

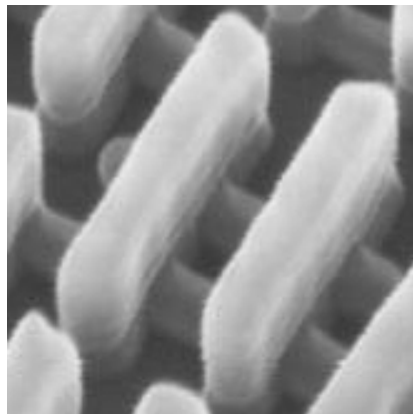
## Digital representasjon

*Dag Langmyhr*  
dag@ifi.uio.no

Hvordan lagre

- tall
- tekst
- bilder
- lyd

som **bit** i en datamaskin



## Binære tall

For å bruke bit (0 og 1) som tall, må vi telle binært. Dette gjøres egentlig på samme måte som vi teller desimalt:

- Øk siste siffer med 1.
- Hvis vi ikke har flere siffer, sett til 0 og gjenta for sifferet til venstre.

Binært	0	1	10	11	100	101	110	111	1000	1001	1010	1011
Desimalt	0	1	2	3	4	5	6	7	8	9	10	11

## Notasjon

Hvis det mulighet for tvil, skriver vi binære tall som  $1001_2$  og desimale tall som  $1001_{10}$ .

## Et matematisk blikk

I det desimale tallsystemet har posisjonene vekt  
 $1, 10, 100, 1000, \dots = 10^0, 10^1, 10^2, 10^3, \dots$

(Notasjonen  $a^n$  (« $a$  i  $n$ -te») betyr  $\underbrace{a \times a \times \dots \times a}_{n \text{ ganger}}$ .)

I det binære tallsystemet har posisjonene vekt  
 $1, 2, 4, 8, \dots = 2^0, 2^1, 2^2, 2^3, \dots$

$$\begin{array}{cccc}
 8 & 4 & 2 & 1 \\
 2^3 & 2^2 & 2^1 & 2^0 \\
 \downarrow & \downarrow & \downarrow & \downarrow \\
 1011_2 = & 1 & 0 & 1 & 1 & = 11_{10}
 \end{array}$$

Vi teller binært med grupper av bit, f eks **byte** som inneholder 8 bit:

0	0	0	0	0	0	0	0	=	$0_{10}$	
0	0	0	0	0	0	0	0	1	=	$1_{10}$
0	0	0	0	0	0	0	1	0	=	$2_{10}$

⋮

0	1	1	1	1	1	1	1	0	=	$126_{10}$
0	1	1	1	1	1	1	1	1	=	$127_{10} = 2^7 - 1$

Vi stopper når vi kommer til øverste bit.

## Negative tall

Et negativt tall  $n$  lagres i 8 bit som  $2^8 + n = 256 + n$ :

1	1	1	1	1	1	1	1	=	$-1_{10}$	=	$256 - 1 = 255$
1	1	1	1	1	1	1	0	=	$-2_{10}$		
1	1	1	1	1	1	0	1	=	$-3_{10}$		

⋮

1	0	0	0	0	0	0	0	1	=	$-127_{10}$
1	0	0	0	0	0	0	0	0	=	$-128_{10}$

Øverste bit («fortegnsbitet») angir om et tall er negativt.

Hvor store tall kan vi da lagre?

## Heltall i Java

Java tilbyr disse tallene:

<b>byte</b>	8 bit	-128 - +127
<b>short</b>	16 bit	-32 768 - +32 767
<b>int</b>	32 bit	-2 147 483 648 - +2 147 483 647
<b>long</b>	64 bit	-9 223 372 036 854 775 808 - +9 223 372 036 854 775 807

Hvor store tall kan vi da lagre?

Hvorfor er dette viktig?

**Overflyt.java**

```
class Overflyt {  
    public static void main(String arg[]) {  
        int v = 1000000, v2 = v*v;  
  
        System.out.println("v=" + v + " og v2=" + v2);  
    }  
}
```

gir dette resultatet

```
$ javac Overflyt.java  
$ java Overflyt  
v=1000000 og v2=-727379968
```

Hvor store tall kan vi da lagre?

