

DSP Design

Numbering Systems Basic Building Blocks Scaling and Round-off Noise

Viktor Öwall

viktor.owall@eit.lth.se



Viktor Öwall, Dept. of Electrical and Information Technology, Lund University, Sweden - www.eit.lth.se

Number Representation

Viktor Öwall, Dept. of Electrical and Information Technology, Lund University, Sweden - www.eit.lth.se

Floating vs. Fixed point

In floating point a value is represented by

- mantissa determining the resolution/precision
- exponent determining the dynamic range

In fixed point we only have a single value

Floating point gives higher dynamic range but the cost is high in

- energy
- area
- calculation time

For energy efficient implementations fixed point is preferred

Viktor Öwall, Dept. of Electrical and Information Technology, Lund University, Sweden - www.eit.lth.se

Binary numbers, unsigned integers

MSB =
Most Significant Bit

LSB =
Least Significant Bit

	2^2	2^1	2^0	
0	0	0	0	(0)
0	0	0	1	(1)
0	0	1	0	(2)
0	0	1	1	(3)
1	0	0	0	(4)
1	0	0	1	(5)
1	1	0	0	(6)
1	1	0	1	(7)

N bits
↓
2^N ord

Viktor Öwall, Dept. of Electrical and Information Technology, Lund University, Sweden - www.eit.lth.se

Dynamic range and Resolution

Nr. of bits	Nr. of levels	Resolution $V_{fs}=0.5V$	Dynamic Range $V_{LSB}=0.03125$
4	16	0.03125V	0.5V
8	256	2mV	8V
12	4096	0.12mV	128V
16	65 536	7.6μV	2042V

How do we use the bits?
Depends on the application!

Unsigned Number Representation

Fixed radix (base) systems

The digits $a \in \{0, 1, 2, \dots, r-1\}$ in a radix r system:

$$\begin{aligned} & \sum_{i=k-1}^{-l} r^i \times a_i = \\ & = r^{k-1}a_{k-1} + r^{k-2}a_{k-2} \cdots r^1a_1 + r^0a_0 + r^{-1}a_{-1} \cdots r^{-l}a_{-l} \\ & \text{described in a fixed point positional number system:} \\ & a_i a_{i-1} \cdots a_1 a_0 \cdot \underbrace{a_{-1} \cdots a_{-l}}_{\text{Fractional part}} \end{aligned}$$

Example: Unsigned Number

$$\begin{aligned} & \sum_{i=k-1}^{-l} 10^i a_i = \{a \in \{0, 1, 2, \dots, 9\} \text{ in radix 10}\} \\ & = 10^{k-1}a_{k-1} + 10^{k-2}a_{k-2} \cdots 10^1a_1 + 10^0a_0 + 10^{-1}a_{-1} \cdots 10^{-l}a_{-l} \end{aligned}$$

$$\begin{aligned} & \sum_{i=k-1}^{-l} 2^i a_i = \{a \in \{0, 1\} \text{ in radix 2}\} \\ & = 2^{k-1}a_{k-1} + 2^{k-2}a_{k-2} \cdots 2^1a_1 + 2^0a_0 + 2^{-1}a_{-1} \cdots 2^{-l}a_{-l} \end{aligned}$$

Example: Unsigned Number

$$\begin{aligned} & \sum_{i=k-1}^{-l} 2^i a_i = \{a \in \{0, 1\} \text{ in radix 2}\} \\ & = 2^{k-1}a_{k-1} + 2^{k-2}a_{k-2} \cdots 2a_1 + a_0 + 2^{-1}a_{-1} \cdots 2^{-l}a_{-l} \end{aligned}$$

$$\begin{aligned} & 1010.0110 \Rightarrow \\ & 1 \cdot 2^3 + 0 \cdot 2^2 + 1 \cdot 2^1 + 0 \cdot 2^0 + 0 \cdot 2^{-1} + 1 \cdot 2^{-2} + 1 \cdot 2^{-3} + 0 \cdot 2^{-4} = \\ & 8 + 2 + \frac{1}{4} + \frac{1}{8} \end{aligned}$$

