ㅏ (U+1161) + ㅑ (U+1175)

Its composition is shown below:

--> [NFC] ㅓᄇᆞ (U+116B)

Likewise, all three sequences in example 3.2 above ended up with the same decomposition as shown below:

[NFD] ㄱ (U+1100) + ㅏ (U+1161) + ㅣ (U+1175)

Its composition is shown below:

--> [NFC] ㄱ (U+1100) + ㅑ (U+1162)

Now let's consider Old Hangeul syllables and Old complex letters. UAX #15 composition does not accept at all Old Hangeul syllables or Old complex letters. In contrast, our new composition process handles Old syllables and Old complex letters just like Modern syllables and Modern complex letters.

Canonical decomposition in example 3.3 above was:

[NFD] ㅓ (U+1107) + ㅏ (U+1109) + ㄱ (U+1100) + ㅏ (U+1161).

Its composition is shown below:

--> [NFC] ㅕ (U+1122) + ㅏ (U+1161)

3.4 A canonical composition process for Hangeul/option-Modern-syl

Since we frequently use Modern Hangeul syllables, we could modify the proposed composition process slightly to utilize Modern Hangeul syllables in Wanseong form included in ISO/IEC 10646. A canonical composition modified as shown in this section will be called "composition process/option-Modern-syl".

If, after composition as described in section 3.3, a sequence of letters is found to be a Modern syllable, we can further compose the sequence of letters into a Modern syllable.

Although the author believes that the composition in section 3.3 is better than the one in this section, sometimes we may find this modified composition process useful.

This modified composition process is basically as follows:

a) For Old Hangeul, we use our new composition process; that is, we compose simple letters into complex letters, if possible.

b) For Modern Hangeul, we use the composition process described in UAX #15; that is, after applying the new composition process, we further compose letters into a Modern syllable, if possible.

4. Conclusions

In this paper, the author proposed new canonical decomposition process (NFD: Normalization Form D) and composition processes (NFC: Normalization Form C) for Hangeul, Korean script. The new processes can handle both Modern and Old Hangeul in a consistent way. Therefore, the new processes are more general than the ones described in UAX #15 [UAX15], which handle only Modern Hangeul syllables.

Since NFD and NFC will play an important role in sorting and searching, Internet-related applications such as Multilingual Domain names, databases, etc., we need to establish good NFD and NFC for Hangeul before they are widely used.

The author hopes that NFD and NFC proposed in this paper will be reviewed and reflected in ISO/IEC 14651 International String Ordering and Comparison [ISO14651], Unicode Standard Annex #15 [UAX15], Nameprep for IDN (Internationalized Domain Name) [Nameprep], and other related documents.

References.

[ISO10646] ISO/IEC 10646-1:2000, Information Technology -- Universal Multiple-Octet Coded Character Set (UCS) -- Part 1: Architecture and Basic Multilingual Plane. 2000-10-05.

[ISO11941] ISO TS (Technical Specification) 11941, Information and Documentation -- Transliteration of Korean script into Latin characters. First edition, 1996-12-01.

[ISO14651] ISO/IEC 14651, International String Ordering and Comparison - Method for comparing character strings and description of the common template tailorable ordering

[KimK99] A story about Hangeul inside computers, Vol. 2. KIM, Kyongsok. Busan National University Press. Dec. 31, 1999 [in Hangeul (Korean)].

[Nameprep] Paul Hoffman, Marc Blanchet, "Preparation of Internationalized Host Names", Feb 2001,
<http://www.ietf.org/internet-drafts/draft-ietf-idn-nameprep-03.txt>

[UAX15] Unicode Standard Annex #15. Technical Report. Unicode Normalization Forms. Mark Davis and Martin Duerst. 2001-03-23.

[Unic30] The Unicode Standard Version 3.0, The Unicode Consortium. 2000.

Appendices :

Appendix 1. Decomposition of complex letters into simple letters.

1) Decomposition of syllable-initial complex letters into simple letters.