### IkkeOverflyt.java

```
class IkkeOverflyt {  
    public static void main(String arg[]) {  
        long v = 1000000, v2 = v*v;  
  
        System.out.println("v=" + v + " og v2=" + v2);  
    }  
}
```

gir riktig svar:

```
$ javac IkkeOverflyt.java  
$ java IkkeOverflyt  
v=1000000 og v2=1000000000000
```



## Hva så med Python?

Hovedregel: Python jonglerer med lagring av tall, så overflyt forekommer ikke.

Men: De fleste biblioteker man bruker i Python, er skrevet i C, og der har man de samme problemene som i Java.

## Heksadesimal notasjon

Det er lett å gjøre feil når man jobber med binære tall:

```
11111110110111001011101010011000...
...01110110010101000011001000010000
```

Det er enklere å erstatte fire og fire binære sifre med **heksadesimale** sifre 0x0-0xF:

<u>1111</u>	<u>1110</u>	<u>1101</u>	<u>1100</u>	<u>1011</u>	<u>1010</u>	<u>1001</u>	<u>1000</u>
0xF	0xE	0xD	0xC	0xB	0xA	0x9	0x8
<u>0111</u>	<u>0110</u>	<u>0101</u>	<u>0100</u>	<u>0011</u>	<u>0010</u>	<u>0001</u>	<u>0000</u>
0x7	0x6	0x5	0x4	0x3	0x2	0x1	0x0

## Oktal notasjon

Tidligere brukte man ofte **oktal** notasjon der man slår sammen tre og tre bit:

$$\begin{array}{cccccccc} \underbrace{111}_{07} & \underbrace{110}_{06} & \underbrace{101}_{05} & \underbrace{100}_{04} & \underbrace{011}_{03} & \underbrace{010}_{02} & \underbrace{001}_{01} & \underbrace{000}_{00} \end{array}$$

(I Java skriver man oktale tall som **0nnn**; i Python 2 skriver man dem også som **0nnn**, men i Python 3 bruker man **0onnn**.)

I dag har hex-notasjonen overtatt.

## Men hva med disse tallene?

- Andromedagalaksen er 24 029 742 100 000 000 000 km unna.
- $\pi = 3,14159265$
- Et H-atom er 0,000 000 096 mm stort.

Vi trenger altså ikke bare heltall men også tall som

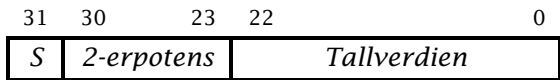
- kan ha veldig små og veldig store verdier
- kan ha desimaler
- ikke behøver å være helt nøyaktige.

Løsningen i det desimale tallsystemet er å oppgi tallet med 10-erpotens:

- Andromedagalaksen er  $2,4 \cdot 10^{19}$  km unna.
- $\pi = 3,14159265$
- Et H-atom er  $9,6 \cdot 10^{-8}$  mm stort.

Så lagrer vi både den justerte tallverdien (9,6) og 10-erpotensen (-8).

På en datamaskin gjør vi det samme, men bruker vi 2-erpotenser i f eks float:



I Java har vi disse flyttallene:

			<b>Minste</b>	<b>Største</b>
<b>float</b>	32 bit	7 sifre	$1,2 \cdot 10^{-38}$	$3,4 \cdot 10^{38}$
<b>double</b>	64 bit	16 sifre	$2,2 \cdot 10^{-308}$	$1,8 \cdot 10^{308}$

Python har bare **double** men kaller den en **float**.

Vi kan konvertere mellom heltall og flyttall i Java:

```
class Konvertering { Konvertering.java
    public static void main(String arg[]) {
        double f1 = 3.94, f2;
        int i;

        i = (int)f1;  f2 = (double)i / 2;
        System.out.println("i=" + i + " og f2=" + f2);
    }
}
```

Svaret blir: i=3 og f2=1.5

Det er tilsvarende i Python:

**konvert.py**

```
f1 = 3.94  
i = int(f1)  
f2 = float(i) / 2  
print "i =", i, "og f2 =", f2
```

Svaret blir:  $i = 3$  og  $f2 = 1.5$



## Er dette viktig?

`FloatTest.java`

```
class FloatTest {  
    public static void main(String arg[]) {  
        double a = 1.000000000000000000000001,  
               b = 1.000000000000000000000008;  
  
        System.out.println("a-b = " + (a-b));  
    }  
}
```

fungerer slik:

```
$ javac FloatTest.java  
$ java FloatTest  
a-b = 0.0
```

`floattest.py`

```
a = 1.000000000000000000000001  
b = 1.000000000000000000000008  
print "a-b =", a-b
```

fungerer slik:

```
$ python floattest.py  
a-b = 0.0
```



Kan vi lagre tegn som tall?

