

'z' → ASCII 7A<sub>16</sub> → Unicode 7A

UTF-8 ~~0~~ 111 1010 → 0 111 1010  
7 data bits                      7 A

UTF-16 7A is in the BMP

00 7A

Big Endian: 00 7A

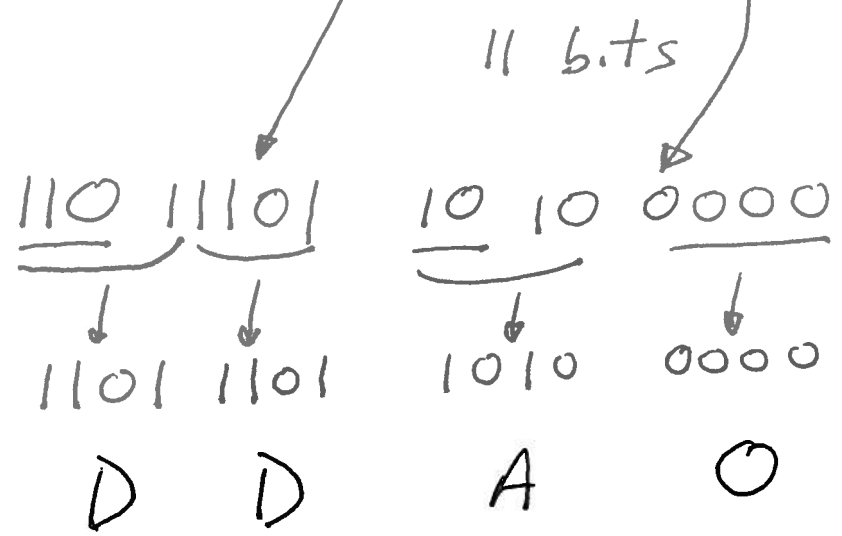
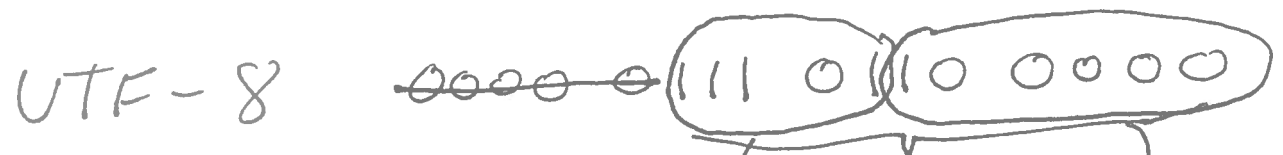
Little Endian: 7A 00

UTF-32 00 00 00 7A

Big Endian: 00 00 00 7A .....

Little Endian: 7A 00 00 00 .....

Unicode 0760<sub>16</sub> (Thousand)



