

# IE1205 Digital Design:

## F6 : Digital aritmetik 2

- Ett tal kan representeras binärt på många sätt.
- De vanligaste taltyperna som skall representeras är:
  - Heltal, positiva heltal (eng. integers)
    - ett-komplementet, två-komplementet, sign-magnitudo
  - Decimala tal med fix tal-område
    - Fix-tal (eng. Fixed point)
  - Decimala tal i olika talområden
    - Flyt-tal (eng. Floating point)

# Heltal



Positiva Heltal:

$-2^7$	$2^6$	$2^5$	$2^4$	$2^3$	$2^2$	$2^1$	$2^0$
0	1	1	0	1	1	0	1

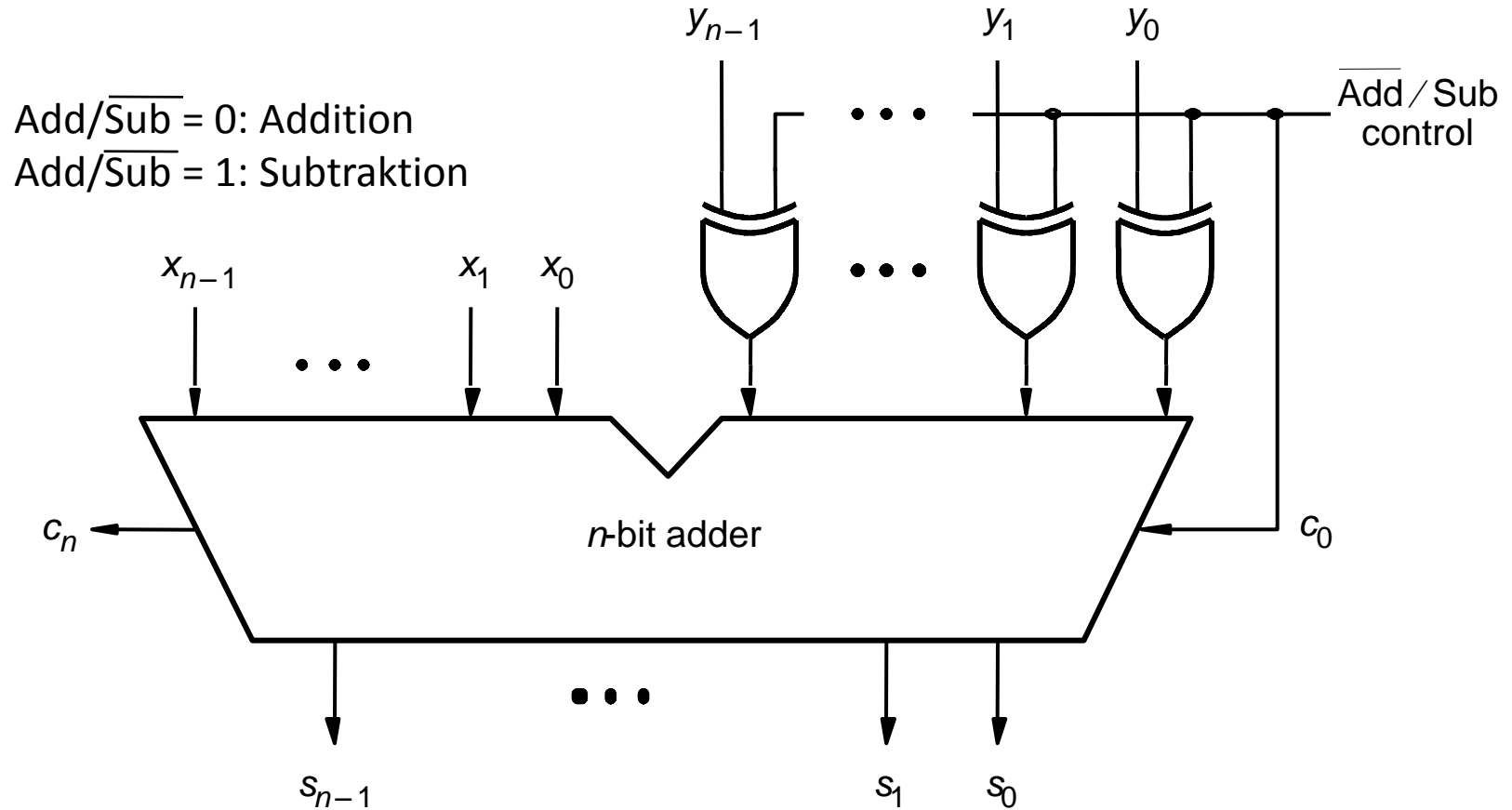
$$=1*2^6 + 1*2^5 + 1*2^3 + 1*2^2 + 1*2^0 = 109$$

Negativa Heltal:

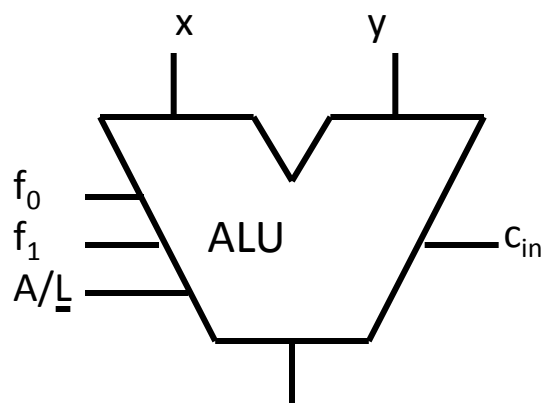
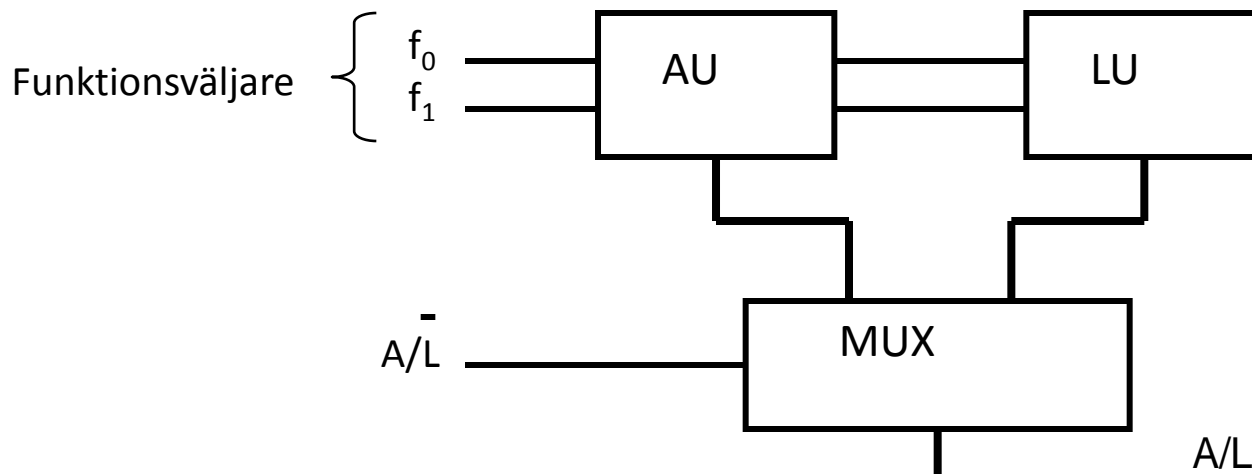
$-2^7$	$2^6$	$2^5$	$2^4$	$2^3$	$2^2$	$2^1$	$2^0$
1	1	1	0	1	1	0	1

$$=-1*2^7 + 1*2^6 + 1*2^5 + 1*2^3 + 1*2^2 + 1*2^0 = -19$$

# Add/sub-enheten



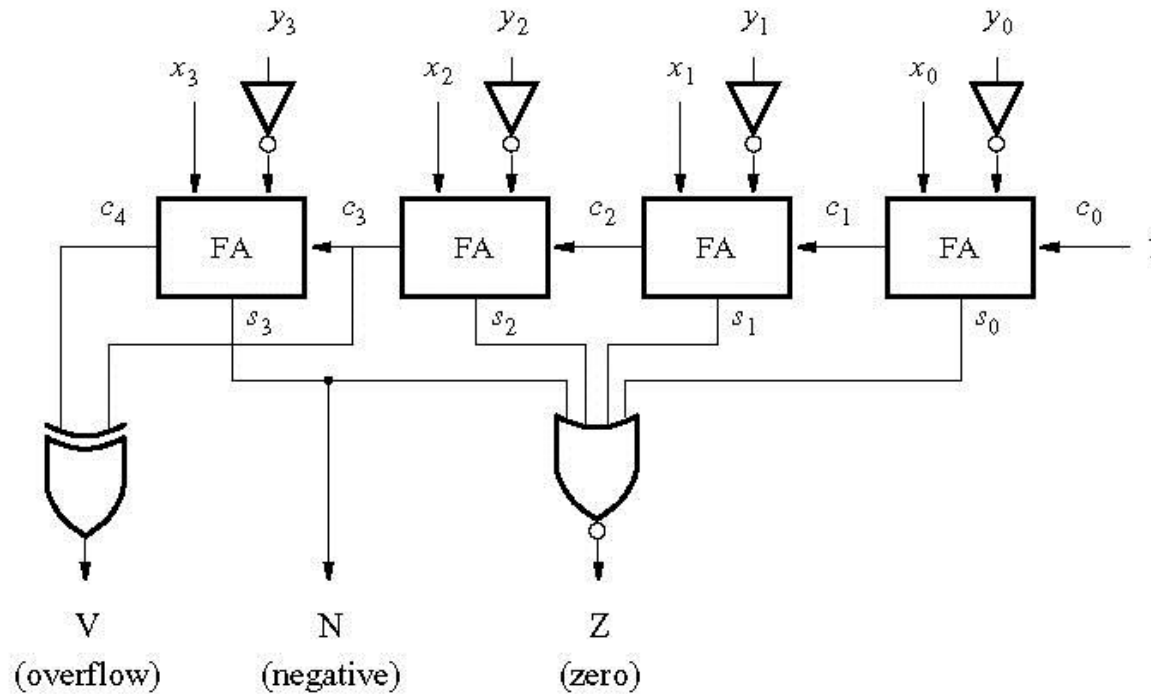
# Arithmetic Logic Unit (ALU)



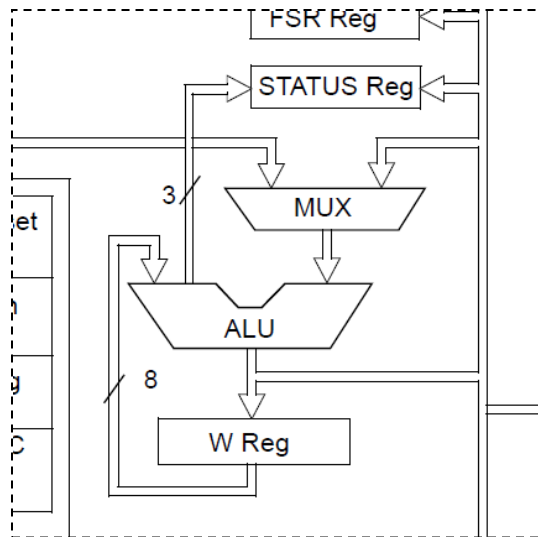
A/L	f1	f0	Funktion
0	0	0	$x+y$
0	0	1	$x+y+C_{in}$
0	1	0	$x-y$
0	1	1	$x-y-\overline{C_{in}}$
1	0	0	$x \text{ or } y$
1	0	1	$x \text{ and } y$
1	1	0	$x \text{ xor } y$
1	1	1	$\text{inv } x$

# Komparator

- Komparatorn implementeras som subtraktionskrets



# Exempel: ALU i mikrokontroller



Microkontroller finns tex i smartcard

## 9.1 Instruction Description

### **ADDWF**      Add W and f

Syntax:      [ *label* ] ADDWF    f,d

Operands:     $0 \leq f \leq 31$   
                   $d \in [0,1]$

Operation:     $(W) + (f) \rightarrow (\text{destination})$

Status Affected: C, DC, Z

Description:    Add the contents of the W register and register 'f'. If 'd' is '0', the result is stored in the W register. If 'd' is '1', the result is stored back in register 'f'.