U+1101 -> U+1100 1100

U+1104 -> U+1103 1103

U+1108 -> U+1107 1107

U+110A -> U+1109 1109

U+110D -> U+110C 110C

U+1113 -> U+1102 1100

U+1114 -> U+1102 1102

U+1115 -> U+1102 1103

U+1116 -> U+1102 1107

U+1117 -> U+1103 1100

U+1118 -> U+1105 1102

U+1119 -> U+1105 1105

U+111A -> U+1105 1112

U+111B -> U+1105 110B

U+111C -> U+1106 1107

U+111D -> U+1106 110B

U+111E -> U+1107 1100

U+111F -> U+1107 1102

U+1120 -> U+1107 1103

U+1121 -> U+1107 1109

U+1122 -> U+1107 1109 1100

U+1123 -> U+1107 1109 1103

U+1124 -> U+1107 1109 1107

U+1125 -> U+1107 1109 1109

U+1126 -> U+1107 1109 110C

U+1127 -> U+1107 110C

U+1128 -> U+1107 110E

U+1129 -> U+1107 1110

U+112A -> U+1107 1111

U+112B -> U+1107 110B

U+112C -> U+1107 1107 110B

U+112D -> U+1109 1100

U+112E -> U+1109 1102

U+112F -> U+1109 1103

U+1130 -> U+1109 1105

U+1131 -> U+1109 1106

U+1132 -> U+1109 1107
U+1133 -> U+1109 1107 1100
U+1134 -> U+1109 1109 1109
U+1135 -> U+1109 110B
U+1136 -> U+1109 110C
U+1137 -> U+1109 110E
U+1138 -> U+1109 110F
U+1139 -> U+1109 1110
U+113A -> U+1109 1111
U+113B -> U+1109 1112
U+113D -> U+113C 113C
U+113F -> U+113E 113E

U+1140 -> U+1107 1103
U+1141 -> U+110B 1100
U+1142 -> U+110B 1103
U+1143 -> U+110B 1106
U+1144 -> U+110B 1107
U+1145 -> U+110B 1109
U+1146 -> U+110B 1140
U+1147 -> U+110B 110B
U+1148 -> U+110B 110C
U+1149 -> U+110B 110E
U+114A -> U+110B 1110
U+114B -> U+110B 1111
U+114D -> U+110C 110B
U+114F -> U+114E 114E

U+1151 -> U+1150 1150
U+1152 -> U+110E 110F
U+1153 -> U+1103 1112
U+1156 -> U+1111 1107
U+1157 -> U+1111 110B
U+1158 -> U+1112 1112

2) Decomposition of syllable-peak complex letters into simple letters.

U+1162 -> U+1161 1175
U+1164 -> U+1163 1175
U+1166 -> U+1165 1175
U+1168 -> U+1167 1175
U+116A -> U+1169 116A
U+116B -> U+1169 116A 1175

U+116C -> U+1169 1175
U+116F -> U+116E 1165
U+1170 -> U+116E 1165 1175
U+1171 -> U+116E 1175
U+1174 -> U+1173 1175
U+1176 -> U+1161 1169
U+1177 -> U+1161 116E
U+1178 -> U+1163 1169
U+1179 -> U+1163 116D
U+117A -> U+1165 1169
U+117B -> U+1165 116E
U+117C -> U+1165 1173
U+117D -> U+1167 1169
U+117E -> U+1167 116E
U+117F -> U+1169 1165
U+1180 -> U+1169 1165 1175
U+1181 -> U+1169 1167 1175
U+1182 -> U+1169 1169
U+1183 -> U+1169 116E
U+1184 -> U+116D 1163
U+1185 -> U+116D 1163 1175
U+1186 -> U+116D 1167
U+1187 -> U+116D 1169
U+1188 -> U+116D 1175
U+1189 -> U+116E 1161
U+118A -> U+116E 1161 1175
U+118B -> U+116E 1165 1173
U+118C -> U+116E 1167 1175
U+118C -> U+116E 1167 1175
U+118D -> U+116E 116E
U+118F -> U+1172 1165
U+1190 -> U+1172 1165 1175
U+1191 -> U+1172 1167
U+1192 -> U+1172 1167 1175
U+1193 -> U+1172 116E
U+1194 -> U+1172 1175
U+1195 -> U+1173 116E
U+1196 -> U+1173 1173
U+1197 -> U+1173 1175 116E
U+1198 -> U+1175 1161
U+1199 -> U+1175 1163
U+119A -> U+1175 1169
U+119B -> U+1175 116E
U+119C -> U+1175 1173

U+119D -> U+1175 119E
U+119F -> U+119E 1165
U+11A0 -> U+119E 116E
U+11A1 -> U+119E 1175
U+11A2 -> U+119E 119E

3) Decomposition of syllable-final complex letters into simple letters.

