

UTF-8

1 Code point = 1-4 code units, 1 code unit = 1 byte

- Starts with 0: 1 byte = 1 code unit → We have 7 bits left
- Starts with 110: 2 bytes = 2 code units → We have 11 bits left
- Starts with 1110: 3 bytes = 3 code units → We have 16 bits left
- Starts with 1111 0: 4 bytes = 4 code units → We have 21 bits left

Examples:

- $\underline{1110} \underline{000}$ $\underline{10011111}$ $\underline{10001101}$ $\underline{10001110}$
↓ data data data data
4 byte unit code because it starts with 11110

→ 000 011111, 001101 001110 code point

→ separate into bytes

0x0 0x1 0xF 0x3 0x4 0xE → U+01F34E

→ apple emoji (unicode) = U+01F34E, (UTF-8) = F09F8D8E

- $\underline{00101011}$
↓ data
1 byte unit code because it starts with a 0

→ 0101011 matches up with ASCII

→ ASCII value is 0x2 0xB → U+002B

→ +

0xxx xxxx

110x xxxx 10xx xxxx

1110 xxxx 10xx xxxx 10xx xxxx

1111 0xxx 10xx xxxx 10xx xxxx 10xx xxxx

- $\underline{11001111}$ $\underline{10000000}$
↓ data
2 bytes code unit

→ 0111100,0000

→ 0x3 0xC 0x0 → U+03C0

→ Lower case π

- $\underline{00100001}$
↓ data
1 byte code unit

→ 010,0001

→ 0x2 0x1 → U+0021

→ !

UTF-16

1 code point = 1-2 code units, 1 code unit = 2 bytes

- First code unit starts with 11010 → We have 11 bits left 1101 1xxx xxxx xxxx
- Second code unit starts with 110111 → We have 10 bits left 1101 11xx xxxx xxxx
- Code points under 65536 are in a single code unit

Examples:

- $\underline{11011000}$ $\underline{00110100}$ $\underline{11011101}$ $\underline{00011110}$
data data

0000,1101,0001,0001,1110

0x0 0xD 0x1 0x1 0xE → U+000D11E

unicode = U+000D11E, UTF-16 = D934DD1E

UTF-32

1 code point = 1 code unit, 1 code unit = 4 bytes

- Always uses 32 bits
- same value as the character