Exempel på instruktion i mikrokontroller

# Multiplikation av två (positiva) heltal



$$\begin{array}{r} \phantom{+} \phantom{0} \phantom{0} \phantom{1} \phantom{0} \\ \phantom{+} \phantom{0} \phantom{0} \phantom{1} \phantom{1} \\ \hline \phantom{+} \phantom{0} \phantom{0} \phantom{1} \phantom{0} \\ \phantom{+} \phantom{0} \phantom{0} \phantom{1} \phantom{0} \\ \phantom{+} \phantom{0} \phantom{0} \phantom{0} \phantom{0} \\ + \phantom{0} \phantom{0} \phantom{1} \phantom{0} \\ \hline 0 \phantom{0} \phantom{1} \phantom{0} \phantom{1} \phantom{1} \phantom{0} \end{array} \quad \begin{array}{r} 2 \\ 11 \\ \\ \\ \\ \\ 22 \end{array}$$



# Multiplikation med en teckenbit

Teckenförläng!  
(Sign Extension)

$$\begin{array}{r}
 \phantom{0000} 1011 \quad -5 \\
 * \phantom{0000} 0010 \quad 2 \\
 \hline
 \phantom{0000} 0000 \\
 \phantom{0000} 111011 \\
 \phantom{0000} 0000 \\
 + \phantom{0000} 0000 \\
 \hline
 1110110 \quad -10
 \end{array}$$



# Eller så gör vi det enkelt för oss...



- Använd enbart positiva tal i multiplikationen
  - Konvertera till positiva tal
  - Håll reda på resultatets tecken

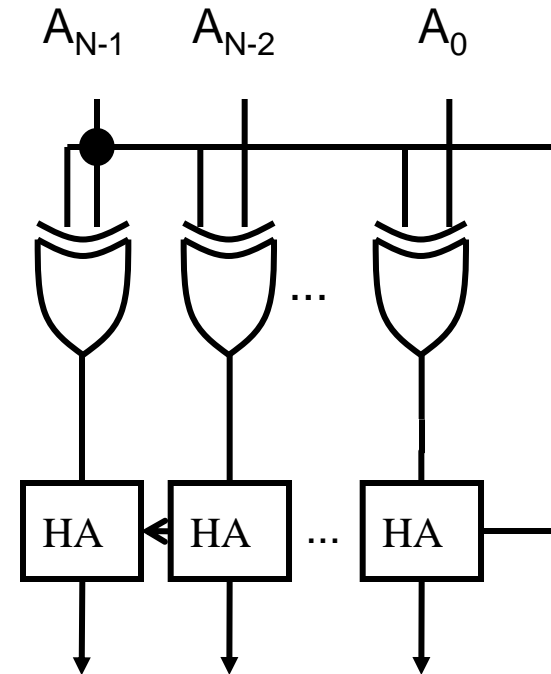
$$(+,+) \Rightarrow +; (+,-) \Rightarrow -; (-,+) \Rightarrow -; (-,-) \Rightarrow +$$

- Två-komplementera till negativt tal om nödvändigt

# En enkel lösning (forts.)

		Sign <sub>A</sub>	
		0	1
Sign <sub>B</sub>	0	0	1
	1	1	0

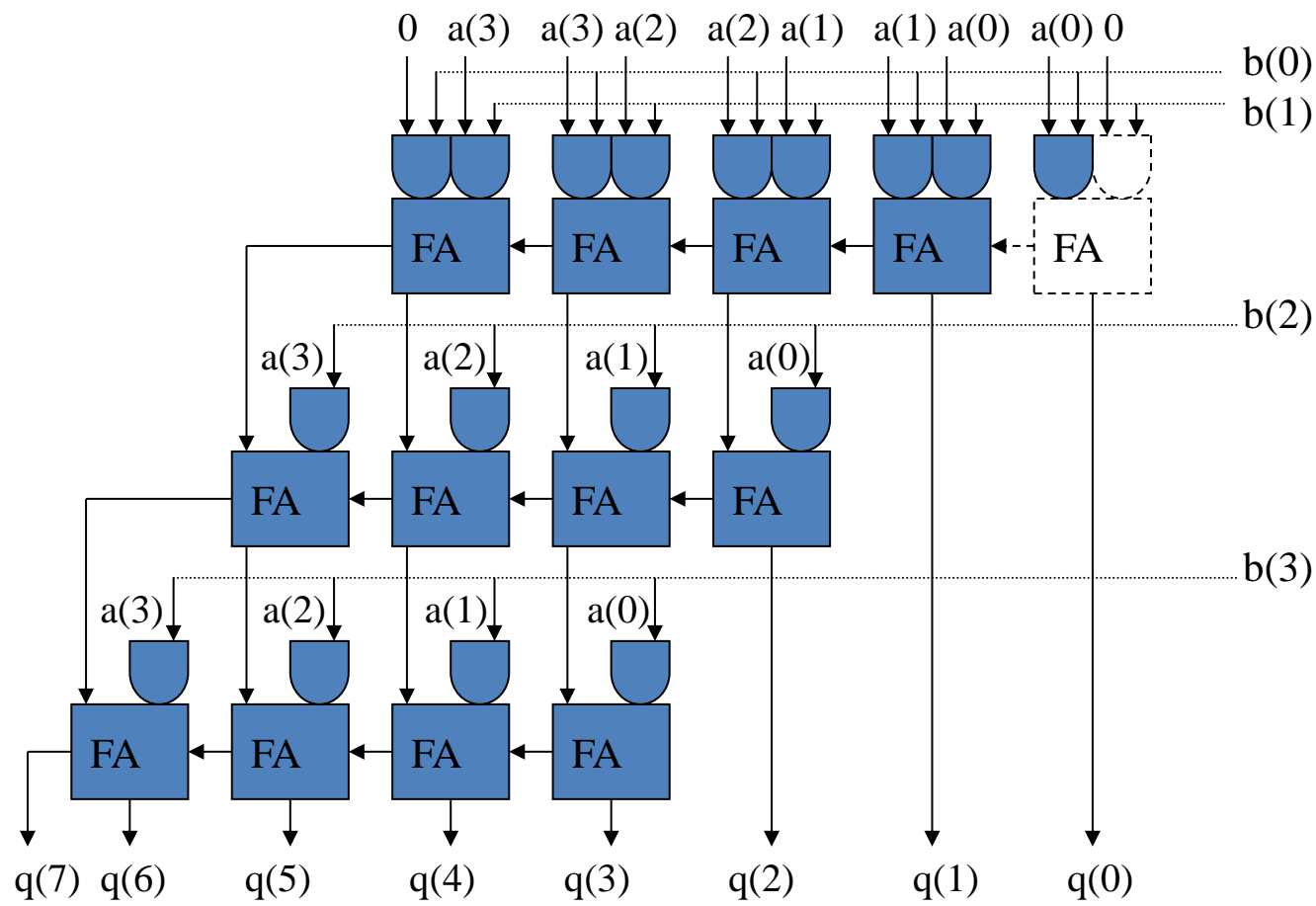
$$F_{\text{invert}} = S_A \text{ xor } S_B$$



2's complement of A



# Multiplikatorn (två positiva tal)



# Snabbfråga repetition



*Varför är en carry look ahead adder snabbare än en ripple adder?*

A: Den har fler antal grindar

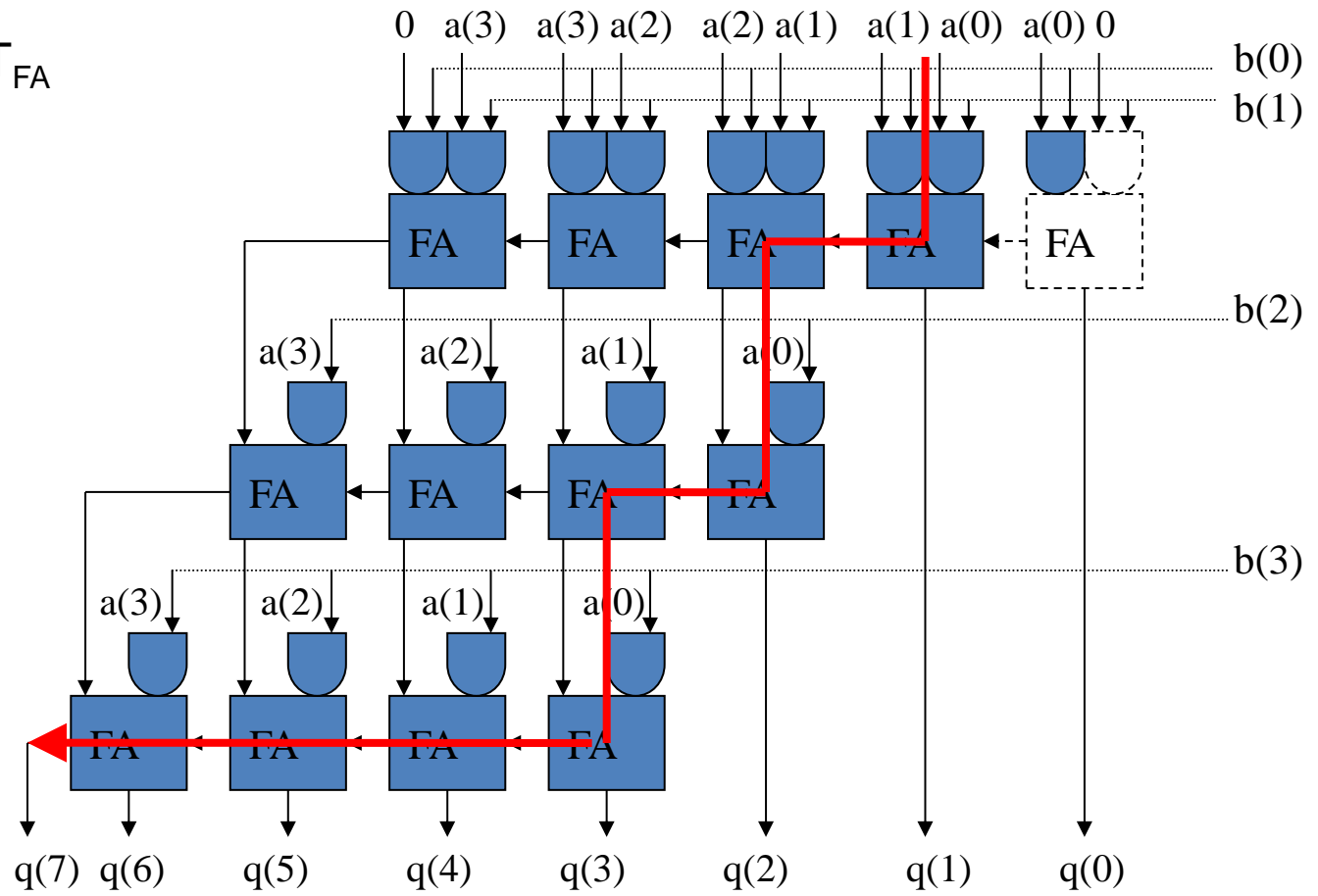
B: Den har färre antal grindar i signalvägen