## Tekster

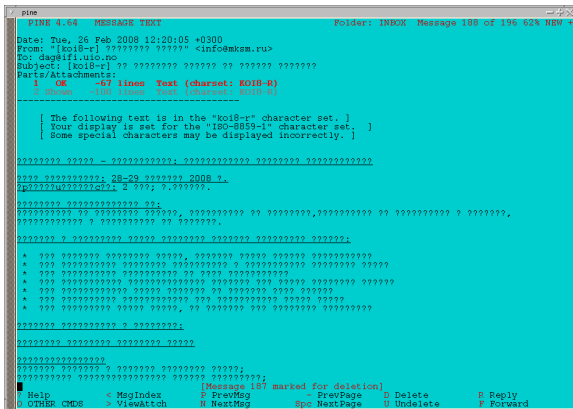
Hvordan lagrer vi tegn i datamaskinen? Det er lett:

Sett opp en tabell over tegn vi trenger og gi hvert tegn et nummer.

Er det virkelig så enkelt?

## Problemet

Hva hvis ikke alle bruker samme tabell?



```
pine
-----
FINE 4.64    MESSAGE TEXT    Folder: INBOX    Message 187 of 198 624 NEW
Date: Tue, 26 Feb 2008 12:20:05 +0300
From: "[ko18-r] ????????? ??????" <info@kms.ru>
To: dag@ifi.uio.no
Subject: [ko18-r] ?? ????????? ?????? ?? ?????? ??????
Parts/Attachments:
  1 OK    -67 lines    Text (charset: KOI8-R)
  2 image    4160 lines    Text (charset: KOI8-R)
-----
[ The following text is in the "ko18-r" character set. ]
[ Your display is set for the "ISO-8859-1" character set. ]
[ Some special characters may be displayed incorrectly. ]

????????? ?????? - ??????????????; ?????????????? ?????????? ??????????????

???? ?????????????? 28-29 ????????? 2008 ?
?p????????????????? 2 ???; ?;????????

????????? ?????????????? ??:
????????????? ?? ?????????? ??????, ?????????? ?? ??????????,????????????? ?? ?????????????? ? ???????,
????????????? ? ?????????????? ?? ??????????

????????? ? ?????????????? ?????? ?????????? ?????????? ?????????? ??????????

* ??? ?????????? ?????????? ??????, ?????????? ?????? ?????????? ??????????????
* ??? ?????????????? ?????????????? ?????????????? ? ?????????????? ?????????? ??????
* ??? ?????????????? ?????????????? ?? ?????? ??????????
* ??? ?????????????? ?????????????? ?????? ?????????? ?? ?????????? ?????? ??????
* ??? ?????????????? ?????????????????? ??? ?????????????? ?????? ??????
* ??? ?????????????? ?????? ??????, ?? ?????????? ??? ?????????? ??????????

????????? ?????????????? ? ??????????:
????????????? ?????????????? ??????

?????????????????????
????????? ?????????? ? ?????????? ?????????? ??????;
????????????? ?????????????????????????????? ?????????? ??????????????
■ Help    < MsgIndex    P PrevMsg    - PrevPage    D Delete    R Reply
0 OTHER CMDS    > ViewAttach    N NextMsg    Spc NextPage    U Undelete    F Forward
```



Mange gode forsøk

## ASCII

Første vellykkede forsøk på standardisering var ASCII (American standard code for information interchange) i 1963:

ASCII Code Chart

	0	1	2	3	4	5	6	7	8	9	A	B	C	D	E	F
0	NUL	SOH	STX	ETX	EOT	ENQ	ACK	BEL	BS	HT	LF	VT	FF	CR	SO	SI
1	DLE	DC1	DC2	DC3	DC4	NAK	SYN	ETB	CAN	EM	SUB	ESC	FS	GS	RS	US
2		!	"	#	\$	%	&	'	( )	*	+	,	-	.	/	
3	0	1	2	3	4	5	6	7	8	9	:	;	<	=	>	?
4	@	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O
5	P	Q	R	S	T	U	V	W	X	Y	Z	[	\	]	^	_
6	`	a	b	c	d	e	f	g	h	i	j	k	l	m	n	o
7	p	q	r	s	t	u	v	w	x	y	z	{		}	~	DEL

- + ASCII ble etter hvert brukt av de fleste.
- + Det er lett å sjekke om et tegn er et siffer eller en bokstav.
- + Det er lett å konvertere fra liten til stor bokstav.
- Mangler ÆØÅ og andre bokstaver og tegn.