Signed Digit Number Representation

The digits $a \in \{-\alpha, \dots, 0, \dots, r-1-\alpha\}$ in a radix r system:

$$\sum_{i=k}^{-l} r^i \times a_k$$

Example Radix 10: $a \in \{-4, -3, \dots, 0, \dots, 4, 5\}$

$$(3 \ -1 \ 5)_{10} = 10^2 \times 3 - 10^1 \times 1 + 10^0 \times 5 = 300 - 10 + 5 = 295$$

$$(3 \ .-1 \ 5)_{10} = 10^0 \times 3 - 10^{-1} \times 1 + 10^{-2} \times 5 = 3 - 0.1 + 0.05 = 2.95$$

Viktor Öwall, Dept. of Electrical and Information Technology, Lund University, Sweden - www.eit.lth.se

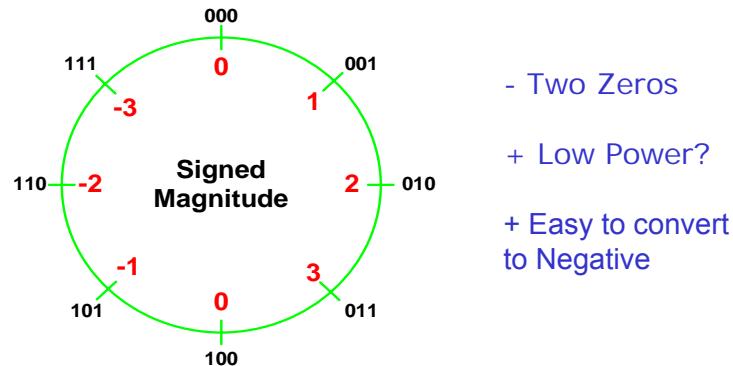
Signed Number Representation

Sign Magnitude
One's Complement
Two's Complement

Viktor Öwall, Dept. of Electrical and Information Technology, Lund University, Sweden - www.eit.lth.se

Signed Magnitude

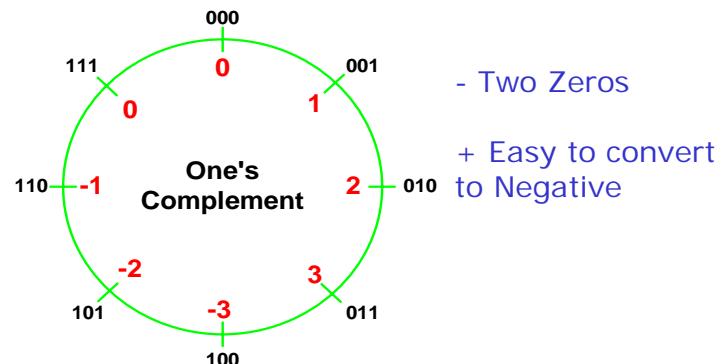
Unsigned numbers with a sign-bit



Viktor Öwall, Dept. of Electrical and Information Technology, Lund University, Sweden - www.eit.lth.se

One's Complement

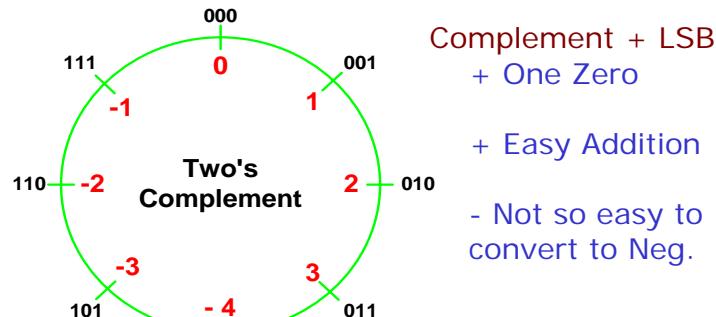
Signed numbers by inverting (Complement)



Viktor Öwall, Dept. of Electrical and Information Technology, Lund University, Sweden - www.eit.lth.se