C: Den kan implementeras med en bättre teknologi



# Multiplikatorn (två positiva tal)

$$T_{MUL} \approx (3N - 4) T_{FA}$$

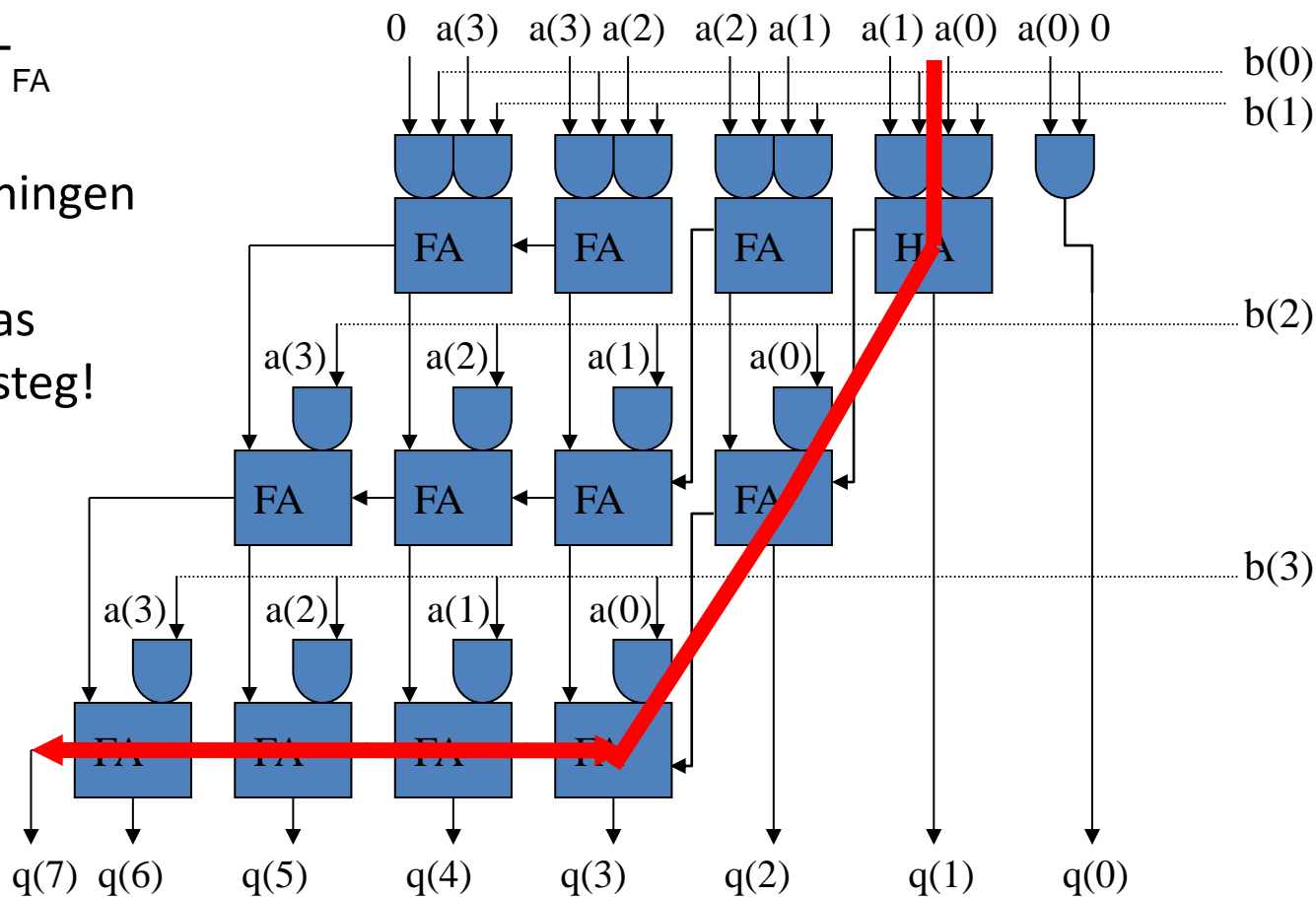




# En snabbare lösning - Carry-Save Adders (BV: sida 311)

$$T_{MUL} \approx (2 \cdot N - 2) \cdot T_{FA}$$

Minskar fördröjningen  
eftersom carry  
utgången kopplas  
direkt till nästa steg!



# Multiplikation med 2



$0101 * 2 = 1010$	$(5 * 2 = 10)$
$1010 * 2 = 10100$	$(-6 * 2 = -12)$
$01010101 * 2 = 010101010$	$(85 * 2 = 190)$
$10010101 * 2 = 100101010$	$(-107 * 2 = -214)$

jfr multiplikation med 10 i basen 10:

$$63 * 10 = 630, -63 * 10 = -630 \text{ etc.}$$

# Multiplikation med $2^n$



$$0101 * 2 = 1010$$

$$(5 * 2 = 10)$$

$$0101 * 2^2 = 10100$$

$$(5 * 4 = 20)$$

$$0101 * 2^3 = 101000$$

$$(5 * 8 = 40)$$

$$0101 * 2^4 = 1010000$$

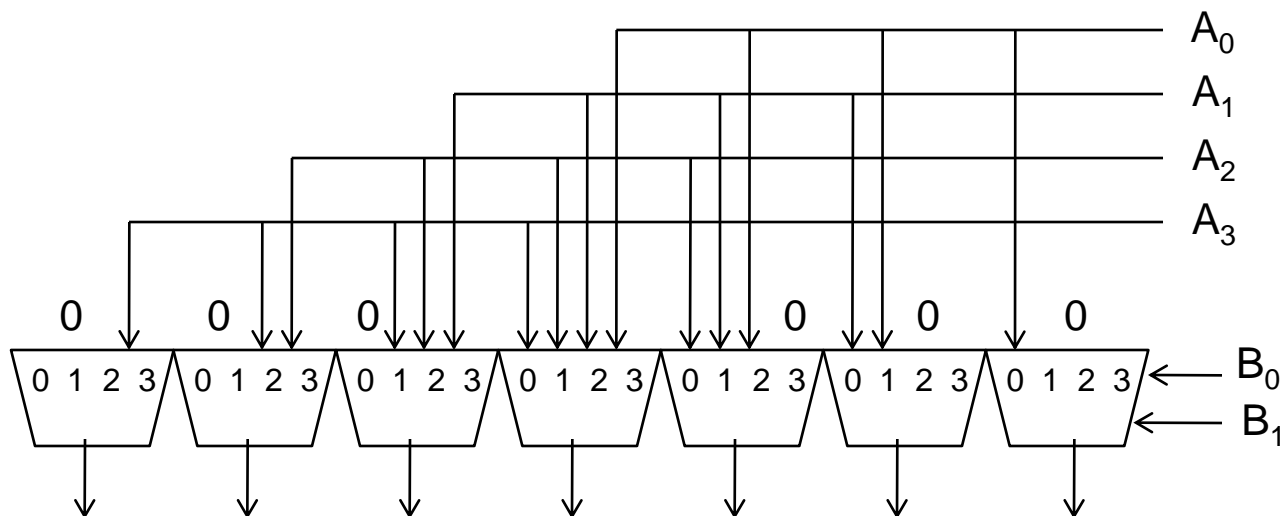
$$(5 * 16 = 80)$$

jmfr multiplikation med 10 i basen 10:

$$6 * 10 = 60, 6 * 100 = 600, 6 * 1000 = 6000 \text{ etc.}$$

- En multiplikation med  $2^n$  kan göras med genom att skifta alla bitar  $n$  steg till vänster och att fylla på med nollor
- $13 * 8$  kan beräknas genom att skifta (01011) tre bitar till höger
  - Resultat: 01011000 motsvarar  $(104)_{10}$
  - Observera att man behöver flera bitar för att representera resultatet!

# Barrel-shiftern



Multiplikation med  $(2^0, 2^1, 2^2, 2^3)$  – dvs 1, 2, 4, 8

# Division mellan två (positiva) heltal (BV sid 693 – Fig 10.21 b)



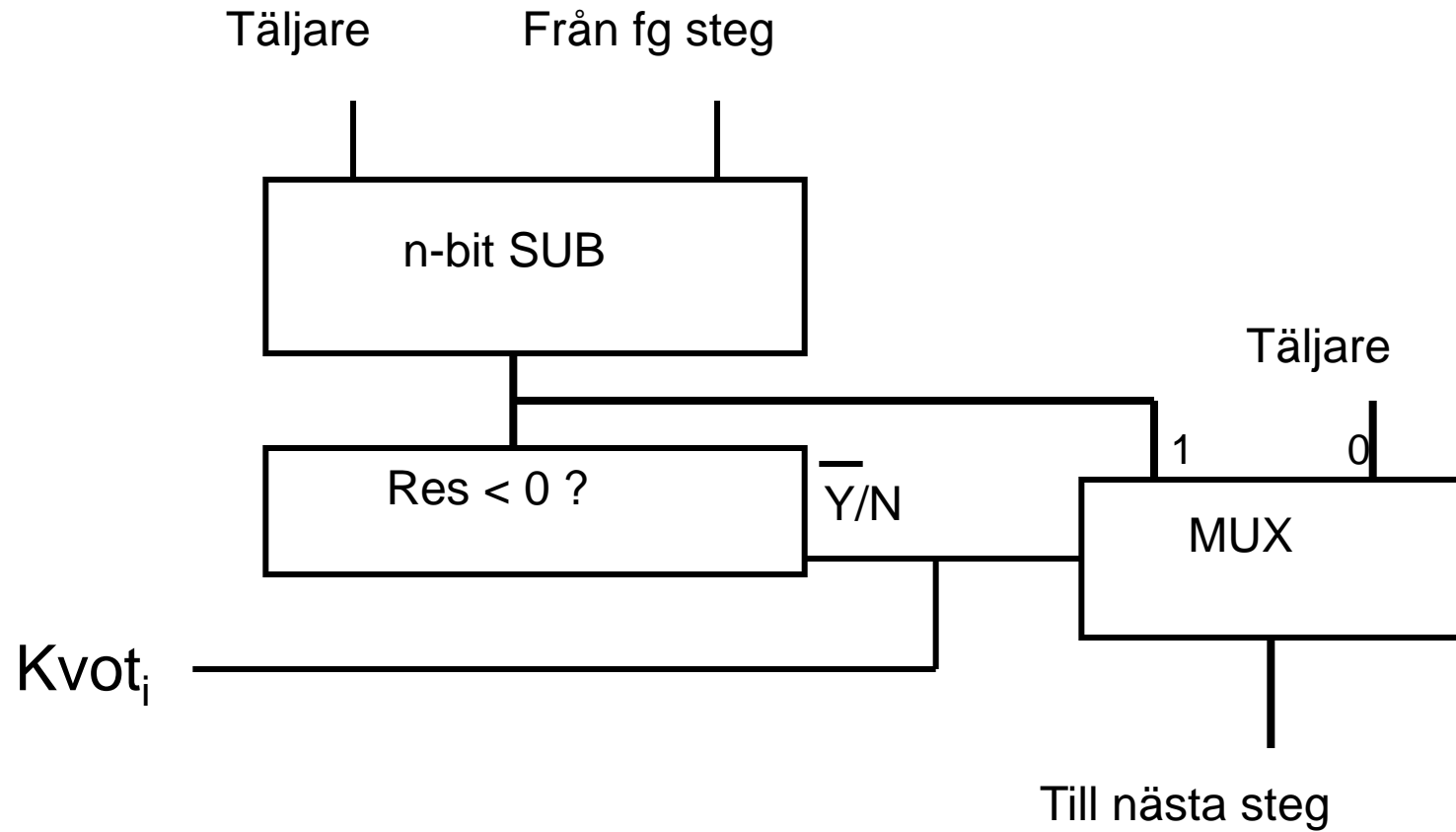
$$\begin{array}{r} \phantom{10} 0101 \\ 10 \overline{) 1011} \\ \underline{- 10} \phantom{1} \\ \phantom{10} 001 \\ \phantom{10} \underline{- 00} \\ \phantom{10} 011 \\ \phantom{10} \phantom{0} \underline{- 10} \\ \phantom{10} \phantom{0} \phantom{0} 1 \end{array}$$