Det siste ga opphav til et utall av lokale varianter der for eksempel [\\|}] ble erstattet av **ÆØÅæøå**. Dette førte til «Gj|vik-syndromet».

Mange gode forsøk

## Latin-1

En klart bedre løsning ble ISO 8859-1 (Latin-1).

Den (og varianter) er ennå i bruk.

## ISO 8859-1

000000	@	000000	.	000000	A	000000	à
000001	!	000001	A	000001	a	000001	á
000002	"	000002	B	000002	b	000002	â
000003	#	000003	C	000003	c	000003	ã
000004	\$	000004	D	000004	d	000004	ä
000005	%	000005	E	000005	e	000005	å
000006	&	000006	F	000006	f	000006	æ
000007	'	000007	G	000007	g	000007	ç
000008	(	000008	H	000008	h	000008	è
000009	)	000009	I	000009	i	000009	é
000010	*	000010	J	000010	j	000010	ê
000011	+	000011	K	000011	k	000011	ë
000012	,	000012	L	000012	l	000012	ì
000013	-	000013	M	000013	m	000013	í
000014	.	000014	N	000014	n	000014	î
000015	/	000015	O	000015	o	000015	ï
000016	0	000016	P	000016	p	000016	ð
000017	1	000017	Q	000017	q	000017	ñ
000018	2	000018	R	000018	r	000018	ò
000019	3	000019	S	000019	s	000019	ó
000020	4	000020	T	000020	t	000020	ô
000021	5	000021	U	000021	u	000021	õ
000022	6	000022	V	000022	v	000022	ö
000023	7	000023	W	000023	w	000023	÷
000024	8	000024	X	000024	x	000024	ø
000025	9	000025	Y	000025	y	000025	ù
000026	:	000026	Z	000026	z	000026	ú
000027	;	000027	[	000027	{	000027	û
000028	<	000028	\	000028		000028	ü
000029	=	000029	]	000029	}	000029	ý
000030	>	000030	^	000030	~	000030	ÿ
000031	?	000031		000031		000031	

## Unicode

Men det er bare én løsning på sikt: en tegnkoding som omfatter alle skriftspråk i verden. Den heter **Unicode** og er nå stort sett ferdig.

- Unicode skal omfatte alle skriftspråk som brukes eller har vært brukt

π Я 音 æ∞

- Det er plass til drøyt 1 000 000 tegn.
- 110 187 tegn er foreløbig definert.
- Mer og mer programvare (som Java og Python) støtter Unicode.



## UTF-8

Hvordan lagre Unicode-tekst uten å bruke for mye plass på disk eller over nettet? **UTF-8** bruker fra 1 til 4 byte til å lagre et tegn:

\$	U+0024	<u>00</u> 100100
¢	U+00A2	<u>11</u> 000010 <u>10</u> 100010
€	U+20AC	<u>11</u> 100010 <u>10</u> 000010 <u>10</u> 101100
⚡	U+10384	<u>11</u> 110000 <u>10</u> 010000 <u>10</u> 001110 <u>10</u> 000100

Unicode og UTF-8 er nå standard ved Ifi.

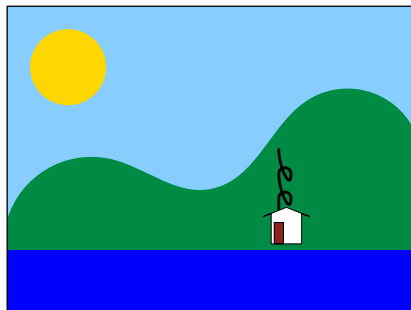


Hva er et bilde?

## Bilder

Jeg har laget en vakker tegning:

Hvordan kan jeg lagre den som bit i en datamaskin?