Two's Complement

Most widely used fixed point numbering system



Viktor Öwall, Dept. of Electrical and Information Technology, Lund University, Sweden - www.eit.lth.se

Two's Complement

The digits $a \in \{0,1\}$ in a radix 2 system:

$$\begin{aligned} & -2^{k-1} \times a_{k-1} + \sum_{i=k-2}^{-l} 2^i \times a_i = \\ & = -2^{k-1} a_{k-1} + 2^{k-2} a_{k-2} \dots 2^1 a_1 + 2^0 a_0 + 2^{-1} a_{-1} \dots 2^{-l} a_{-l} \\ & \text{described in a fixed point positional number system:} \\ & a_{k-1} a_{k-2} \dots a_1 a_0 \cdot \underbrace{a_{-1} \dots a_{-l}}_{\text{Fractional part}} \end{aligned}$$

Sign Bit

Viktor Öwall, Dept. of Electrical and Information Technology, Lund University, Sweden - www.eit.lth.se

Example: 2's complement

$$1010.0110 \Rightarrow$$

$$\begin{aligned} & -1 \cdot 2^3 + 0 \cdot 2^2 + 1 \cdot 2^1 + 0 \cdot 2^0 + 0 \cdot 2^{-1} + 1 \cdot 2^{-2} + 1 \cdot 2^{-3} + 0 \cdot 2^{-4} = \\ & -8 + 2 + \frac{1}{4} + \frac{1}{8} \end{aligned}$$

Viktor Öwall, Dept. of Electrical and Information Technology, Lund University, Sweden - www.eit.lth.se

Sign Extension in Two's Complement

$$\begin{aligned} & -2^{k-1} a_{k-1} + 2^{k-2} a_{k-2} \dots 2 a_1 + a_0 = \\ & -2^k a_{k-1} + 2^{k-1} a_{k-1} + 2^{k-2} a_{k-2} \dots 2 a_1 + a_0 = \\ & -2^{k+1} a_{k-1} + 2^k a_{k-1} + 2^{k-1} a_{k-1} + 2^{k-2} a_{k-2} \dots 2 a_1 + a_0 \end{aligned}$$

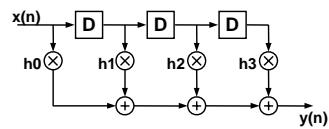
Example:

$$10010 = 110010 = 1110010 = 11110010 = \dots$$

$$00010 = 000010 = 0000010 = 00000010 = \dots$$

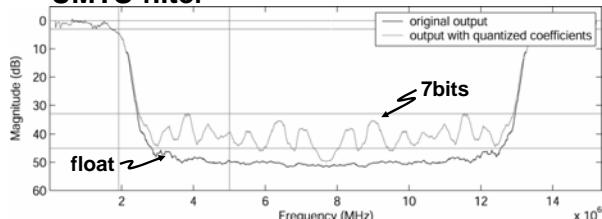
Viktor Öwall, Dept. of Electrical and Information Technology, Lund University, Sweden - www.eit.lth.se

The Wordlength, i.e. nr of bits

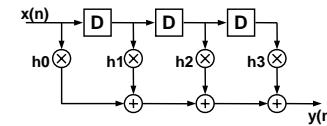


Every extra bit costs
 • energy/power
 • delay
 • area
 \Rightarrow the wordlength has to be reduced

UMTS-filter

Viktor Öwall, Dept. of Electrical and Information Technology, Lund University, Sweden - www.eit.lth.se

The Wordlength, i.e. nr of bits



The output of

- adder output needs an extra bit to be sure of no overflow, e.g.
decimal: $2+2 = 4 \Rightarrow$ binary: $10+10=100$
- multiplier $M \times N$ bits $\Rightarrow M+N$ bits for full precision

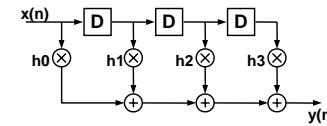
\Rightarrow Precision has to be limited

Viktor Öwall, Dept. of Electrical and Information Technology, Lund University, Sweden - www.eit.lth.se

Basic Building Blocks

Viktor Öwall, Dept. of Electrical and Information Technology, Lund University, Sweden - www.eit.lth.se