$$11/2 = 5$$

$$\text{Rest} = 1$$

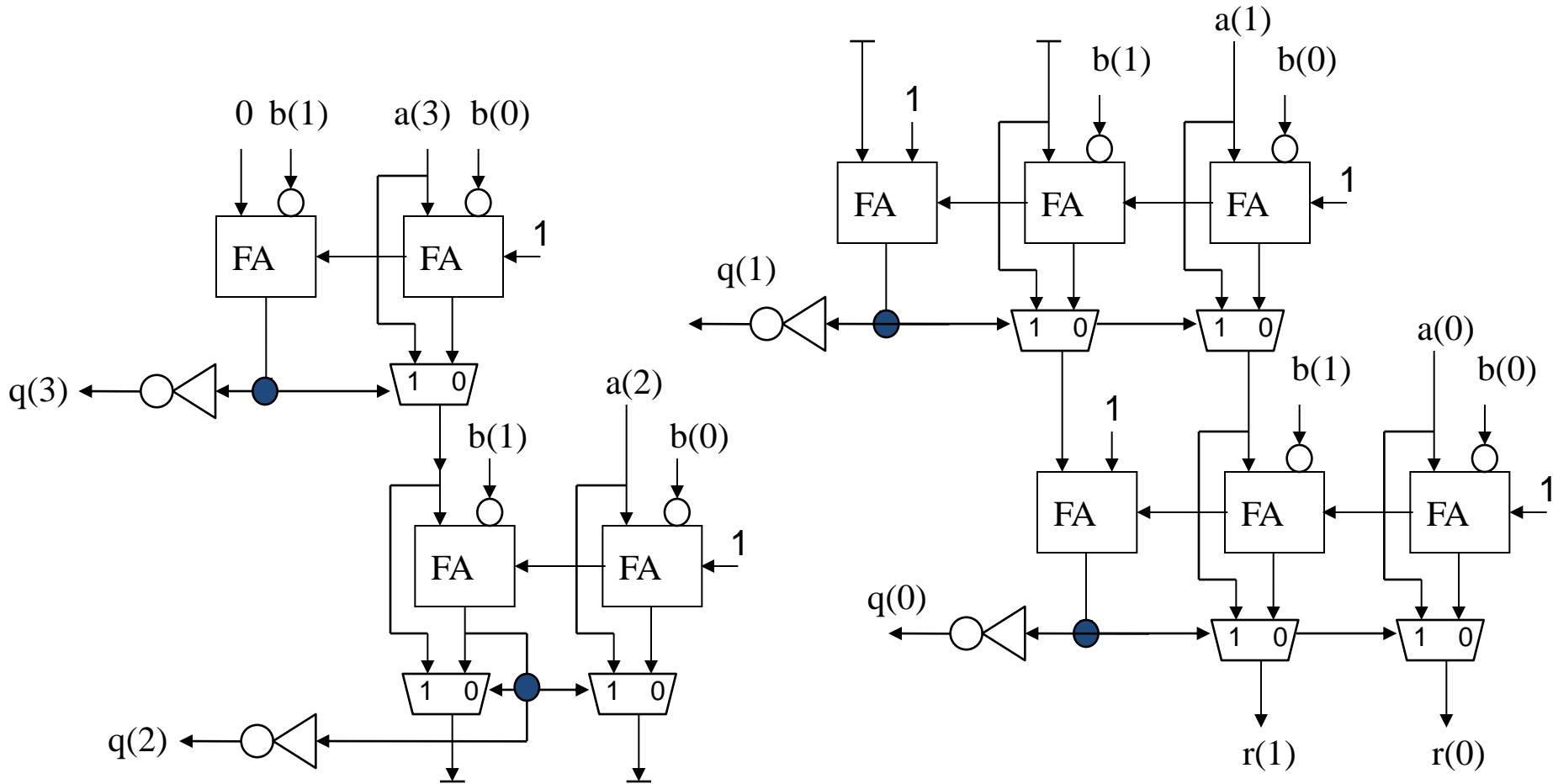


# Divideraren





# Divideraren (forts.)

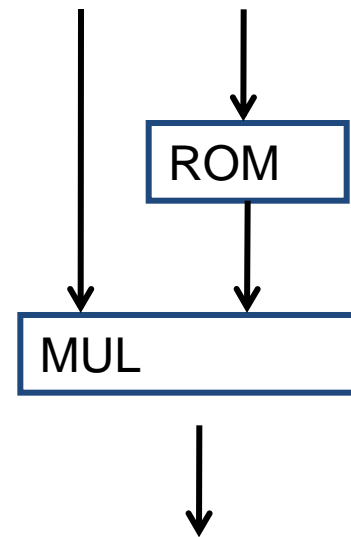


# Using a multiplier and a ROM for division



$$Q=A/B=A*(1/B)$$

- Properties
  - Fast
  - **VERY** big!!



# Division med negativa heltal



- Division med negativa tal är ganska knepigt.
- Ett sätt att utföra divisionen ändå
  - Konvertera till positiva tal
  - Håll reda på resultatets tecken

$$(+,+) \Rightarrow +, (+,-) \Rightarrow -, (-,+) \Rightarrow -, (-,-) \Rightarrow +$$

- Två-komplementera till negativt tal om nödvändigt

# Division med 2



$$\begin{array}{l} 01010/2=00101 \\ \longrightarrow 10100/2=11010 \end{array} \qquad \begin{array}{l} (10/2=5) \\ (-12/2=-6) \end{array}$$

jmfr division med 10 i basen 10:  
 $630/10 = 63$ ,  $-630/10=-63$  etc.

# Logiskt vs Aritmetisk shift höger

- Man skiljer mellan logisk och aritmetiskt shift
  - Logisk shift höger shiftar bara till höger. Bitarna skall ej tolkas som ett tal. Man fyller bara på med 0:or.
  - Aritmetisk shift höger behandlar bitarna som ett tal. Teckenbiten behålls vid shift. Man fyller på till vänster med teckenbiten.

# Division med $2^n$



$1010/2=101$	$(10/2=5)$
$10100/2^2=101$	$(20/4=5)$
$101000/2^3=101$	$(40/8=5)$
$1010000/2^4=101$	$(80/16=5)$

jmfr division med 10 i basen 10:

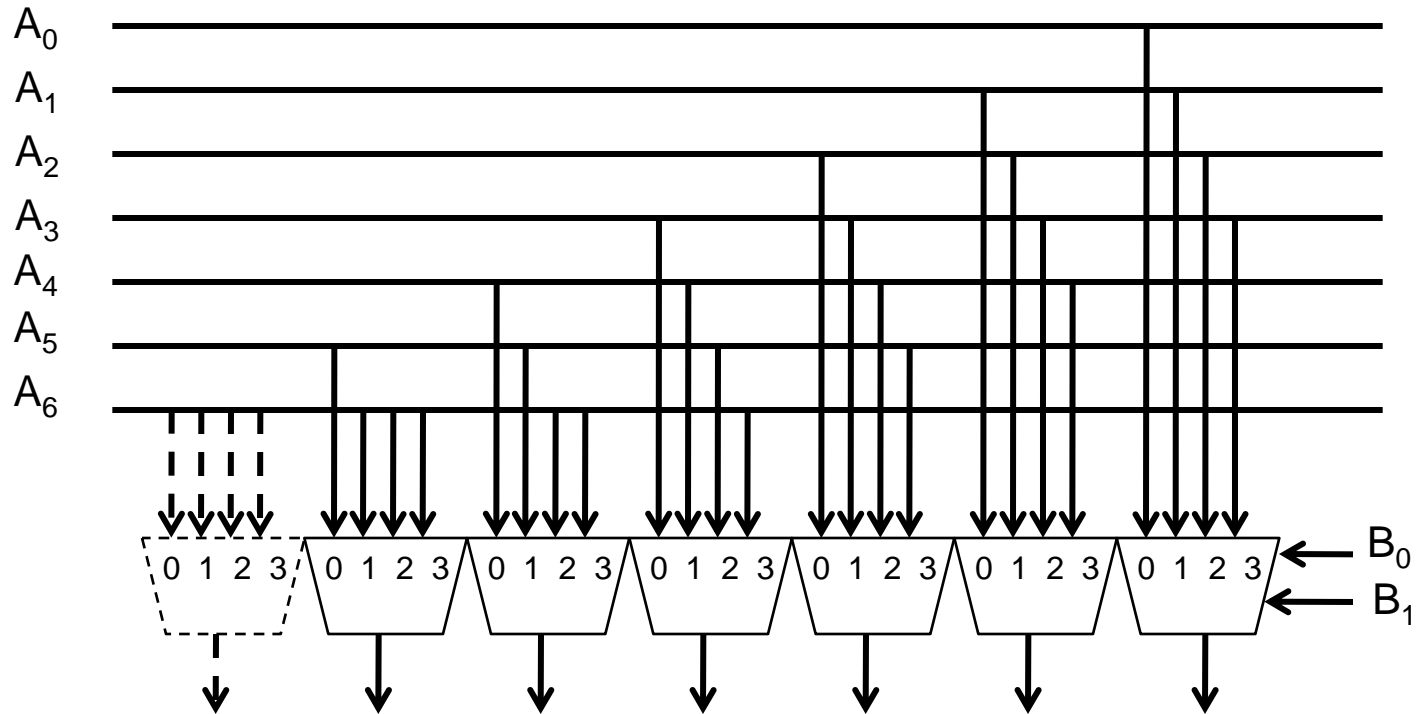
$$60/10 = 6, 600/100=6, 6000/1000=6 \text{ etc.}$$

# Division med $2^n$



- En division med  $2^n$  kan göras med genom att skifta alla bitar  $n$  steg till höger och att fylla på med nollor
- Observera att resultatet behöver inte vara korrekt, eftersom man egentligen behöver bitar "efter kommatecknet"
- $17/4$  motsvarar att skifta  $010001$  2-bitar till höger
  - Resultat:  $000100 = (4)_{10}$
  - Eftersom  $(0.25)_{10}$  inte kan representeras är resultatet inte korrekt!

