

# Chapter 3 - Unicode - boolean - identifiers

12/5/22

Unicode

## 3.3.4

code point = code value associated with a character

Unicode has 17 code planes. 1st plane is Basic Multilingual Plane. Other 16 planes hold "supplementary characters".

BMP =  $U+0000$  to  $U+FFFF$

Supp =  $U+10000$  to  $U+10FFFF$

Each character in BMP represented by a single 16-bit code unit (char)  
surrogates: BMP has a range of unused values called the surrogates area:

See UTF-16 encoding details for this

$0xD800$  to  $0xDBFF$ ,  $0xDC00$  to  $0xFFFF$  (2048 values)  
1st code unit — 2nd code unit

Supplementary characters are created by combining surrogates from each sub-range as shown above

Easy to tell if value encodes character in a single code unit or part of a supplementary character (and which surrogate)

Avoid working with chars. Work w/ code points.

## 3.3.5

boolean: false or true

Cannot convert to integer + vice versa (not directly)

## 3.4

identifiers { letters, digits, currency symbols, punctuation connectors }  
← cannot start w/ digit

ANY LANGUAGE, e.g.  $\pi$  is a letter

Don't use \$ ← intended for compiler-generated names

Use `Character.isJavaIdentifierStart()` + `isJavaIdentifierPart()` to determine what's allowed

Can't use Java keywords as names.

Can't use a single underscore '\_' as a name.

Write a program to take an int in the range  $0x10000$  to  $0x10FFFF$  and return the high + low surrogate pairs to encode it in UTF-16. Iterate over the range + print the character, the code point value in hex, and the high + low surrogate values in hex. Iterate over all BMP + supplementary characters + determine if they can be used in Java identifiers (2 whether they can be used at the beginning of the identifier).