Basic Building Blocks



In the FIR filter

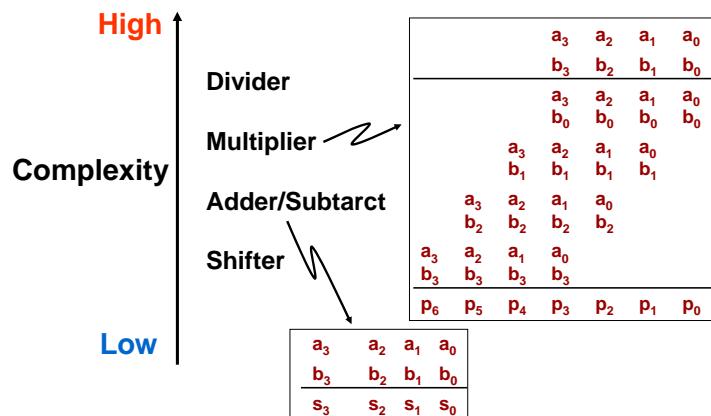
- adders
- multipliers
- registers

in other algorithms also: shift, minus, division,...

- left shift is multiply by 2
 - right shift is a divide by 2
- but is low complexity!

Viktor Öwall, Dept. of Electrical and Information Technology, Lund University, Sweden - www.eit.lth.se

Comparing Basic Building Blocks



Viktor Öwall, Dept. of Electrical and Information Technology, Lund University, Sweden - www.eit.lth.se

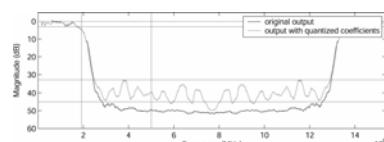
Scaling and Round-off Noise

Viktor Öwall, Dept. of Electrical and Information Technology, Lund University, Sweden - www.eit.lth.se

Quantization

Two Types

- Coefficient Quantization
 - Non-Ideal Transfer Function
 - Compare to analog component variations
- Signal Quantization
 - Round-off Noise
 - Limit Cycles



Viktor Öwall, Dept. of Electrical and Information Technology, Lund University, Sweden - www.eit.lth.se

Quantization

Round-off Noise

- Affect the output as a random disturbance

Limit Cycle Oscillations

- Undesired periodic components
- Due to non-linear behavior in the feedback (rounding or overflow)

Viktor Öwall, Dept. of Electrical and Information Technology, Lund University, Sweden - www.eit.lth.se

Quantization Analysis

Using “real” rounding, truncation, and overflow

- Give exact result
- Tricky - need integer representation

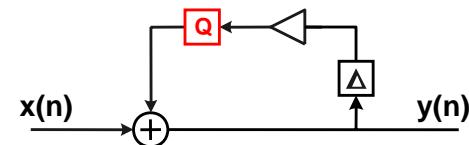
Using noise models

- Floating point representation can still be used
- Suitable for Matlab, C/C++ ...

Viktor Öwall, Dept. of Electrical and Information Technology, Lund University, Sweden - www.eit.lth.se

Rounding Truncation

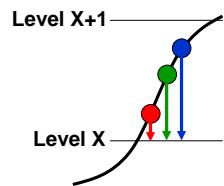
Rounding/Truncation is “always” there!
Especially necessary in recursive systems



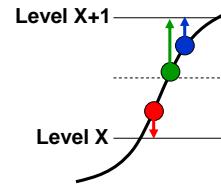
Without quantization - infinite wordlength
Multiplication $\Rightarrow n+m$ output bits
Addition $\Rightarrow n+1$ output bits

Viktor Öwall, Dept. of Electrical and Information Technology, Lund University, Sweden - www.eit.lth.se