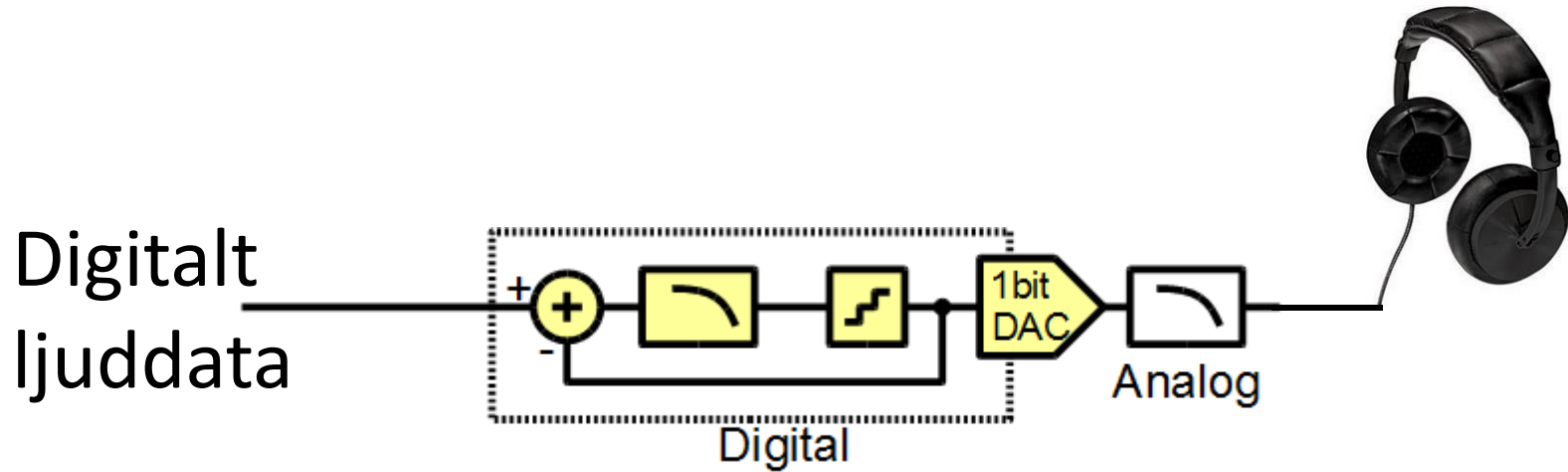
# Barrel-shiftern (2)



Division med  $(2^0, 2^1, 2^2, 2^3)$  – dvs 1, 2, 4, 8



# Exempel signalbehandling



Digital till Analog  
omvandlare (DAC)

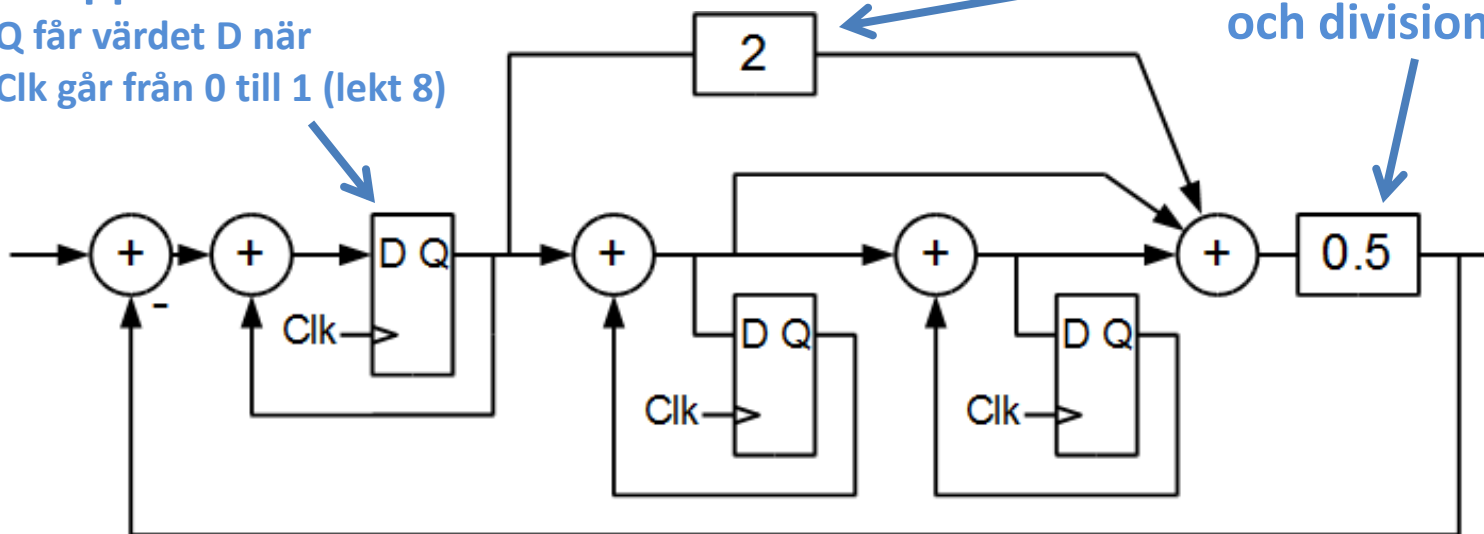
# Exempel signalbehandling

Matematisk funktion:

$$H_{qn}(z) = \frac{(1 - z^{-1})^3}{1 - z^{-1} + 0.5z^{-2}}$$

D-vippa "låser" värdet  
Q får värdet D när  
Clk går från 0 till 1 (lekt 8)

Endast multiplikation  
och division med 2

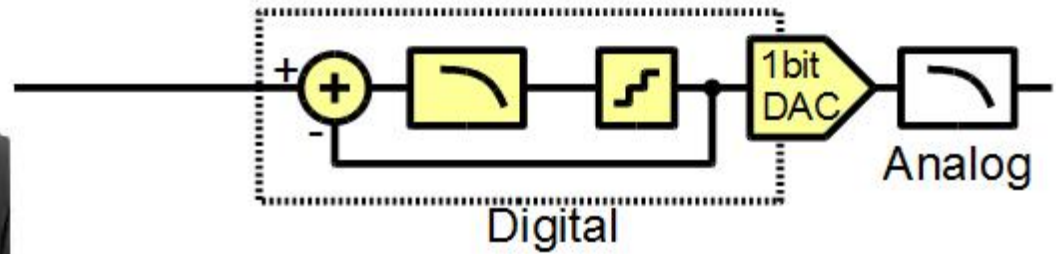


Implementering av funktion i digital hårdvara:

