

Sek 4.3 i kompendiet

Hvordan tekst representeres på en datamaskin

Faktaboks 4.9:

På en datamaskin representeres ^{1 tekst} tegn ved hjelp av heltallskoder. En tabell blir brukt for å assosiere ethvert tegn med en slik kode.

Vi trenger tre ting:

1. Et tegnsett (A, B, C, ..., a, b, c, ..., 1, 2, 3, ...)

2. En tabell med en kode for hver tegn
%, \$, #

3. En regel for hvordan maskinen skal representere koden ved hjelp av bytes. (1 byte = 8 bits)
→ code point i kompendiet.

ASCII-tegnsettet:

128-tegn
kodene fra 0 til 127 (trenger 7 bits)
maskinen bruker 1 byte per tegn, setter det
mest signifikante bitet til 0. $(\underbrace{0111111}_{127})$

A, B, C, ...	:	65, 66, ...	- 90
a, b, c, ...	:	97, 98, ...	- 122
0, 1, 2, ..., 9	:	48 - 57	
#, \$, %, &	:	35, 36, 37, 38	

Legg merke til at Q , ø , å IKKE er med i ASCII-tegnsettet!

Trenger et større tegnsett for at vi skal kunne skrive dokumenter på norsk.

ISO - Latin tegnsett (definert på 80-tallet) har forskjellige varianter for hvert språk inneholder 256 tegn

ISO-Latin 1 (western) inneholder Q , ø , å

197	Å	229	å
198	Æ	230	Q
216	Ø	248	ø

kodeene for 0-127 er de samme som i ASCII

90-tallet: Firmaet Unicode definerte tegnsett som inneholder de fleste tegn for de fleste språk. 100 000 tegn!

Unicode - tabellen: Koder fra 0 til 1114 111 bare 10% av kodene for 1114 111 er i bruk.

Gjør det lett å legge til tegn senere. 10ffff_{16}

ASCII, ISO-LATIN 1: Tegnet representeres med 8 bits = 1 byte ved hjelp av koden i tabellen.

Alle tegn i Unicode-tabellen kan representeres med 21 bits, og vi trenger da å bruke 3 bytes.

I praksis ville vi koste borte plass hvis vi allokerte 3 bytes for alle tegn. Derfor har man laget "enkodinger" som kan kode tegn, ved å bruke forskjellig antall bytes per tegn

↓
regel for å skrive en kode om til et sett av bytes.

4.3.4 UTF-8 enkoding (av Unicode - tegnsettet).

Regel: Første bits i enhver byte er en av

0 , 10 , 110 , 1110 , 11110
 ↓
 ASCII
 ↓
 byte 2, 3 eller 4
 i tegn som blir rep.
 med flere bytes.
 ↓
 første bits når
 vi bruke 2 bytes (3, 4)
 til å rep. tegn

Eksempel 4.12

norsk Å : har Unicode-kode : 197

$$197 = 128 + 64 + 4 + 1 = \frac{1100}{2^7} \frac{0101}{2^0} = \frac{1100}{16} \frac{0101}{5} = C5_{16}$$

vi ser at $128 \leq 197 \leq 2047$, slik at regel 2 for UTF-8 enkoding skal brukes (\Rightarrow to bytes)

regel 2: vi representerer med 11 bits:

$$\begin{array}{cccc} 000 & 1100 & 0101 & \\ d_{10} d_9 & \dots & & d_0 \end{array}$$

UTF-8 sier at dette skal kodes som

$$110 d_{10} d_9 d_8 d_7 d_6 \quad 10 d_5 d_4 d_3 d_2 d_1 d_0 = \frac{11000011}{8+4=12=C_{12}} \frac{10000101}{8_{16} \quad 5_{16}} = \underline{\underline{C385_{16}}}$$

Viktigst for dere å huske:

1. UTF-8 koder med 1, 2, 3, eller 4 bytes.
2. UTF-16 koder med 2, eller 4 bytes.

Kontrollspørsmål: Hvis en tekst er kodet med 3413 bytes. Hvilke av UTF-8 eller UTF-16 kan denne ha vært kodet med?

python: # coding = UTF-8 /
 Eps 4.4 4b be 75 74
 Knu t
 Mørken

ø kan ikke lagres som ren ascii.

latin ø har unicode-kode f8
 (hvis koder med latin1, så skriver u f8 for ø).

UTF-8: f8 = 1111 1000 = 000 111 1000

⇒ kodes som $\underbrace{11000011}_c \underbrace{10111000}_b = c3b8_{16}$

når koder med UTF-8: ø kodes som c3b8₁₆
 eller kodes alt som med ISO-latin1.

UTF-16: kodes alt som det står; med to bytes.

ø (f8) kodes som 00f8.