Truncation and Rounding



Truncation
All values approximated in the same direction
Max error = 1 LSB

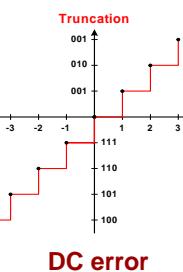


Avrundning
Values approximated up or down
Max error = 1/2 LSB

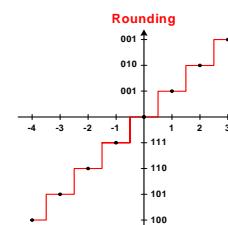
Viktor Öwall, Dept. of Electrical and Information Technology, Lund University, Sweden - www.eit.lth.se

Rounding Truncation

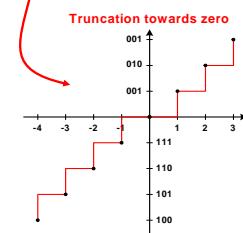
No energy added to the system
Often used in recursive algorithms



DC error



"Rounded to even"



Add LSB before truncation if negative

Viktor Öwall, Dept. of Electrical and Information Technology, Lund University, Sweden - www.eit.lth.se

Scaling

Adjust signal range to fit the hardware

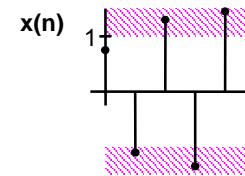
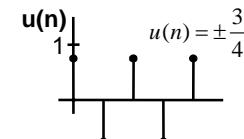
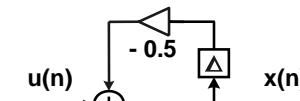
Unchanged transfer function (Scaled coefficients might move the pole-zeros)

Trade-off

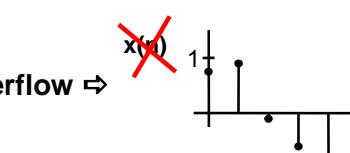
- Scale up to reduce roundoff noise
- Scale down to avoid overflow

But you loose precision!

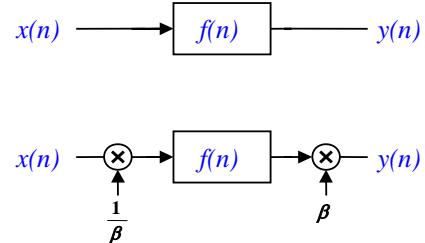
Example Where Scaling is Needed



Overflow \Rightarrow



Scaling



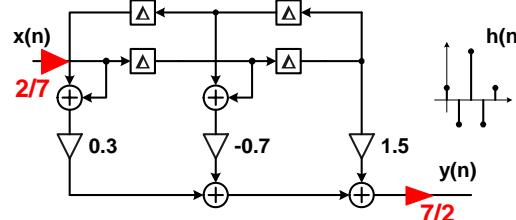
$$\text{Safe scaling if } \beta = \sum_{i=0}^{\infty} |f(i)|$$

Where $f(i)$ is the unit sample response

Example: Safe Scaling

$$\frac{2}{7}x(n) \text{ and } \frac{7}{2}y(n) \text{ give safe scaling}$$

$$\beta = \sum_{i=0}^{\infty} |f(i)| = 0.3 \times 2 + |-0.7 \times 2| + 1.5 = 3.5$$



Increased roundoff noise
Internal scaling might improve

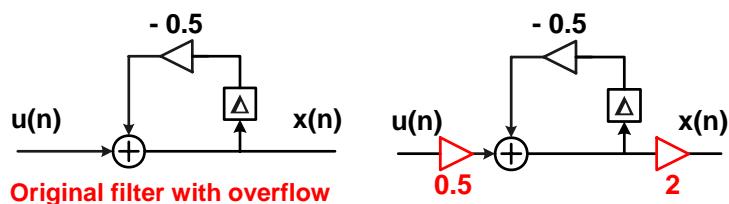
(Linear phase FIR. Note the strength reduction)

Example: Safe Scaling

$$\beta = \sum_{i=0}^{\infty} |f(i)| = \sum_{i=0}^{\infty} |(-0.5)^i| =$$

Geometric series

$$|(-0.5)^0| + |(-0.5)^1| + |(-0.5)^2| + \dots = \frac{1}{1 - 0.5} = 2$$

Viktor Öwall, Dept. of Electrical and Information Technology, Lund University, Sweden - www.eit.lth.se