# Exempel signalbehandling

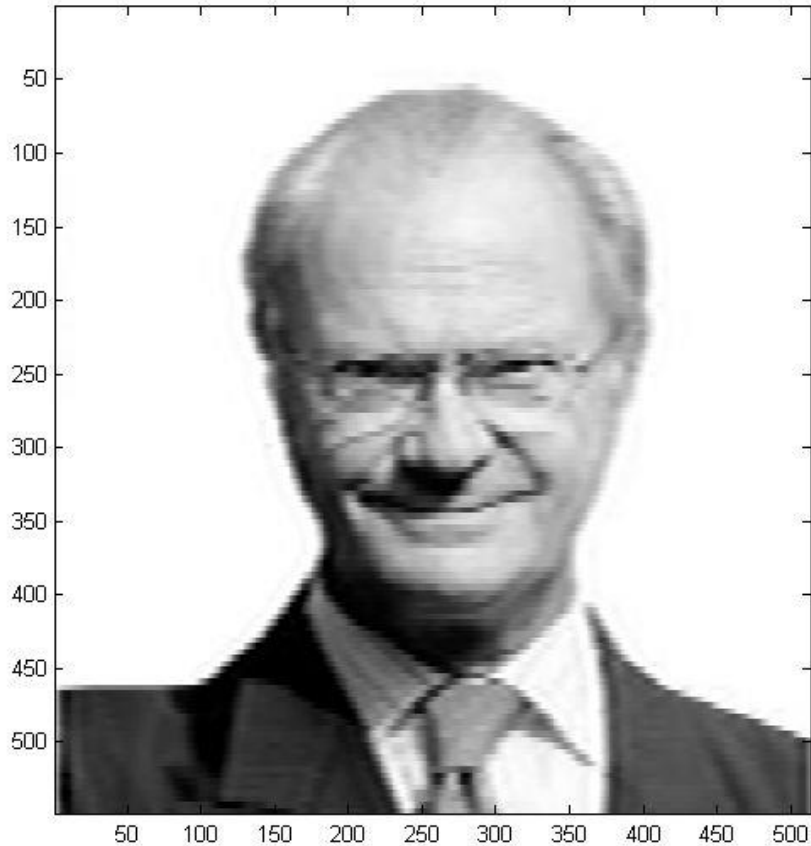


Carl XVI Gustaf

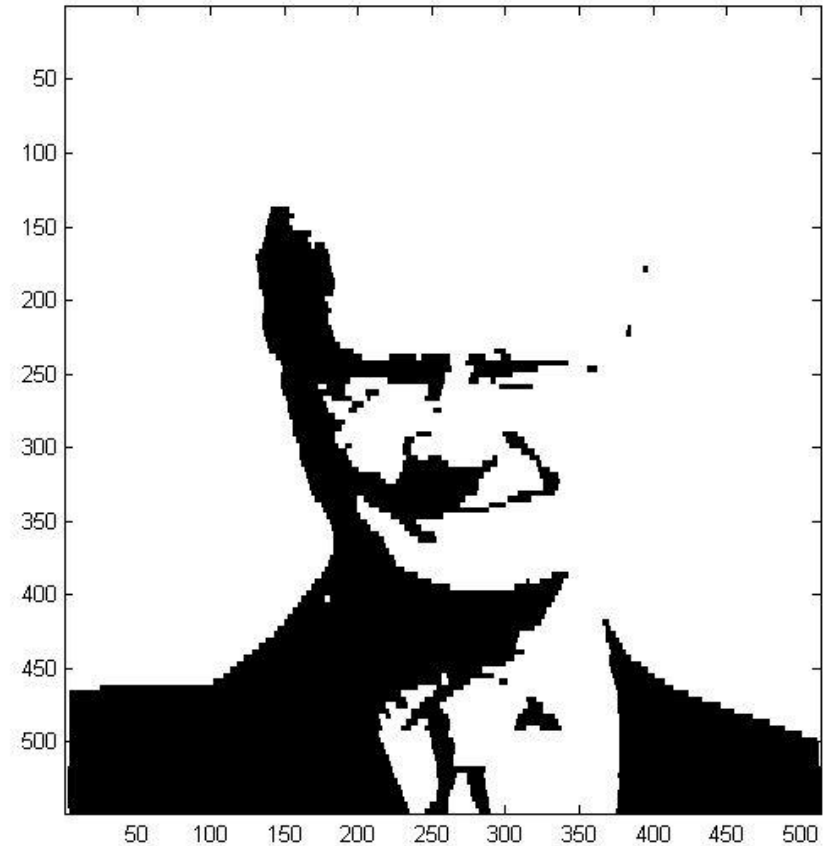


# Exempel signalbehandling

Original



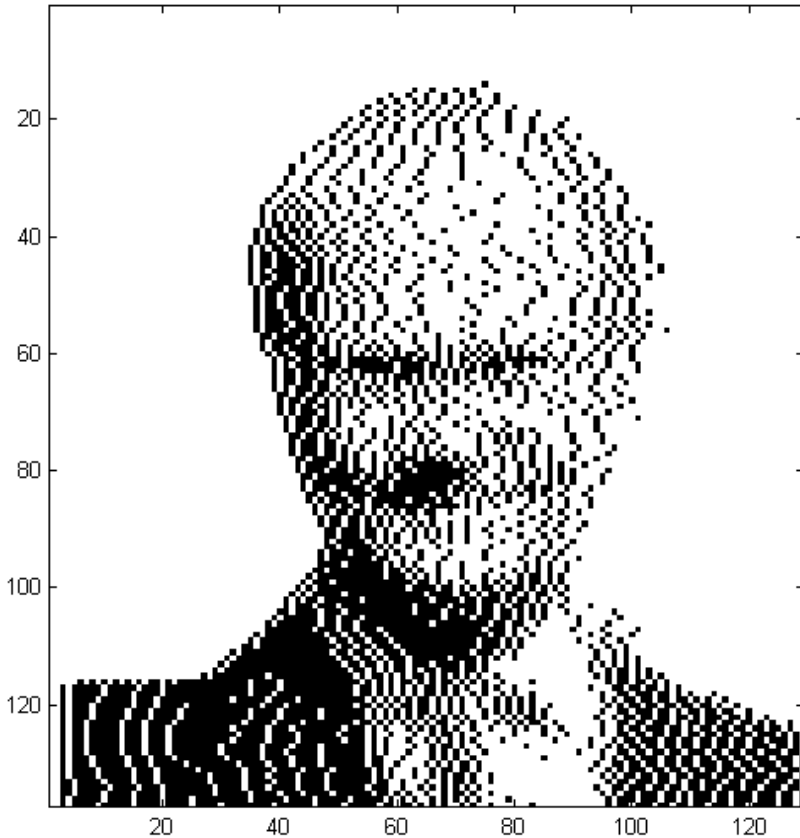
2 level quantization



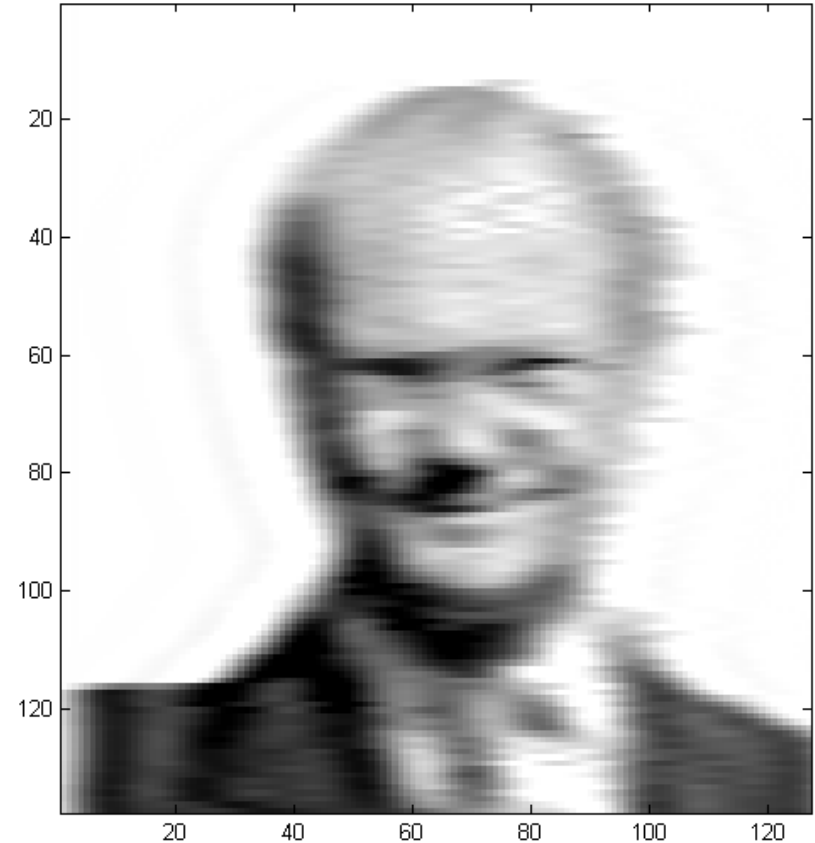
En digital utgång har bara två diskreta värden

# Exempel signalbehandling

1st Order Delta Sigma



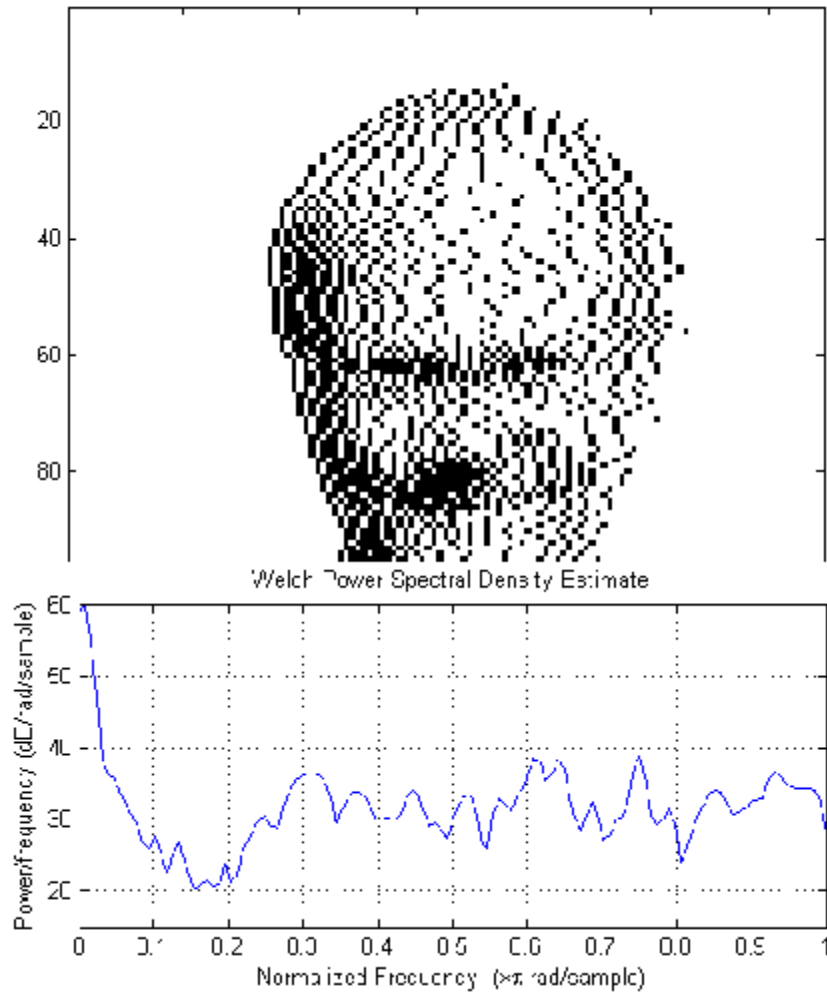
Low pass filtered



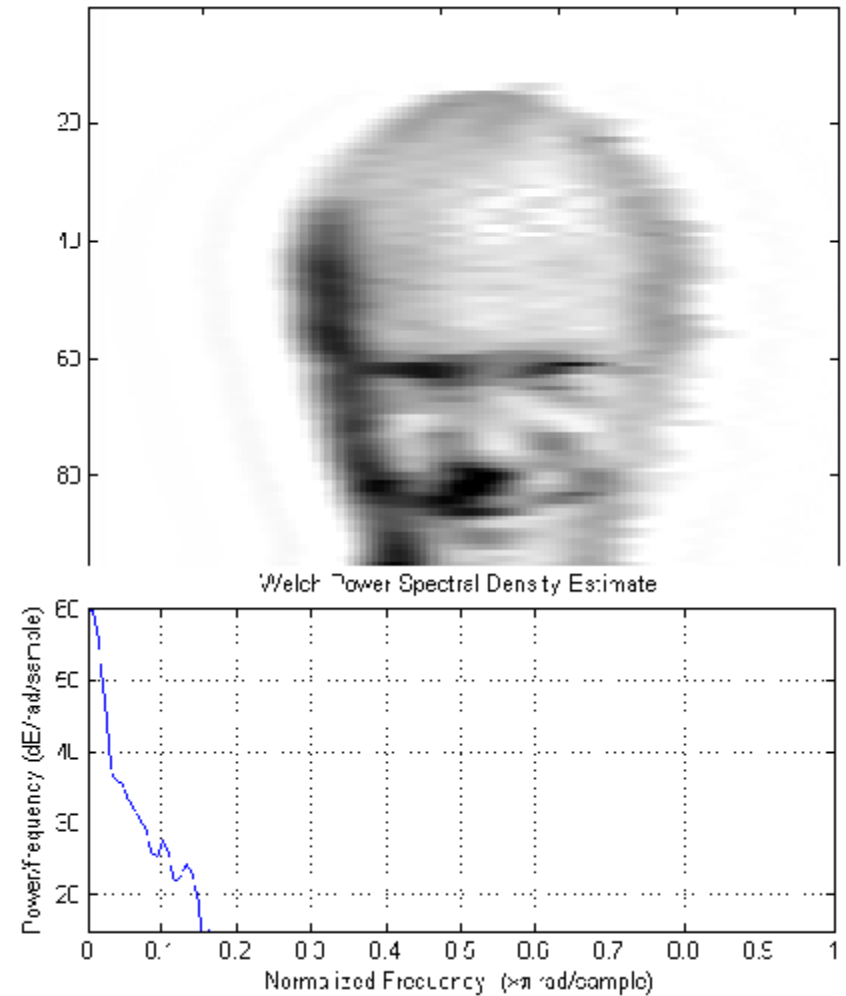
Genom digital signalbehandling kan signalen brusformas  
Efter filtrering återskapas den analoga signalen