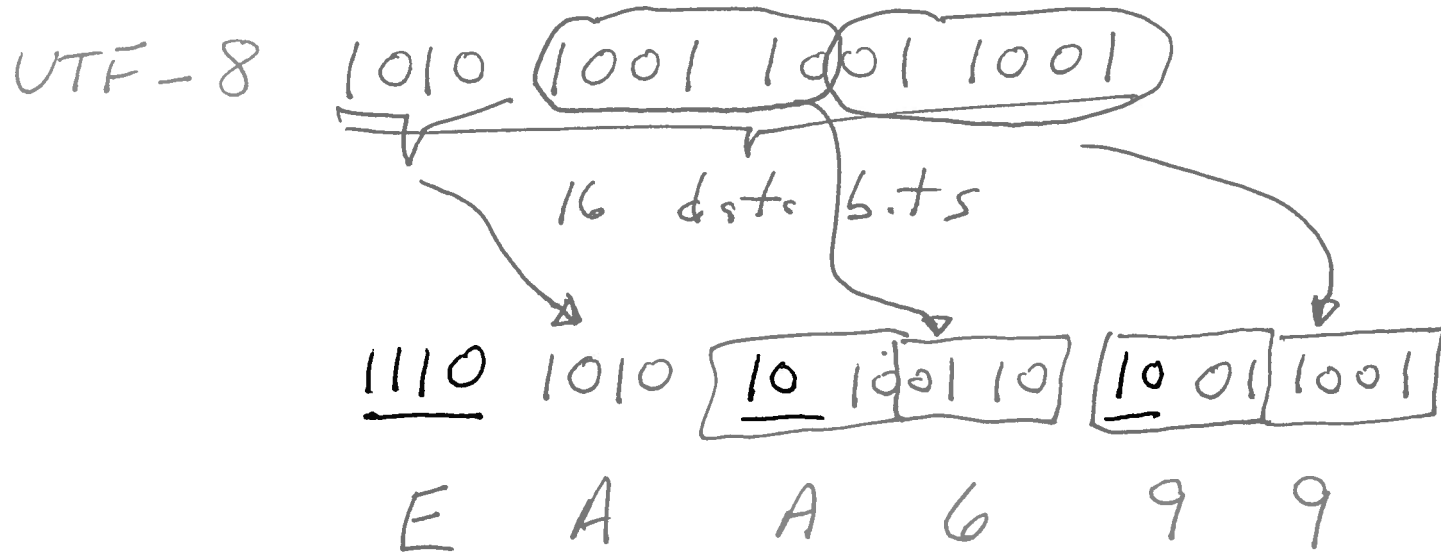
DD AO

UTF-16 0760 → Big Endian 07 60  
 Little Endian 60 07

UTF-32 00000760

Unicode A999<sub>16</sub> (Japanese)

3



UTF-16 A999 → A9 99 B.E.  
↳ 99 A9 L.E.

UTF-32 00 00 A9 99

Unicode 12345<sub>16</sub>

4

UTF-8 ~~0001~~ 0010 0011 0100 0101

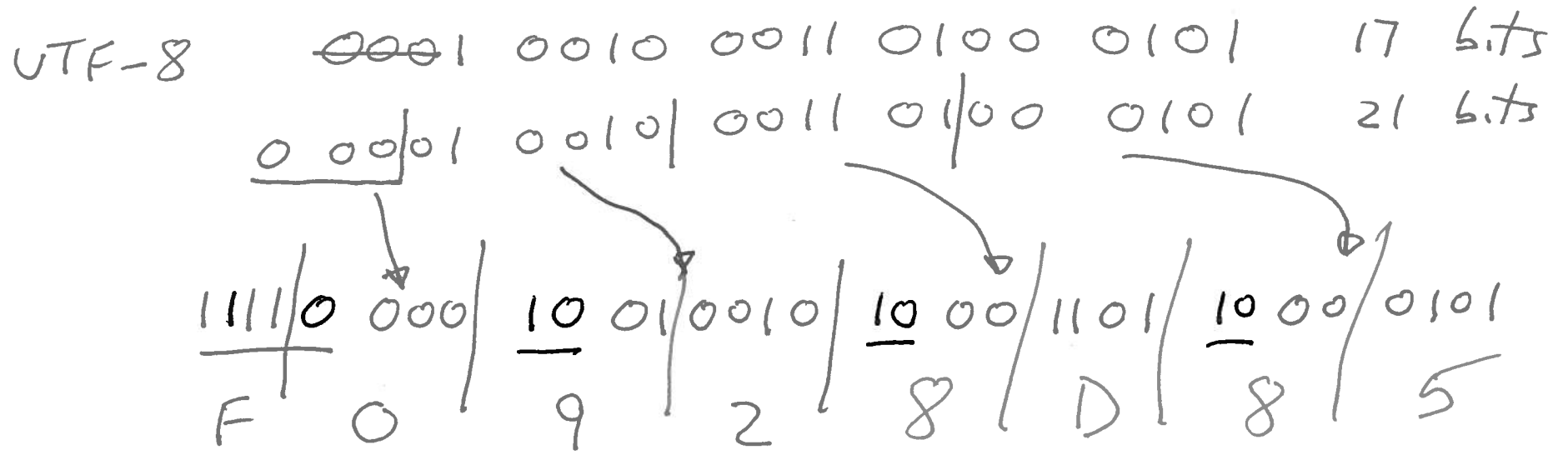
17 data bits

encode in (21) bits in 4 bytes

~~11110 100 1010 0011 10~~

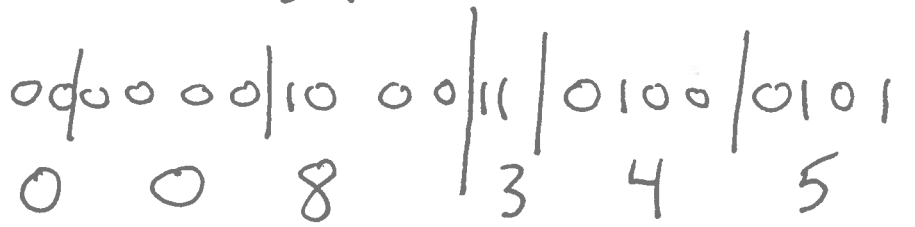
need to start over!

