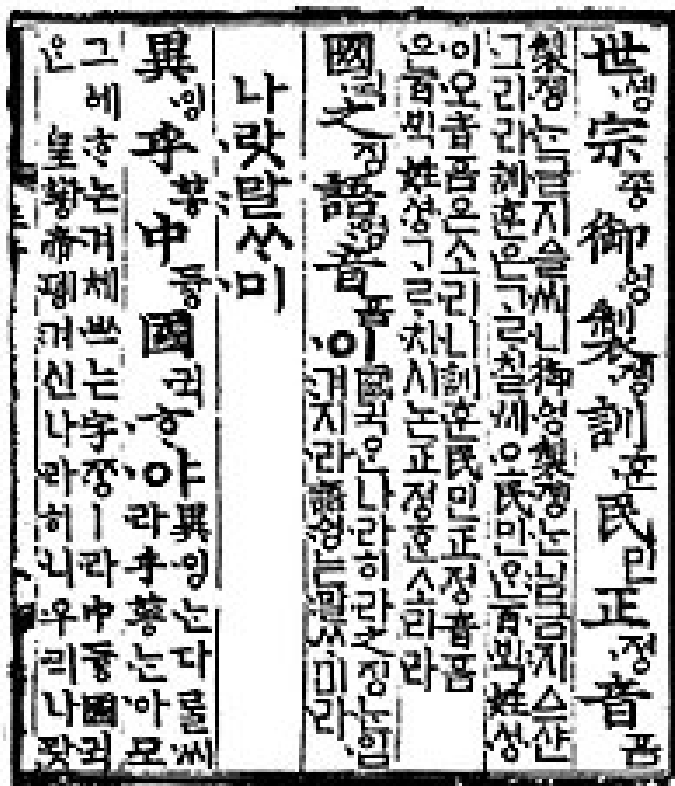


ITP20003 Java Programming

Hangul Encoding in Java



"Because the speech of this country is different from that of China, it [the spoken language] does not match the [Chinese] letters. Therefore, even if the ignorant want to communicate, many of them in the end cannot state their concerns. Saddened by this, I have [had] 28 letters newly made. It is my wish that all the people may easily learn these letters and that [they] be convenient for daily use."

<https://en.wikipedia.org/wiki/Hunminjeongeum>

Character Encoding

- A character encoding is a way to characters of languages in binary representations
 - mapping from digital code to character sets
- Computer systems use different character encodings
 - e.g., ASCII, Windows-1252, ISO/IEC-8859, UTF-8, UTF-16, KS C 5601, EUC-KR
- UTF-8 (Unicode Transformation Format using 8-bit bytes) is one of the most popular encodings (de facto standard)
 - include comprehensive characters of different languages
 - used by more than 92% of all the websites today

UTF-8 *

- Use one to four bytes (8-bits to 32 bits) which defines 1,112,064 code points
- Layout

Number of bytes	Bits for codepoint	First codepoint	Last codepoint	Byte 1	Byte 2	Byte 3	Byte 4	
1	7	U+0000	U+007F	0xxxxxxx				ASCII
2	11	U+0080	U+07FF	110xxxxx	10xxxxxx			Latin alphabets
3	16	U+0800	U+FFFF	1110xxxx	10xxxxxx	10xxxxxx		Basic Multilingual Plane (e.g., frequent CJK)
4	21	U+10000	U+10FFFF	11110xxx	10xxxxxx	10xxxxxx	10xxxxxx	Less common CJK, math symbols, emoji

* <http://www.unicode.org/>

* <https://en.wikipedia.org/wiki/UTF-8>

UTF-16 and Java

- UTF-16 uses one or two 16-bit numbers to encode total 1112064 code points
 - extends UCS-2 which is a fixed size, 16-bits encoding system
- Java represents a char as a 16-bit number whose encoding is the same as the part of UTF-16 (or UCS-2).
 - each 16-bit number for a character is called as a code point
- In addition, Java uses UTF-16 to represent Strings and streamed communications

Korean Characters in UTF-16

- All Hangul letters are in as 16-bits code points of UTF-16
- A Hangul syllable is represented as a combination of a Initial Jaeum, a Medial Moeum, and a Final Jaeum
 - 19 initial Jaeums
 - 21 Moeums
 - 27 final Jaeums (excluding the empty case)
- Code points
 - Jaeums: **U+3131** to **U+314E**
 - Ones for the initial and ones for the final are not the same
 - Moeums: **U+314F** to **U+3164**
 - **11172** syllables: **U+AC00** to **U+D7A3** ($19 \times 21 \times 28 = 11172$)
 - <http://www.programminginkorean.com/programming/hangul-in-unicode>

Example

```
public class HangulEncoding
{
    public static void main(String [] args) {
        String s = args[0] ;
        for (int i = 0 ; i < s.length() ; i++) {
            int cp = s.charAt(i) ;
            System.out.println("'" + cp + "' " + Character.toChars(cp)[0]) ;
        }
    }
}
```

Programming Exercise: HangulLinearize

- Receive a Korean sentence that consist of compound characters and then print out the linearized/normalized sequence of single-letter characters

```
$ java HangulLinearize "여기까지가 끝인가 보오"  
ㅇ ㄷ ㄱ | ㄱ ㅏ ㅈ | ㄱ ㅏ ㄱㅡㅌㅇ | ㄴ ㄱ ㅏ ㅂㅏㅇㅏ  
$
```


Readings

- Understanding Hangul Encoding
 - <https://d2.naver.com/helloworld/19187>
 - <https://d2.naver.com/helloworld/76650>