# Exempel signalbehandling

1st Order Delta Sigma

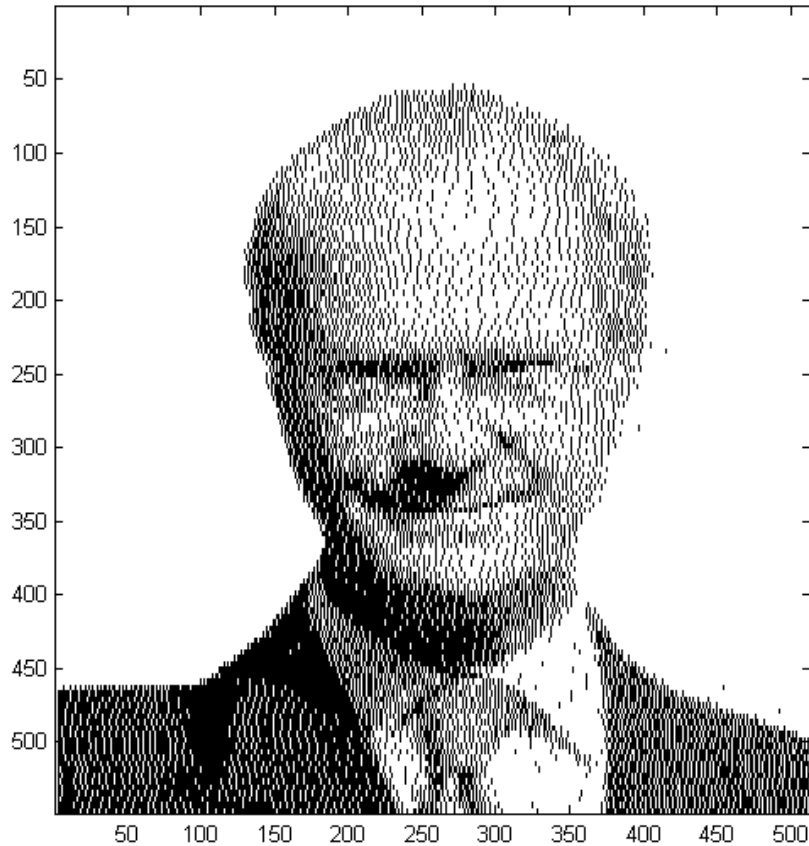


Low pass filter ec

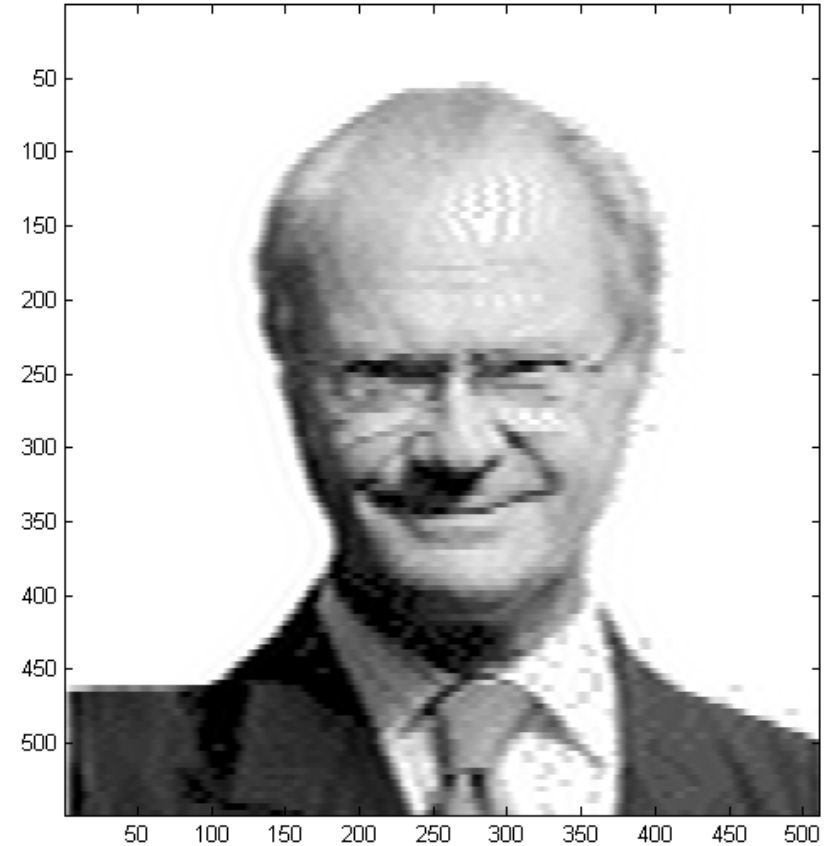


# Exempel signalbehandling

1st Order Delta Sigma

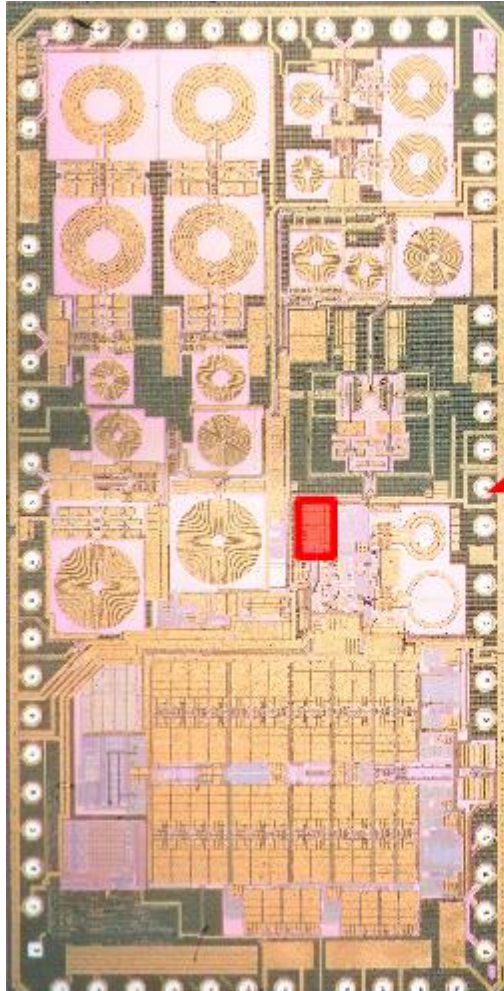


Low pass filtered



Högre klockhastighet möjliggör högre bandbredd hos signalen

# Exempel signalbehandling



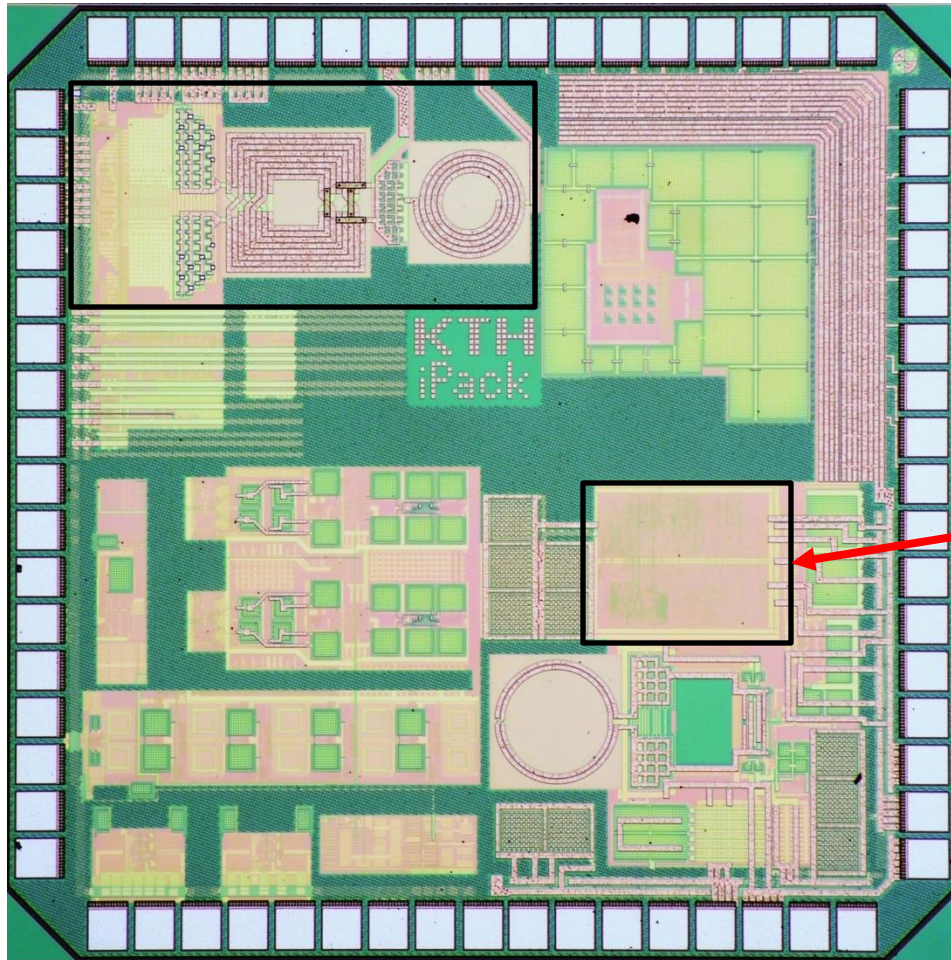
Exempel:

Brusformning i TV tuner  
 $\Delta\Sigma$  modulated Synthesizer  
75 MHz – 1.8 GHz

Tack vare att vi endast multiplicerar/  
dividerar med 2 blir ytan och  
strömförbrukningen mycket liten



# Exempel signalbehandling



Digital

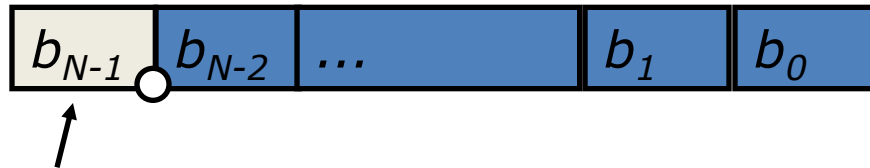
I denna krets är digitaldelen inte lika optimerad (innehåller en division). Den tar större yta och har högre strömförbrukning.

# Fix-tal (Fixed-point numbers)



## Två-komplementsrepresentation

$$B = b_{N-1} . b_{N-2} \dots b_1 b_0 \quad \text{where } b_i \in \{0, 1\}$$



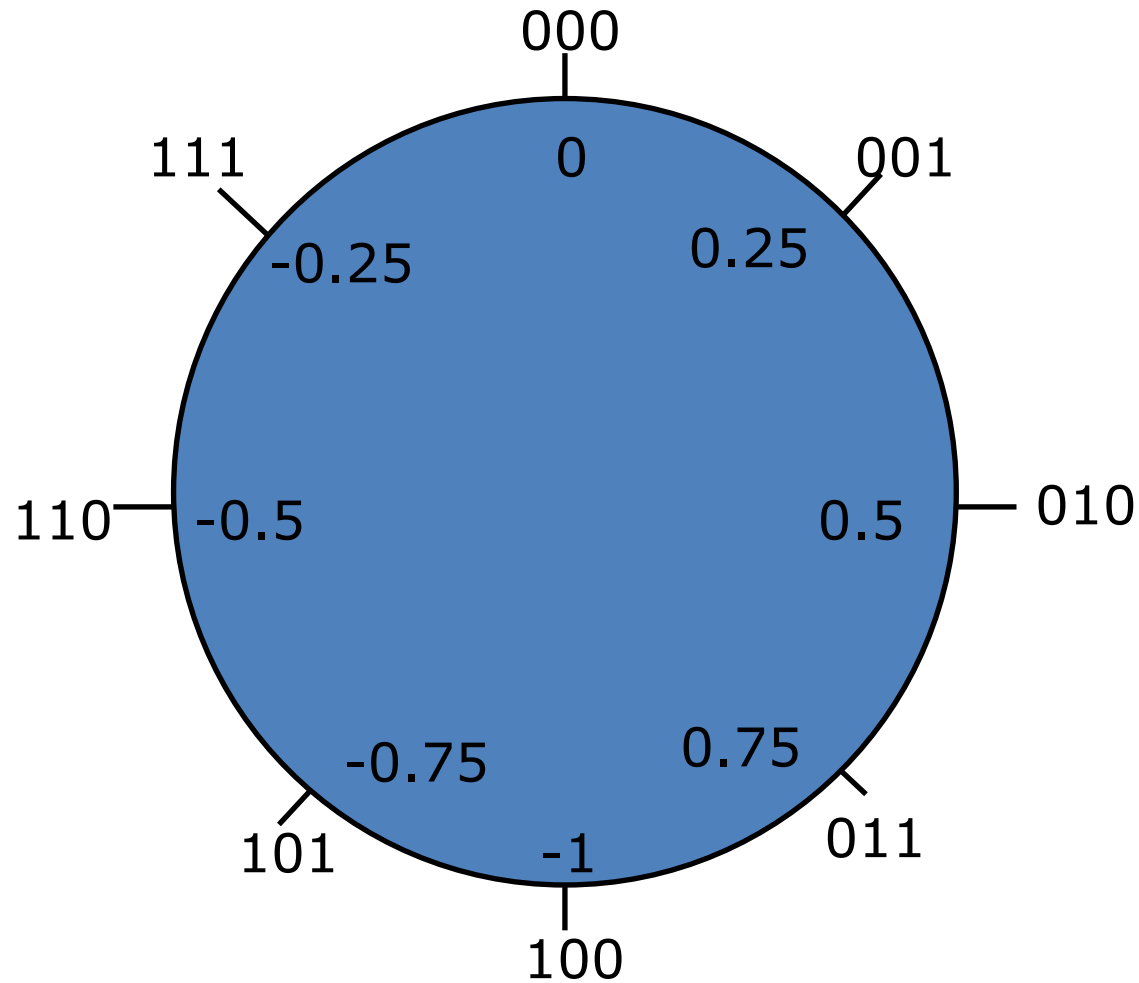
Sign Bit

## Decimalt värde

$$FiP(B) = -b_{N-1} 2^0 + b_{N-2} 2^{-1} + \dots + b_1 2^{-(N-2)} + b_0 2^{-(N-1)}$$

