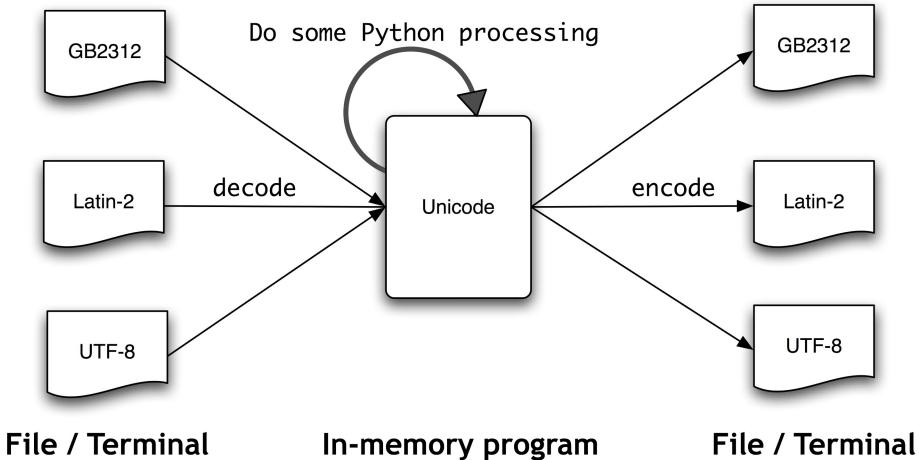


Unicode for korean characters

Wednesday, May 27, 2020 11:31 PM



<https://www.nltk.org/book/ch03.html>

3.3 Text Processing with Unicode

Unicode supports over a million characters including Korean, Chinese, Arabic, etc.

- Code point: Each character is assigned a number, called a **code point**.
- From a Unicode perspective, characters are abstract entities which can be realized (visualized) as one or more **glyphs**.
- Only glyphs appear on your screen or be printed on paper.
- A **font** is a mapping from characters to glyphs.

<https://docs.python.org/3/howto/unicode.html>

Introduction to Unicode

- A code point value is an integer in **[0, 0x10FFFF]**
- In the standard, a code point is written using the notation U+3C01 to mean the character with code point value 0x3C01.

UTF-8

- UTF: Unicode Transformation Format
- 8 means that 8-bit values are used in the encoding.
- One of the most commonly used encoding, and python 3 default.
- Rule:
 1. If the code point is < 128, it's represented by the corresponding byte value (ASCII)
 2. If the code point is ≥ 128 , it's turned into a sequence of two, three, or four bytes, where each byte of the sequence is between 128 and 255.
- Properties:
 1. It can handle any Unicode code point
 2. NULL character is U+0000 only. It can be processed by C functions such as strcpy().
 3. ASCII is valid.

4. Fairly compact. The majority of commonly used characters can be represented with one or two bytes.
5. If bytes are corrupted or lost, it's possible to determine the start of the next UTF-8-encoded code point and resynchronize.

https://en.wikipedia.org/wiki/Korean_language_and_computers

https://en.wikipedia.org/wiki/Hangul_Syllables

Hangul in Unicode

- Hangul Syllables (완성형 한글 코드)
 - Hangul Syllables is a Unicode block containing precomposed Hangul syllable blocks for modern Korean.
 - Range: U+AC00 ... U+D7AF
 - Source Standards: KS C 5601-1992

| | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | A | B | C | D | E | F |
|--------|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| U+AC0x | 가 | 각 | 깎 | 갚 | 간 | 갓 | 꺗 | 갇 | 갈 | 갉 | 꺃 | 꺃 | 꺃 | 꺃 | 꺃 | 꺃 |
| U+AC1x | 감 | 갑 | 값 | 갓 | 갓 | 강 | 갓 | 깟 | 깍 | 깍 | 꺃 | 꺃 | 꺃 | 꺃 | 꺃 | 꺃 |
| U+AC2x | 갠 | 갰 | 깁 | 깬 | 깰 | 깰 | 깁 | 깁 | 깁 | 깁 | 꺃 | 꺃 | 꺃 | 꺃 | 꺃 | 꺃 |
| U+AC3x | 갰 | 갱 | 갯 | 깁 | 깰 | 깁 | 깁 | 깁 | 깁 | 깁 | 꺃 | 꺃 | 꺃 | 꺃 | 꺃 | 꺃 |
| U+AC4x | 꺃 | 꺃 | 꺃 | 꺃 | 꺃 | 꺃 | 꺃 | 꺃 | 꺃 | 꺃 | 꺃 | 꺃 | 꺃 | 꺃 | 꺃 | 꺃 |
| U+AC5x | 꺃 | 꺃 | 꺃 | 꺃 | 꺃 | 꺃 | 꺃 | 꺃 | 꺃 | 꺃 | 꺃 | 꺃 | 꺃 | 꺃 | 꺃 | 꺃 |
| U+AC6x | 꺃 | 꺃 | 꺃 | 꺃 | 꺃 | 꺃 | 꺃 | 꺃 | 꺃 | 꺃 | 꺃 | 꺃 | 꺃 | 꺃 | 꺃 | 꺃 |
| U+D74x | ᡥ | ᡥ | ᡥ | ᡥ | ᡥ | ᡥ | ᡥ | ᡥ | ᡥ | ᡥ | ᡥ | ᡥ | ᡥ | ᡥ | ᡥ | ᡥ |
| U+D75x | ᡥ | ᡥ | ᡥ | ᡥ | ᡥ | ᡥ | ᡥ | ᡥ | ᡥ | ᡥ | ᡥ | ᡥ | ᡥ | ᡥ | ᡥ | ᡥ |
| U+D76x | ᡥ | ᡥ | ᡥ | ᡥ | ᡥ | ᡥ | ᡥ | ᡥ | ᡥ | ᡥ | ᡥ | ᡥ | ᡥ | ᡥ | ᡥ | ᡥ |
| U+D77x | ᡥ | ᡥ | ᡥ | ᡥ | ᡥ | ᡥ | ᡥ | ᡥ | ᡥ | ᡥ | ᡥ | ᡥ | ᡥ | ᡥ | ᡥ | ᡥ |
| U+D78x | ᡥ | ᡥ | ᡥ | ᡥ | ᡥ | ᡥ | ᡥ | ᡥ | ᡥ | ᡥ | ᡥ | ᡥ | ᡥ | ᡥ | ᡥ | ᡥ |
| U+D79x | ᡥ | ᡥ | ᡥ | ᡥ | ᡥ | ᡥ | ᡥ | ᡥ | ᡥ | ㅎ | ㅎ | ㅎ | ㅎ | ㅎ | ㅎ | ㅎ |
| U+D7Ax | ᡥ | ᡥ | ᡥ | ᡥ | ㅎ | ㅎ | ㅎ | ㅎ | ㅎ | ㅎ | ㅎ | ㅎ | ㅎ | ㅎ | ㅎ | ㅎ |

[https://en.wikipedia.org/wiki/Hangul_Jamo_\(Unicode_block\)](https://en.wikipedia.org/wiki/Hangul_Jamo_(Unicode_block))

조합형 한글을 위한 Unicode code points

Hangul Jamo ([Korean](#)): 한글 자모, Korean pronunciation: [haŋgul t̪əmo]) is a [Unicode block](#) containing positional (Choseong, Jungseong, and Jongseong) forms of the Hangul consonant and vowel clusters. They can be used to dynamically compose

syllables that are not available as [precomposed Hangul syllables in Unicode](#), specifically archaic syllables containing sounds that have since merged phonetically with other sounds in modern pronunciation.