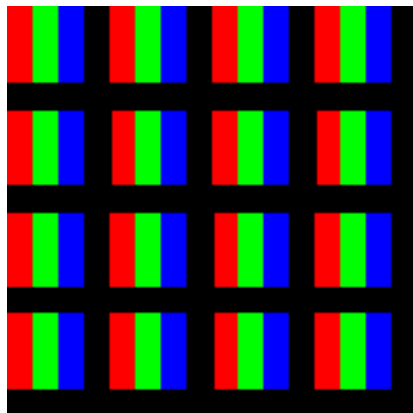


## Hva er et bilde?

På en fargeskjerm består hvert bildepunkt («pixel») av tre farger

- Rød
- Grønn
- Blå

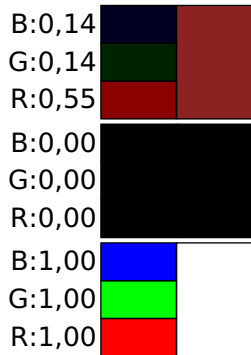
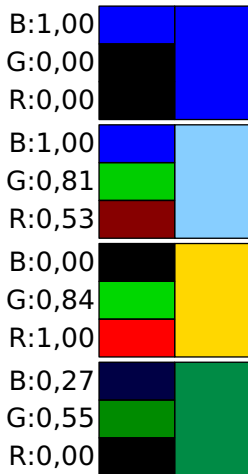
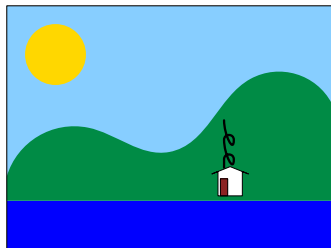
som kan lyse sterkt eller svakt.



(Utskrift på papir bruker **CMYK** (Cyan, Magenta, Yellow, black) i stedet.)

Hva er et bilde?

Hvilke farger trenger vi?

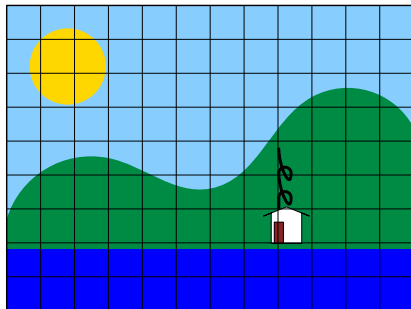


Hva er et bilde?

Da er det bare å legge et rutenett på bildet og notere mengden R, G og B i hver rute. Om vi bruker 1 byte til hver fargemengde, får vi:

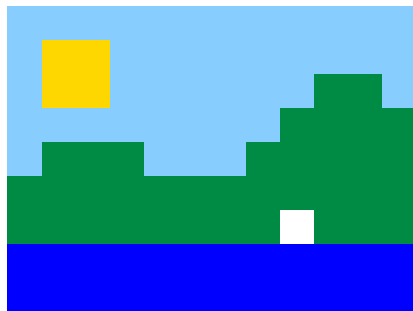
```

135 206 255
135 206 255
135 206 255
:
0 0 255
    
```



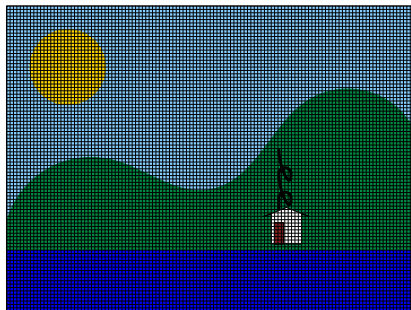
Hva er et bilde?

... og bildet er lagret:



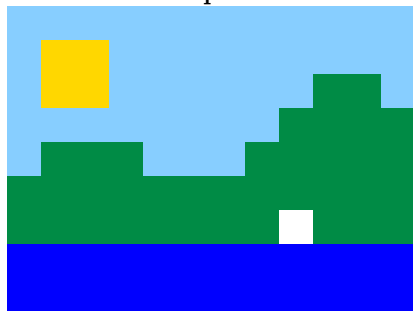
Hva er et bilde?

Vi bør nok prøve med et mer finmasket rutenett:



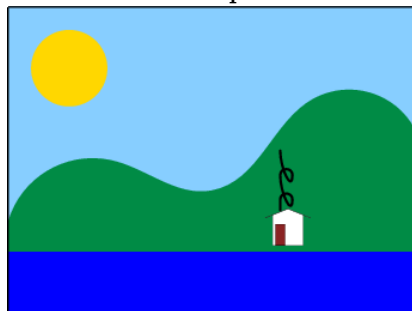
Hva er et bilde?