Detta format kallas också för  $Q_{N-1}$ -format eller *fractional representation*

# Fixed-Point Q2-Representation

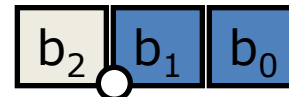


# Multiplikation i Q-formatet

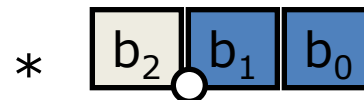
- Multiplikation orsakar inte overflow i Q-formatet, men kan resultera i förlust av precision

$$\begin{aligned}
 & (1.10)_2 * (0.11)_2 \\
 & = (1.1010)_2 \quad (\text{Q4-format}) \\
 & = (1.10)_2 \quad (\text{Q2-format})
 \end{aligned}$$

Generellt:

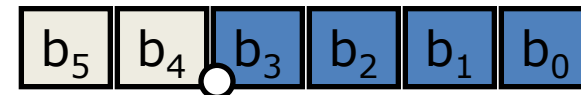


Q2



Q2

Utökade teckenbitar



Q4

# Multiplikation i Q-formatet



Sign Extension

Stryks

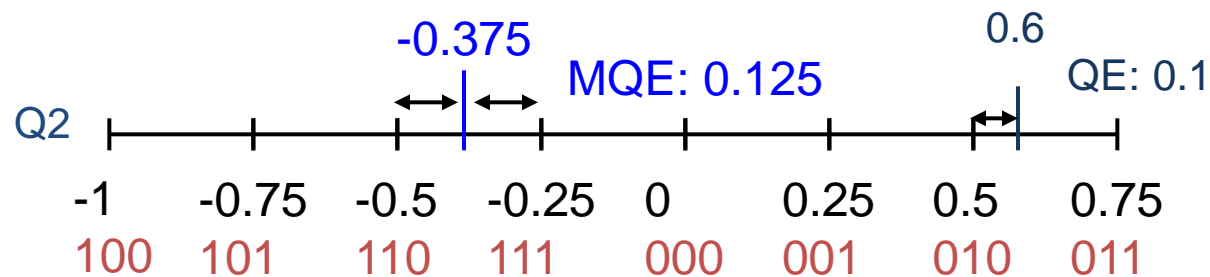
$$\begin{array}{r} 0.11 * 1.10 \\ \hline 1.110 \\ 1.1110 \\ \hline 11.1010 \end{array}$$

Negation  
2-komplement

$$\begin{array}{r} 1.10 * 0.11 \\ \hline 1.01 \\ 0.011 \\ \hline 1.101 \end{array}$$

Glöm inte att expandera teckenbiten!

# Kvantiseringsfel



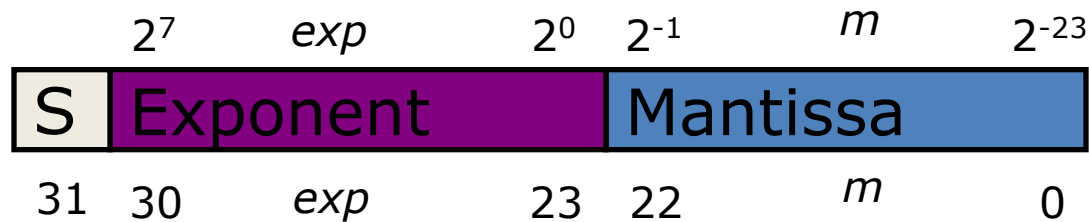
- Maximalt kvantiseringsfel:
  - Eftersom inte alla tal kan representeras uppstår kvantiseringsfel
  - Kvantiseringsfelet blir mindre om vi använder fler bitar

- Ett flyt-tal representeras med en *tecken-bit*, *exponent-bitar* och en mantissa (ä.k. *fraktions-bitar*)



- Värdet beräknas som
$$FIP(B) = (-1)^s * (1.m) * 2^{\pm exp}$$
- Ofta är *exp* biaserad (har en offset), vilket då ger
$$FIP(B) = (-1)^s * (1.m) * 2^{exp-(bias)}$$

- Flyttals-standarden IEEE-754 definierar ett 32-bit flyttal som

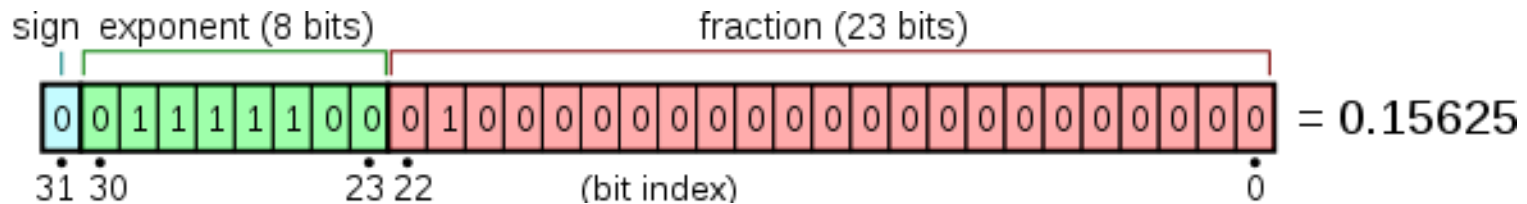


- Värdet beräknas för en 8-bitars exponent enligt nedan  

$$FIP(B) = (-1)^s * (1.m) * 2^{exp-(127)}$$
- Speciala bit pattern har reserverats för att representera negativ och positiv nolla



# Exempel



Sign = 0: Positivt

Exp:  $(64 + 32 + 16 + 8 + 4) - 127 = 124 - 127 = -3$

Mantissa:

Binary 1.01 =  $1 + 1 / 4 = 1.25$

*Första biten är alltid ett eftersom mantissan är normaliserad*

$1.25 * 2^{-3} = 1.25 * 0.125 = 0.15625$

# Floating-Point Numbers (IEEE)



- Exponent Values 1 to 254: normalized non-zero floating-point numbers; biased exponent (-126...+127)
- Exponent of zero and fraction of zero: positive or negative zero
- Exponent of ones (255) and fraction of zero: positive or negative infinity
- Exponent of zero and fraction of non-zero: Denormalized number (true exponent is -126), represent numbers from 0 to  $2^{-126}$

$$FIP(B) = (-1)^s * (\mathbf{0}.m) * 2^{-126}$$

- Exponent of ones (255) with a non-zero fraction: *NotANumber* (Exception Condition)
- There is also a standard for a 64-bit number

# Addition med flyttal



- Givet två flyttal:

$$a = a_{frac} \cdot 2^{a_{exp}}$$

$$b = b_{frac} \cdot 2^{b_{exp}}$$

- Summan av dessa tal är:

$$c = a + b$$

$$= \begin{cases} (a_{frac} + (b_{frac} \cdot 2^{-(a_{exp} - b_{exp})})) * 2^{a_{exp}} & , \text{if } a_{exp} \geq b_{exp} \\ (b_{frac} + (a_{frac} \cdot 2^{-(b_{exp} - a_{exp})})) * 2^{b_{exp}} & , \text{if } b_{exp} \geq a_{exp} \end{cases}$$

- Givet två flyttal:

$$a = a_{frac} \cdot 2^{a_{exp}}$$

$$b = b_{frac} \cdot 2^{b_{exp}}$$

- Differensen mellan dessa tal är:

$$c = a - b$$

$$= \begin{cases} (a_{frac} - (b_{frac} \cdot 2^{-(a_{exp} - b_{exp})})) * 2^{a_{exp}} & , \text{if } a_{exp} \geq b_{exp} \\ (b_{frac} - (a_{frac} \cdot 2^{-(b_{exp} - a_{exp})})) * 2^{b_{exp}} & , \text{if } b_{exp} \geq a_{exp} \end{cases}$$

# Multiplikation med flyttal

- Givet två flyttal:

$$a = a_{frac} \cdot 2^{a_{exp}}$$

$$b = b_{frac} \cdot 2^{b_{exp}}$$

- Produkten av dessa tal är:

$$\begin{aligned} c &= a * b \\ &= \left( a_{frac} * b_{frac} \cdot 2^{a_{exp} + b_{exp}} \right) \end{aligned}$$

- Givet två flyttal:

$$a = a_{frac} \cdot 2^{a_{exp}}$$

$$b = b_{frac} \cdot 2^{b_{exp}}$$

- Kvoten mellan dessa tal är:

$$\begin{aligned} c &= a / b \\ &= \left( a_{frac} / b_{frac} \cdot 2^{a_{exp} - b_{exp}} \right) \end{aligned}$$

- När en flyttals-operation är klar måste den normaliseras
  - Mantissans skiftas tills dess första bit är 1
  - För varje skift-steg så räknas exponenten upp eller ned med ett.
  - Mantissans bitar till höger om den första ettan sparas
- Om exponenten är noll är mantissans första bit 0

$$FIP(B) = (-1)^s * (1.m) * 2^{exp-(127)}$$

$$FIP(B) = (-1)^s * (0.m) * 2^{-(126)}$$

# Fixed-Point vs. Floating-Point



- Fixed-Point operationer fungerar på samma sätt som heltals-operations och är snabbare
- Fixed-point värden behöver skalas, vilket ofta leder till förlust av precision
- Kostnaden för att bygga hårdvara är signifikant större för flyttals-processorer/räknare



# Snabbfråga

*Hur representeras 0.75 i Q4 fixed point format?*

A: 01001011

B: 00001100

C: 1100



# Information om laborationer



- Kontrollera labbtid och kunskapsgrupp i Daisy
- Gör förberedelseuppgifter och fyll i web-formulär

- **Multiplikation och division av heltal**
  - Konvertera negativa tal till sitt positiva ditto.
  - Utför multiplikationen eller divisionen
  - Håll reda på vilket tecken resultatet skall ha
  - Konvertera positivt resultat till sitt negativa ditto om resultatet skall vara negativt
- **Multiplikation med potenser av 2 (mul med  $2^k$ )**
  - Implementeras som ett skift till vänster med  $k$  steg
- **Division med potenser av 2 (div med  $2^k$ )**
  - Implementeras som ett (aritmetiskt) skift till höger med  $k$  steg. Teckenbiten kopieras till vänster.