

Representasjon av tegn og tekst

se gjen. 4.3

Tegn representeres i i komp.
en datamaskin vha. heltall som
peker inn i en tegn tabell.

ASCII-tabellen, 7 bits, 128 tegn.

ISO LATIN - 8 bits, 256 tegn

ISO LATIN1 inneholder norske og andre
vest europeiske spesialtegn

Unicode - tabell som inneholder alle verdens tegn. Unicode har plass til ca 10^6 tegn, men bare ca. 100000 er i bruk.

Hvis Unicode-kodene brukes direkte må hvert tegn ha 3 bytes.

En Unicode-kode kalles et kodepunkt.

I UTF-8 kodes kodepunktene på en slik måte at ASCII-tegn trenger en byte, andre latinske tegn 2 bytes, kinesiske tegn 3 bytes, og noen andre 4 bytes.

Tapstri kompresjon

I den løse kompresjon er at de hyppigste tegnene får de korteste kodene.

Størrelse på filer.

Eks Lyd. På CD.

Lydsignalet måles 44100 nr. sekund i stereo, og hver måling lagres som et heltall med 2 bytes.

$$44100 \cdot 2 \cdot 2 = 176400 \text{ bytes nr. sekund}$$

10MB = 10 000 000 bytes nr. min

ca 40 MB for en låt, 600 MB på en CD.

Eks film. Klassisk TV-signal (PAL)

$$112 \text{ GB} = 112 \cdot 10^9 \text{ bytes nr. time.}$$

Stort behov for kompakt lagring.

Huffman koding seksjon 7.2 kompendiet.

$\bar{x} = \{x_1, x_2, \dots, x_m\}$ - tekst

Hver x_i er tatt fra et alfabet A ,
og $A = \{\alpha_1, \alpha_2, \dots, \alpha_n\}$.

Antall forekomster av α_i i \bar{x} er
frekvensen $f(\alpha_i)$.

Kompressjon: finn kode $c(\alpha_i)$ for α_i
og lagre x ved at hver x_j lagres
ved sin kode.

Ex. Anta at $\bar{x} = DBACDBD$

$f(A) = 1, f(B) = 2, f(C) = 1, f(D) = 3$

Kort kode til vanligste tegn.

$c(D) = 0, c(B) = 1, c(C) = 01, c(A) = 10$

Lagre kompent: $Z = 011001010$

Ex. Annen koding

$c(D) = 1, c(B) = 01, c(C) = 001, c(A) = 000$

Da blir $Z = 1010000011011$

Koder må ha prefixegenskapen,

En kode må ikke være starten på en
annen kode.

Binær træ

Binært træ T