U+11A9 -> U+11A8 11A8
U+11AA -> U+11A8 11BA
U+11AC -> U+11AB 11BD
U+11AD -> U+11AB 11C2








U+11B0 -> U+11AF 11A8
U+11B1 -> U+11AF 11B7
U+11B2 -> U+11AF 11B8
U+11B3 -> U+11AF 11BA
U+11B4 -> U+11AF 11C0
U+11B5 -> U+11AF 11C1
U+11B6 -> U+11AF 11C2
U+11B9 -> U+11B8 11BA
U+11BB -> U+11BA 11BA
U+11C3 -> U+11A8 11AF
U+11C4 -> U+11A8 11BA 11A8
U+11C5 -> U+11AB 11A8
U+11C6 -> U+11AB 11AE
U+11C7 -> U+11AB 11BA
U+11C8 -> U+11AB 11EB
U+11C9 -> U+11AB 11C0
U+11CA -> U+11AE 11A8
U+11CB -> U+11AE 11AF
U+11CC -> U+11AF 11A8 11BA
U+11CD -> U+11AF 11AB
U+11CE -> U+11AF 11AE
U+11CF -> U+11AF 11AE 11C2
U+11D0 -> U+11AF 11AF
U+11D1 -> U+11AF 11B7 11A8
U+11D2 -> U+11AF 11B7 11BA
U+11D3 -> U+11AF 11B8 11BA
U+11D4 -> U+11AF 11B8 11C2
U+11D5 -> U+11AF 11B8 11BC
U+11D6 -> U+11AF 11BA 11BA
U+11D7 -> U+11AF 11EB

U+11D8 -> U+11AF 11BF
U+11D9 -> U+11AF 11F9
U+11DA -> U+11B7 11A8
U+11DB -> U+11B7 11AF
U+11DC -> U+11B7 11B8
U+11DD -> U+11B7 11BA
U+11DE -> U+11B7 11BA 11BA
U+11DF -> U+11B7 11EB
U+11E0 -> U+11B7 11BE
U+11E1 -> U+11B7 11C2
U+11E2 -> U+11B7 11BC
U+11E3 -> U+11B8 11AF
U+11E4 -> U+11B8 11C1
U+11E5 -> U+11B8 11C2
U+11E6 -> U+11B8 11BC
U+11E7 -> U+11BA 11A8
U+11E8 -> U+11BA 11AE
U+11E9 -> U+11BA 11AF
U+11EA -> U+11BA 11B8
U+11EC -> U+11BC 11A8
U+11ED -> U+11BC 11A8 11A8
U+11EE -> U+11BC 11BC
U+11EF -> U+11BC 11BF
U+11F1 -> U+11F0 11BA
U+11F2 -> U+11F0 11EB
U+11F3 -> U+11C1 11B8
U+11F4 -> U+11C1 11BC
U+11F5 -> U+11C2 11AB
U+11F6 -> U+11C2 11AF
U+11F7 -> U+11C2 11B7
U+11F8 -> U+11C2 11B8
* * *

Appendix 2. Hangeul letters: U+1100 - 117F [KS X 1005-1] (contiuned)