12×9 piksler



324 byte

348×257 piksler



363 682 byte

Kan vi spare plass?

## Kan vi lagre bildet på mindre plass?

### Fargetabell

Vi bruker 3 byte = 24 bit til hvert piksel, men vi har bare 7 ulike farger i bildet. Da kan vi sette opp en tabell

0	0 0 255	mørkeblå
1	135 206 255	lyseblå
2	255 214 0	gul
3	0 140 69	skoggrønn
4	140 36 36	brun
5	0 0 0	svart
6	1 1 1	hvit

og da trenger vi bare 3 bit per pixel.

363 682 byte  $\Rightarrow$  45 483 byte



## «Run-length»-koding

I mange rasterbilder er det ofte mange like piksler på rad. Da kan vi lagre fargen og hvor mange slike piksler det er på rad.

363 682 byte  $\Rightarrow$  45 483 byte  $\Rightarrow$  2 917 byte

Formater som PNG og GIF bruker slike teknikker.

For fotografier gjelder dette:

- På fotografier er det sjelden brå overganger.
- Vi mennesker kan bare skjelne et begrenset antall nyanser.
- Vi søker automatisk etter mønstre.



## JPEG-formatet

JPEG benytter dette til å lage en forenklet versjon av bildet.

Ekte rasterbilde	40,7 MB
JPEG 100%	5,5 MB
JPEG 50%	0,94 MB
JPEG 25%	0,61 MB
JPEG 10%	0,34 MB
JPEG 5%	0,24 MB

(Dette er komprimering med **tap**. Det er umulig å komme tilbake til det opprinnelige bildet.)

Oversikt



Tall



Tekst



Bilder



Lyd



Slutt



Originalbildet 100%



Kvalitet 50%











## En sammenligning



Kvalitet 100%



Kvalitet 5%

## Vektorgrafikk

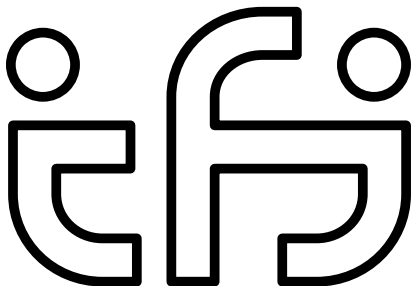
Det er også mulig å lagre bilder som linjer og kurver:

ifl-logo.eps

```

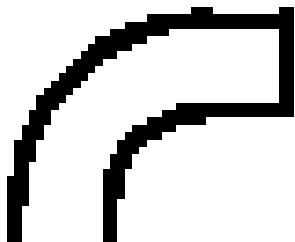
:
newpath
38.4094 2.83464 moveto
38.4094 15.44878 lineto
28.3464 15.44878 lineto
21.03242 15.44878 14.93857
21.1286 14.93857 28.3464 curveto
14.93857 35.77315 lineto
36.56685 35.77315 lineto
36.56685 48.38728 lineto
2.32443 48.38728 lineto
2.32443 28.3464 lineto
2.32443 14.16295 14.0667
2.83464 28.3464 2.83464 curveto
closepath stroke
:

```

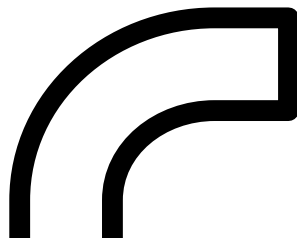


Men hva med vektorer?

Raster



Vektor



## Anbefalinger

- Bruk vektorgrafikk om mulig: SVG, EPS, PDF.
- For fotografier bruk JPEG i så god kvalitet som mulig.
- For annen rastergrafikk bruk PNG med så mange piksler som mulig.

Hva er lyd?