Scaling

- Safe scaling is pessimistic
 - Alternative is scaling with
- $\beta = \sqrt{\sum_{i=0}^{\infty} (|f(i)|^2)}$

- In practice: Scaling with $\beta = 2^{\pm n}$

- Easy to do - a shift

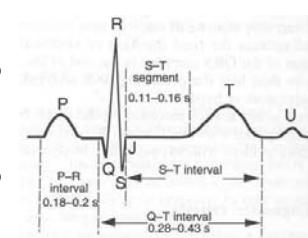
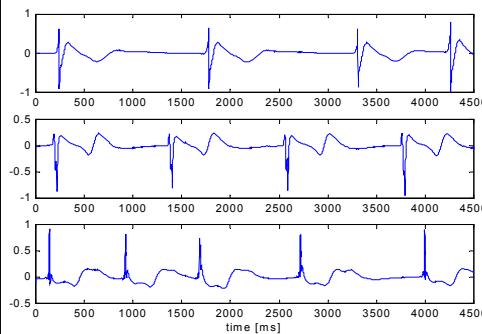
- Increased internal wordlength an alternative

Viktor Öwall, Dept. of Electrical and Information Technology, Lund University, Sweden - www.eit.lth.se

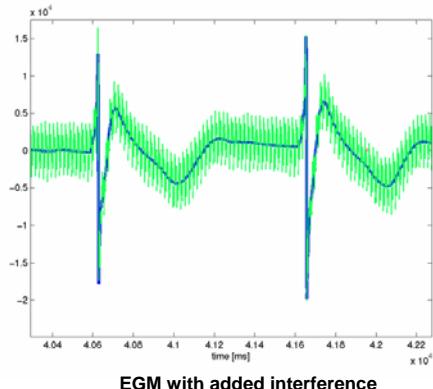
Pacemaker example

Viktor Öwall, Dept. of Electrical and Information Technology, Lund University, Sweden - www.eit.lth.se

The Electrocardiogram (EGM)

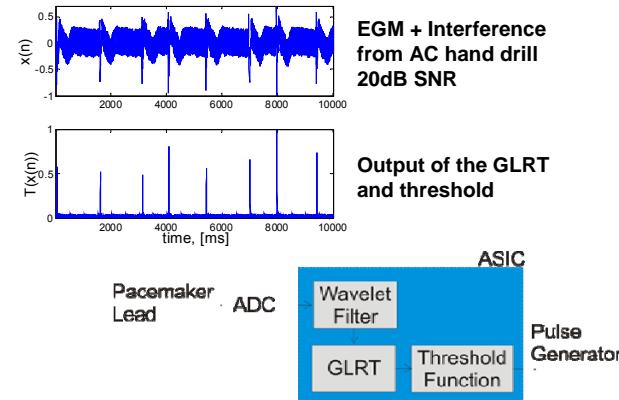
Viktor Öwall, Dept. of Electrical and Information Technology, Lund University, Sweden - www.eit.lth.se

The Interfered signal



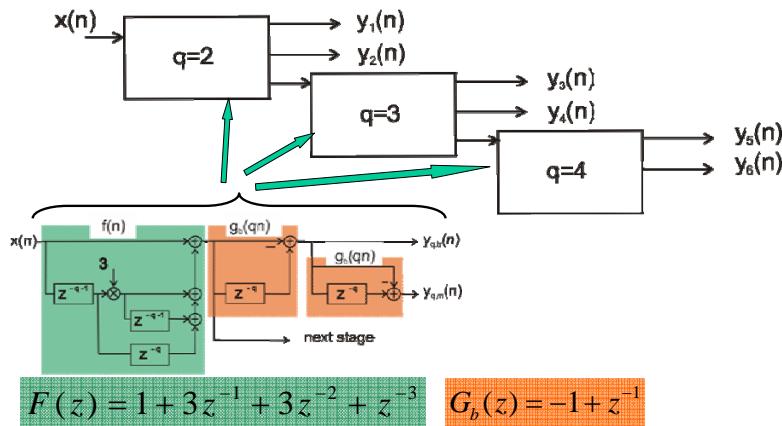
Viktor Öwall, Dept. of Electrical and Information Technology, Lund University, Sweden - www.eit.lth.se