표 35 - 행 11: 한글 자모

Table 35 - Row 11: Hangeul Jamo

	110	111	112	113	114	115	116	117
0	ㄱ 1100	ㅌ 1110	ㅊ 1120	ㅋ 1130	△ 1140	ㅈ 1150	 1160	계 1170
1	기 1101	토티 1111	차 1121	사 1131	오 1141	ㅉ 1151	ㅊ 1161	기 1171
2	ㄴ 1102	ㅇ 1112	ㅅ 1122	ㅆ 1132	ㅇ 1142	ㅈ 1152	ㅊ 1162	ㅊ 1172
3	ㄷ 1103	ㄴ 1113	ㅅ 1123	ㅆ 1133	ㅇ 1143	ㅈ 1153	ㅊ 1163	ㅊ 1173
4	ㄸ 1104	ㄴ 1114	ㅅ 1124	ㅆ 1134	ㅇ 1144	ㅈ 1154	ㅊ 1164	ㅊ 1174
5	ㄹ 1105	ㄴ 1115	ㅅ 1125	ㅆ 1135	ㅇ 1145	ㅈ 1155	ㅊ 1165	ㅊ 1175
6	ㅁ 1106	내 1116	ㅅ 1126	ㅆ 1136	ㅇ 1146	ㅈ 1156	ㅊ 1166	ㅊ 1176
7	ㅂ 1107	ㄷ 1117	ㅅ 1127	ㅆ 1137	ㅇ 1147	ㅈ 1157	ㅊ 1167	ㅊ 1177
8	ㅃ 1108	ㄴ 1118	ㅅ 1128	ㅆ 1138	ㅇ 1148	ㅈ 1158	ㅊ 1168	ㅊ 1178
9	ㅅ 1109	ㄴ 1119	ㅅ 1129	ㅆ 1139	ㅇ 1149	ㅈ 1159	ㅊ 1169	ㅊ 1179
A	ㅆ 110A	ㄴ 111A	ㅅ 112A	ㅆ 113A	ㅇ 114A		ㅊ 116A	ㅊ 117A
B	ㅇ 110B	ㄴ 111B	ㅅ 112B	ㅆ 113B	ㅇ 114B		ㅊ 116B	ㅊ 117B
C	ㅈ 110C	ㅊ 111C	ㅅ 112C	ㅆ 113C	ㅇ 114C		ㅊ 116C	ㅊ 117C
D	ㅉ 110D	ㅊ 111D	ㅅ 112D	ㅆ 113D	ㅇ 114D		ㅊ 116D	ㅊ 117D
E	ㅊ 110E	ㅊ 111E	ㅅ 112E	ㅆ 113E	ㅇ 114E		ㅊ 116E	ㅊ 117E
F	ㅋ 110F	ㅌ 111F	ㅊ 112F	ㅋ 113F	ㅉ 114F	 115F	계 116F	겨 117F

G = 00
P = 00

Hangeul letters: U+1180 - 11FF [KS X 1005-1]

표 36 - 행 11: 한글 자모

Table 36 - Row 11: Hangeul Jamo

	118	119	11A	11B	11C	11D	11E	11F
0	계 1180	개 1190	구 11A0	리 11B0	ㄷ 11C0	르 11D0	ㅌ 11E0	ㅇ 11F0
1	계 1181	겨 1191	기 11A1	려 11B1	도 11C1	랴 11D1	트 11E1	야 11F1
2	ㅊ 1182	계 1192	'' 11A2	래 11B2	ㅎ 11C2	랴 11D2	몽 11E2	ㅇ 11F2
3	누 1183	푸 1193		리 11B3	기 11C3	랴 11D3	비 11E3	표 11F3
4	ㅍ 1184	기 1194		려 11B4	끼 11C4	랴 11D4	비 11E4	풍 11F4
5	배 1185	두 1195		르 11B5	니 11C5	령 11D5	비 11E5	하 11F5
6	벼 1186	= 1196		르 11B6	니 11C6	랴 11D6	빙 11E6	하 11F6
7	ㅆ 1187	쿠 1197		리 11B7	니 11C7	랴 11D7	시 11E7	하 11F7
8	피 1188	기 1198	기 11A8	비 11B8	니 11C8	랴 11D8	시 11E8	하 11F8
9	기 1189	기 1199	기 11A9	비 11B9	니 11C9	랴 11D9	시 11E9	하 11F9
A	개 118A	고 119A	가 11AA	나 11BA	다 11CA	마 11DA	사 11EA	
B	겨 118B	구 119B	나 11AB	사 11BB	다 11CB	마 11DB	사 11EB	
C	계 118C	기 119C	나 11AC	ㅇ 11BC	랴 11CC	매 11DC	ㅇ 11EC	
D	두 118D	! 119D	나 11AD	스 11BD	리 11CD	마 11DD	ㅇ 11ED	
E	기 118E	! 119E	나 11AE	스 11BE	리 11CE	마 11DE	ㅇ 11EE	
F	기 118F	기 119F	리 11AF	쿠 11BF	랴 11CF	마 11DF	ㅇ 11EF	

G = 00
P = 00

Appendix 3. Hangeul syllables: U+AC00 – ACAF [KS X 1005-1]
 (Out of 11,172 Hangeul syllables, only 176 syllables are shown below.)