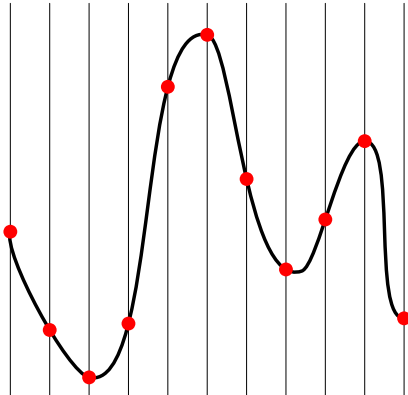
## Hvordan lagre lyd

Lyd er bølger i luft, men de kan overføres som strøm:

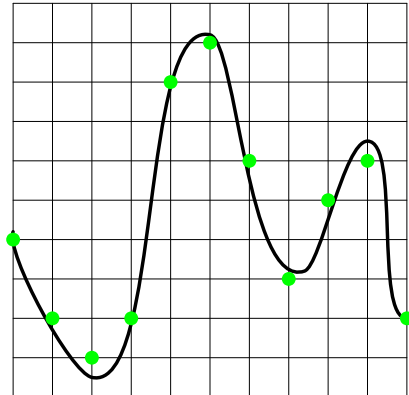


Hvordan kan vi lagre lyd?

Vi kan lagre lyden ved å måle strømmen med jevne mellomrom:

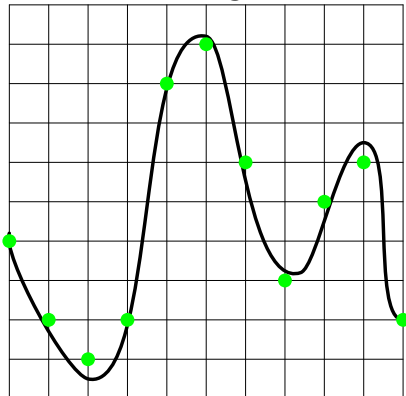


Men for å lagre digitalt må vi måle styrken i faste intervaller:

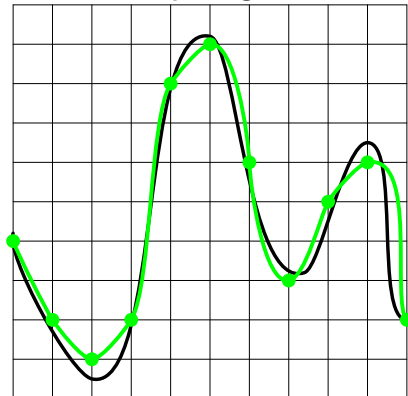


Hvordan kan vi lagre lyd?

Med disse målingene



kan vi gjenskape lyden (men ikke helt nøyaktig):



Kvaliteten på lyden blir da avhengig av

- hvor ofte vi måler
- hvor mange trinn vi benytter til målingen

## CD

En vanlig CD har 74 minutter spilletid med god lyd:

- 44 100 målinger («samples») per sekund
- $2^{16} = 65\,536$  intervaller (dvs 2 byte)
- 2 kanaler
- + 37% feilkorreksjonsdata

En CD må derfor ha plass til 783 MB.



Er det mulig å spare plass?

- Høyre og venstre kanal er stort sett nesten like. Det er lurere å lagre venstre kanal samt forskjellen.
- I stedet for å lagre hver måling med sin verdi, holder det å lagre forskjellen.

Men: På grunn av feil og spoling må man av og til lagre den ekte verdien.

Enda mer plass kan vi spare om vi tar hensyn til hvordan vi mennesker hører:

- Vi kan ikke høre lyder under 20 Hz og over 20 000 Hz.
- Om vi hører en kraftig lyd med én frekvens, hører vi ikke litt svakere lyder med noe høyere frekvens.
- Etter å ha hørt en sterk lyd, hører vi dårligere en tid etterpå (inntil 0,2 s).

Moderne standarder som

**MP2** (brukt i DAB)

**MP3**

**AAC** (brukt i DAB+, YouTube, iTunes, ...)

utnytter dette og tillater til dels sterk komprimering:

Ekte CD	1410 kb/s
MP2	ca 256 kb/s
MP3	ca 160 kb/s
AAC	ca 128 kb/s

men kvaliteten går ned om det komprimeres ennå mer.

## Oppsummering

- Alt er bit i en datamaskin.
- Det finnes ulike typer tallverdier, og programmereren må velge riktig.
- Det er mange ulike tegnkodinger å forholde seg til (ennå).
- Rasterbilder og vektorbilder er nyttige til hvert sitt formål.
- Lydkoding med MP2, MP3 og AAC er blitt standard, men vi bør velge rett kvalitet.