| Hangul Jamo ^[1] | | | | | | | | | | | | | | | | | | | |
|--|------|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|------|--|
| Official Unicode Consortium code chart  (PDF) | | | | | | | | | | | | | | | | | | | |
| | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | A | B | C | D | E | F | | | |
| U+110x | ㄱ | ㄲ | ㄴ | ㄷ | ㄸ | ㄹ | ㅁ | ㅂ | ㅃ | ㅅ | ㅆ | ㅈ | ㅉ | ㅊ | ㅋ | ㅌ | ㅍ | ㅎ | |
| U+111x | ㅌ | ㅍ | ㅎ | ㄴ | ㄴ | ㄴ | ㄴ | ㄴ | ㄴ | ㄴ | ㄴ | ㄴ | ㄴ | ㄴ | ㄴ | ㄴ | ㄴ | ㄴ | |
| U+112x | ㅂ | ㅂ | ㅂ | ㅂ | ㅂ | ㅂ | ㅂ | ㅂ | ㅂ | ㅂ | ㅂ | ㅂ | ㅂ | ㅂ | ㅂ | ㅂ | ㅂ | ㅂ | |
| U+113x | ㅅ | ㅅ | ㅅ | ㅅ | ㅅ | ㅅ | ㅅ | ㅅ | ㅅ | ㅅ | ㅅ | ㅅ | ㅅ | ㅅ | ㅅ | ㅅ | ㅅ | ㅅ | |
| U+114x | ㅈ | ㅈ | ㅈ | ㅈ | ㅈ | ㅈ | ㅈ | ㅈ | ㅈ | ㅈ | ㅈ | ㅈ | ㅈ | ㅈ | ㅈ | ㅈ | ㅈ | ㅈ | |
| U+115x | ㅊ | ㅊ | ㅊ | ㅊ | ㅊ | ㅊ | ㅊ | ㅊ | ㅊ | ㅊ | ㅊ | ㅊ | ㅊ | ㅊ | ㅊ | ㅊ | ㅊ | HJ F | |
| U+116x | HJ F | ㅏ | ㅓ | ㅗ | ㅓ | ㅓ | ㅓ | ㅓ | ㅓ | ㅓ | ㅓ | ㅓ | ㅓ | ㅓ | ㅓ | ㅓ | ㅓ | ㅓ | |
| U+117x | ㅑ | ㅓ | ㅕ | ㅡ | ㅓ | ㅓ | ㅓ | ㅓ | ㅓ | ㅓ | ㅓ | ㅓ | ㅓ | ㅓ | ㅓ | ㅓ | ㅓ | ㅓ | |
| U+118x | ㅕ | ㅕ | ㅕ | ㅕ | ㅕ | ㅕ | ㅕ | ㅕ | ㅕ | ㅕ | ㅕ | ㅕ | ㅕ | ㅕ | ㅕ | ㅕ | ㅕ | ㅕ | |
| U+119x | ㅕ | ㅕ | ㅕ | ㅕ | ㅕ | ㅕ | ㅕ | ㅕ | ㅕ | ㅕ | ㅕ | ㅕ | ㅕ | ㅕ | ㅕ | ㅕ | ㅕ | ㅕ | |
| U+11Ax | ㅏ | ㅓ | ㅗ | ㅓ | ㅓ | ㅓ | ㅓ | ㅓ | ㅓ | ㅓ | ㅓ | ㅓ | ㅓ | ㅓ | ㅓ | ㅓ | ㅓ | ㅓ | |
| U+11Bx | ㅓ | ㅓ | ㅓ | ㅓ | ㅓ | ㅓ | ㅓ | ㅓ | ㅓ | ㅓ | ㅓ | ㅓ | ㅓ | ㅓ | ㅓ | ㅓ | ㅓ | ㅓ | |
| U+11Cx | ㅌ | ㅍ | ㅎ | ㄱ | ㄱ | ㄱ | ㄱ | ㄱ | ㄱ | ㄱ | ㄱ | ㄱ | ㄱ | ㄱ | ㄱ | ㄱ | ㄱ | ㄱ | |
| U+11Dx | ㅎ | ㅎ | ㅎ | ㅎ | ㅎ | ㅎ | ㅎ | ㅎ | ㅎ | ㅎ | ㅎ | ㅎ | ㅎ | ㅎ | ㅎ | ㅎ | ㅎ | ㅎ | |
| U+11Ex | ㆁ | ㆁ | ㆁ | ㆁ | ㆁ | ㆁ | ㆁ | ㆁ | ㆁ | ㆁ | ㆁ | ㆁ | ㆁ | ㆁ | ㆁ | ㆁ | ㆁ | ㆁ | |
| U+11Fx | ㆁ | ㆁ | ㆁ | ㆁ | ㆁ | ㆁ | ㆁ | ㆁ | ㆁ | ㆁ | ㆁ | ㆁ | ㆁ | ㆁ | ㆁ | ㆁ | ㆁ | ㆁ | |

Notes

1.[▲] As of Unicode version 13.0

2. █: Hangul jamo with a green background are modern-usage characters which can be converted into precomposed Hangul syllables under [Unicode normalization form NFC](#).

Hangul jamo with a white background are used for archaic Korean only, and there are no corresponding precomposed Hangul syllables.

"Conjoining Jamo Behavior"  (PDF). *The Unicode Standard*. March 2020.