표 98 - 행 AC: 한글 글자마디

Table 98 - Row AC: Hangeul Syllables

	AC0	AC1	AC2	AC3	AC4	AC5	AC6	AC7	AC8	AC9	ACA
0	가 AC00	감 AC10	갠 AC20	갬 AC30	갈 AC40	각 AC50	갬 AC60	거 AC70	검 AC80	겐 AC90	갬 ACA0
1	각 AC01	갑 AC11	갬 AC21	갹 AC31	갹 AC41	갹 AC51	갹 AC61	거 AC71	겁 AC81	갹 AC91	갹 ACA1
2	각 AC02	갹 AC12	갹 AC22	갹 AC32	갹 AC42	갹 AC52	갹 AC62	거 AC72	갹 AC82	갹 AC92	갹 ACA2
3	갹 AC03	갹 AC13	갹 AC23	갹 AC33	갹 AC43	갹 AC53	갹 AC63	거 AC73	갹 AC83	갹 AC93	갹 ACA3
4	간 AC04	갹 AC14	갹 AC24	갹 AC34	갹 AC44	갹 AC54	갹 AC64	건 AC74	갹 AC84	갹 AC94	갹 ACA4
5	갹 AC05	갹 AC15	갹 AC25	갹 AC35	갹 AC45	갹 AC55	갹 AC65	건 AC75	갹 AC85	갹 AC95	갹 ACA5
6	갹 AC06	갹 AC16	갹 AC26	갹 AC36	갹 AC46	갹 AC56	갹 AC66	건 AC76	갹 AC86	갹 AC96	갹 ACA6
7	간 AC07	갹 AC17	갹 AC27	갹 AC37	갹 AC47	갹 AC57	갹 AC67	건 AC77	갹 AC87	갹 AC97	갹 ACA7
8	갈 AC08	각 AC18	갹 AC28	가 AC38	갹 AC48	갹 AC58	갹 AC68	걸 AC78	갹 AC88	갹 AC98	갹 ACA8
9	갹 AC09	갹 AC19	갹 AC29	갹 AC39	갹 AC49	갹 AC59	갹 AC69	걸 AC79	갹 AC89	갹 AC99	갹 ACA9
A	갹 AC0A	갹 AC1A	갹 AC2A	갹 AC3A	갹 AC4A	갹 AC5A	갹 AC6A	걸 AC7A	갹 AC8A	갹 AC9A	갹 ACA A
B	갹 AC0B	갹 AC1B	갹 AC2B	갹 AC3B	갹 AC4B	갹 AC5B	갹 AC6B	걸 AC7B	갹 AC8B	갹 AC9B	갹 ACA B
C	갹 AC0C	갹 AC1C	갹 AC2C	갹 AC3C	갹 AC4C	갹 AC5C	갹 AC6C	걸 AC7C	갹 AC8C	갹 AC9C	갹 ACA C
D	갹 AC0D	갹 AC1D	갹 AC2D	갹 AC3D	갹 AC4D	갹 AC5D	갹 AC6D	걸 AC7D	갹 AC8D	갹 AC9D	갹 ACA D
E	갹 AC0E	갹 AC1E	갹 AC2E	갹 AC3E	갹 AC4E	갹 AC5E	갹 AC6E	걸 AC7E	갹 AC8E	갹 AC9E	갹 ACA E
F	갹 AC0F	갹 AC1F	갹 AC2F	갹 AC3F	갹 AC4F	갹 AC5F	갹 AC6F	걸 AC7F	갹 AC8F	갹 AC9F	갹 ACA F

G = 00
P = 00



Author's Vita :

KIM, Kyongsok (GIM, Gyeongseog) is a Professor in Division of Computer Science and Engineering at Busan National University, Busan 609-735, South Korea. His e-mail address is gimgs@HANGEUL.pnu.edu. He received Bachelor's and Master's degrees at Seoul National University, South Korea, and Ph. D in Computer Science at University of Illinois at Urbana-Champaign, USA in 1988.

He worked at North Dakota State University as an assistant professor from 1988 to 1992. He has been working at Busan National University since 1992. He stayed in Dept. of Phonetics and Linguistics at University College London, U. K from 1997 to 1998 as a visiting professor.

Currently he is a chairperson of Korea JTC1/SC2 committee on Code; a member of Korea TC46 committee on Document Information and Hangeul Romanization; a member of Korea SC22 (Programming Languages, esp. WG20 - Internationalization); a member of Name Committe at KRNIC (Korea Network Information Center). He is a member of KLS (Korean Language Society), KLIS (Korean Language Information Society), and KISS (Korea Information Science Society).

Information about Graphics Files for Figures :

There are four graphics files as follows:

- app2a.jpg for the first figure in appendix 2
- app2b.jpg for the second figure in appendix 2
- app3.jpg for the figure in appendix 3
- gimgs12w.jpg for author's photo

* * *