Unicode 12345<sub>16</sub> → UTF-32 00012345



UTF-16

1. 12345  
 - 10000  
 -----  
 02345



2. D800  
 + 008  
 -----  
 D808

DC00  
 + 345  
 -----  
 DF45

3. D808 DF45

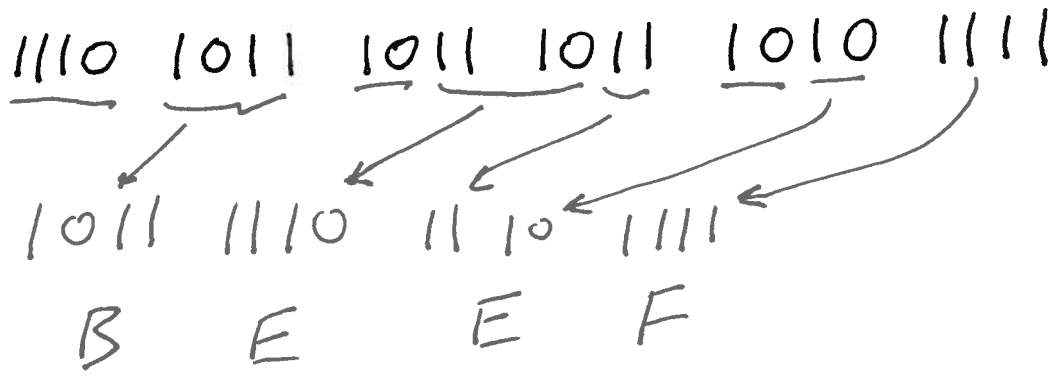
4.

n	D8	08
n+1	08	D8
n+2	DF	45
n+3	45	DF

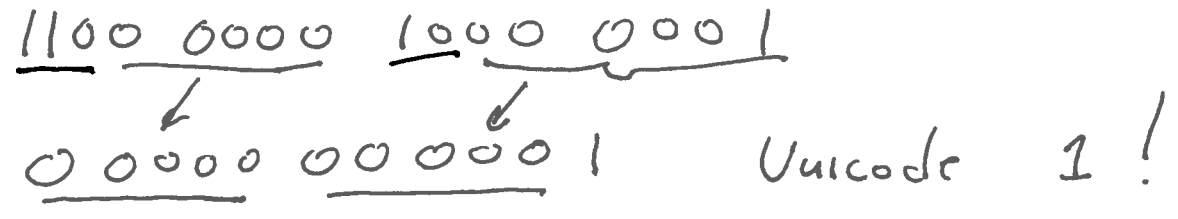
↑ ↑  
 Big Little  
 Endian Endian

UTF-8: EB BB AF

in hex



UTF-8: C0 81  
in hex

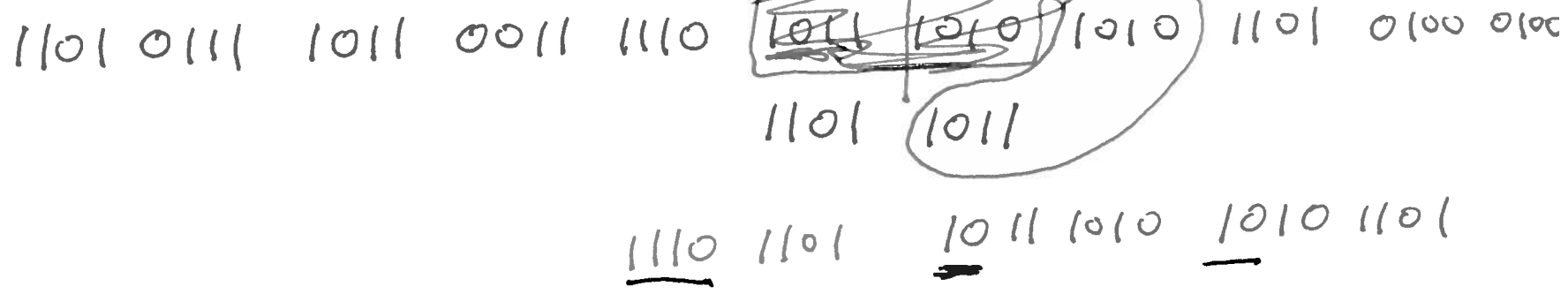


this is an error!

1 should have been encoded in 1 UTF-8 byte  
usually 01

Overlong encoding is an error!

UTF-8: D7 B3 ED BA AD 44  
in hex



self-synchronizing





BOM :

UTF-16

Big Endian  
or Little  
Endian?

FF FE ....

IF ~~B~~ Endian,

FFFE but that's  
not legal character  
it's a noncharacter

IF Little Endian, FE FF That's the BOM!

So the file is Little Endian.