Filtering Performance



Viktor Öwall, Dept. of Electrical and Information Technology, Lund University, Sweden - www.eit.lth.se

Wavelet Filterbank

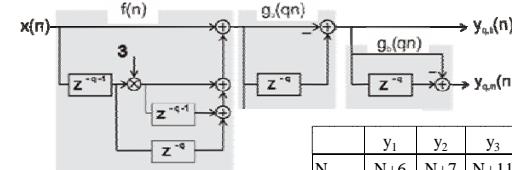


Viktor Öwall, Dept. of Electrical and Information Technology, Lund University, Sweden - www.eit.lth.se

Bit-optimization

- Signals have been monitored to determine the upper bound of the wordlength

Comparison of worst-case wordlength and implemented wordlength at the wavelet output:

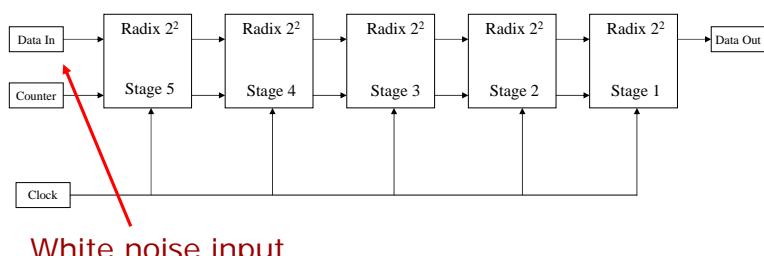


	y_1	y_2	y_3	y_4	y_5	y_6
N_{wc}	$N+6$	$N+7$	$N+11$	$N+12$	$N+14$	$N+15$
N_{imp}	$N+1$	$N+1$	$N+2$	$N+1$	$N+2$	$N+1$

Viktor Öwall, Dept. of Electrical and Information Technology, Lund University, Sweden - www.eit.lth.se

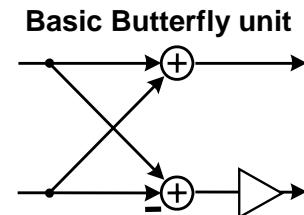
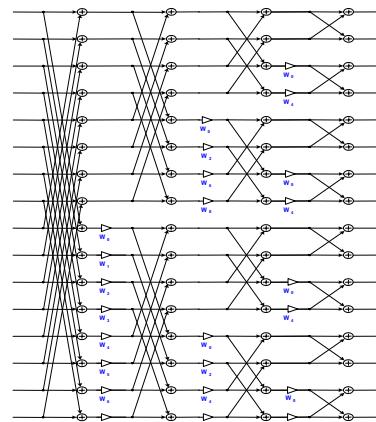
Example: Internal Scaling

- VHDL bit-level simulation
- Compared with Matlab floating-point simulation
- Optimized internal scaling



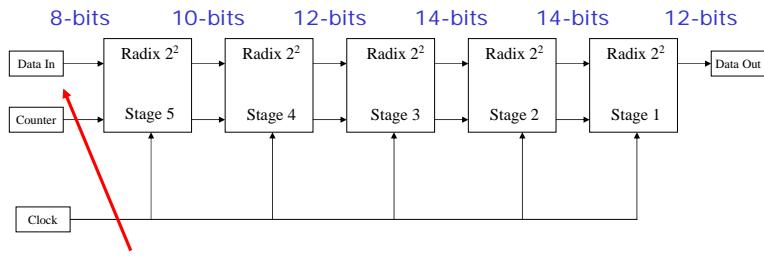
1024-point FFT

A 16-point Radix-2 FFT

Viktor Öwall, Dept. of Electrical and Information Technology, Lund University, Sweden - www.eit.lth.se

Example: Internal Scaling

- VHDL bit-level simulation
- Compared with Matlab floating-point simulation
- Optimized internal scaling