| Number | Lead | Jamo | Character reference | Number | Vowel | Jamo | Character reference | Number | Tail | Jamo | Character reference |
|--------|------|------|---------------------|--------|-------|------|---------------------|--------|------|------|---------------------|
| 1 | G | ㄱ ㄱ | ᄀ | 1 | A | ㅏ ㅏ | ᅡ | 1 | G | ㄱ | ᆨ |
| 2 | GG | ㄲ ㄲ | ᄁ | 2 | AE | ㅐ ㅐ | ᅢ | 2 | GG | ㄲ | ᆩ |
| 3 | N | ㄴ ㄴ | ᄂ | 3 | YA | ㅑ ㅑ | ᅣ | 3 | GS | ㄳ ㄳ | ᆪ |
| 4 | D | ㄷ ㄷ | ᄃ | 4 | YAE | ㅒ ㅒ | ᅤ | 4 | N | ㄴ | ᆫ |
| 5 | DD | ㄸ ㄸ | ᄄ | 5 | EO | ㅓ ㅓ | ᅥ | 5 | NJ | ㄻ ㄻ | ᆬ |
| 6 | R | ㄹ ㄹ | ᄅ | 6 | E | ㅓ ㅓ | ᅦ | 6 | NH | ㄻ ㄻ | ᆭ |
| 7 | M | ㅁ ㅁ | ᄆ | 7 | YEO | ㅕ ㅕ | ᅧ | 7 | D | ㄷ ㄷ | ᆮ |
| 8 | B | ㅂ ㅂ | ᄇ | 8 | YE | ㅘ ㅘ | ᅨ | 8 | L | ㄹ ㄹ | ᆯ |
| 9 | BB | ㅃ ㅃ | ᄈ | 9 | O | ㅗ ㅗ | ᅩ | 9 | LG | ㄺ ㄺ | ᆰ |
| 10 | S | ㅅ ㅅ | ᄉ | 10 | WA | ዋ ዋ | ᅪ | 10 | LM | ㅙ ㅙ | ᆱ |
| 11 | SS | ㅆ ㅆ | ᄊ | 11 | WAE | ㅕ ㅕ | ᅫ | 11 | LB | ㅙ ㅙ | ᆲ |
| 12 | | ㅇ ㅇ | ᄋ | 12 | OE | ㅚ ㅚ | ᅬ | 12 | LS | ㅢ ㅢ | ᆳ |
| 13 | J | ㅈ ㅈ | ᄌ | 13 | YO | ㅙ ㅙ | ᅭ | 13 | LT | ㅚ ㅚ | ᆴ |
| 14 | JJ | ㅉ ㅉ | ᄍ | 14 | U | ㅜ ㅜ | ᅮ | 14 | LP | ㅛ ㅛ | ᆵ |
| 15 | C | ㅊ ㅊ | ᄎ | 15 | WEO | ㅘ ㅘ | ᅯ | 15 | LH | ㅘ ㅘ | ᆶ |
| 16 | K | ㅋ ㅋ | ᄏ | 16 | WE | ㅕ ㅕ | ᅰ | 16 | M | ㅁ ㅁ | ᆷ |
| 17 | T | ㅌ ㅌ | ᄐ | 17 | WI | ㅟ ㅟ | ᅱ | 17 | B | ㅂ ㅂ | ᆸ |
| 18 | P | ㅍ ㅍ | ᄑ | 18 | YU | ㅠ ㅠ | ᅲ | 18 | BS | ㅄ ㅄ | ᆹ |
| 19 | H | ㅎ ㅎ | ᄒ | 19 | EU | ㅡ ㅡ | ᅳ | 19 | S | ㅅ ㅅ | ᆺ |
| 20 | | | | 20 | YI | ㅣ ㅣ | ᅴ | 20 | SS | ㅆ ㅆ | ᆻ |
| 21 | | | | 21 | I | ㅣ ㅣ | ᅵ | 21 | NG | ㆁ ㆁ | ᆼ |
| 22 | | | | 22 | J | ㅈ ㅈ | ᆽ | | | | |
| 23 | | | | 23 | C | ㅊ ㅊ | ᆾ | | | | |
| 24 | | | | 24 | K | ㅋ ㅋ | ᆿ | | | | |
| 25 | | | | 25 | T | ㅌ ㅌ | ĜO; | | | | |
| 26 | | | | 26 | P | ㅍ ㅍ | ᇁ | | | | |
| 27 | | | | 27 | H | ㅎ ㅎ | ᇂ | | | | |

참고:

- [완성형 한글 코드](https://namu.wiki/w/%EC%99%84%EC%84%B1%ED%98%95)
- [조합형/완성형/유니코드의 모든 것, 원제: 한글 및 한국어 정보 처리 코드 \(한국어 정보처리의 과제\), 전상훈, 1999. 01. 01 konan2@chollian.net, klipl@klipl.com](https://tapito.tistory.com/529)
- [옛한글](https://namu.wiki/w/%EC%98%9B%ED%95%9C%EA%B8%80)

```
▶ M↳ ☰→☒
```

```
cv = ord('각')

basevalue = ord('ㄱ')
print(basevalue, hex(basevalue), cv)
jong = (cv - basevalue) % 28
jung = ((cv - basevalue)/28) % 21 if jong >= 28 else 0
cho = (((cv - basevalue)/28)/21)%19 if jung >= 21 else 0

cho, jung, jong
```

```
44032 0xac00 44033
```

```
(0, 0, 1)
```

```
▶ M↳ ☰→☒
```

```
a, b, c = chr(cho+0x1100), chr(jung+0x1161), chr(jong+0x11A8 - 1)
a, b, c

('ㄱ', 'ㅏ', 'ᆨ')
```

```
▶ M↳ ☰→☒
```

```
a,b,c = ord(a) - 0x1100, ord(b) - 0x1161, ord(c) - 0x11A8 + 1
syllable = 0xAC00 + a*588 + b*28 + c
print(syllable, f'0x{syllable:X}', chr(syllable))
```

```
44033 0xAC01 ㄱㅏ
```