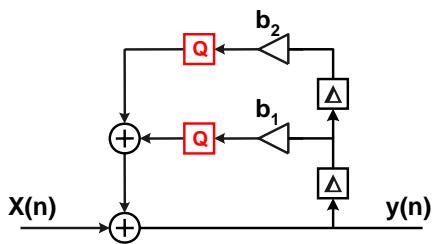
1024-point FFT

Limit Cycles

Viktor Öwall, Dept. of Electrical and Information Technology, Lund University, Sweden - www.eit.lth.se

Limit Cycles

Example: zero input oscillations in 2nd order IIR



$$b_1 = \frac{489}{256} = 1.91015625; b_2 = -\frac{15}{16} = -0.9375$$

Viktor Öwall, Dept. of Electrical and Information Technology, Lund University, Sweden - www.eit.lth.se

Example: zero input oscillations

Limit Cycles

Zero Input



Rounding after multiplication

Truncation after multiplication

with DC offset

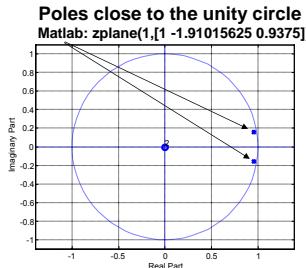
Source: Lars Wanhammar, "DSP Integrated circuits"

Viktor Öwall, Dept. of Electrical and Information Technology, Lund University, Sweden - www.eit.lth.se

Limit Cycles

Zero input oscillations

- Often not accepted in audio



Very difficult problem

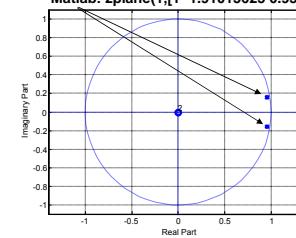
- In general, no solutions for structures > 2nd order
- Can be limited by increased internal wordlength
- Can in some 2nd order structures be eliminated by pole positioning
- 2nd order Wave Digital Filters are free from parasitic oscillations

Viktor Öwall, Dept. of Electrical and Information Technology, Lund University, Sweden - www.eit.lth.se

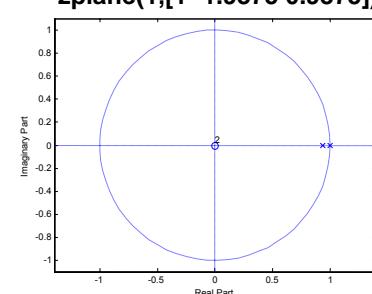
Limit Cycles

Changing the precision move the poles!

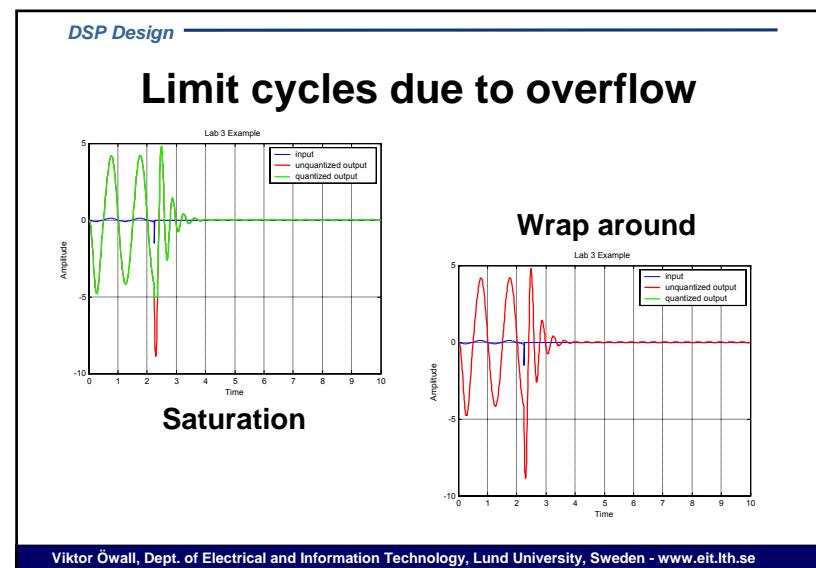
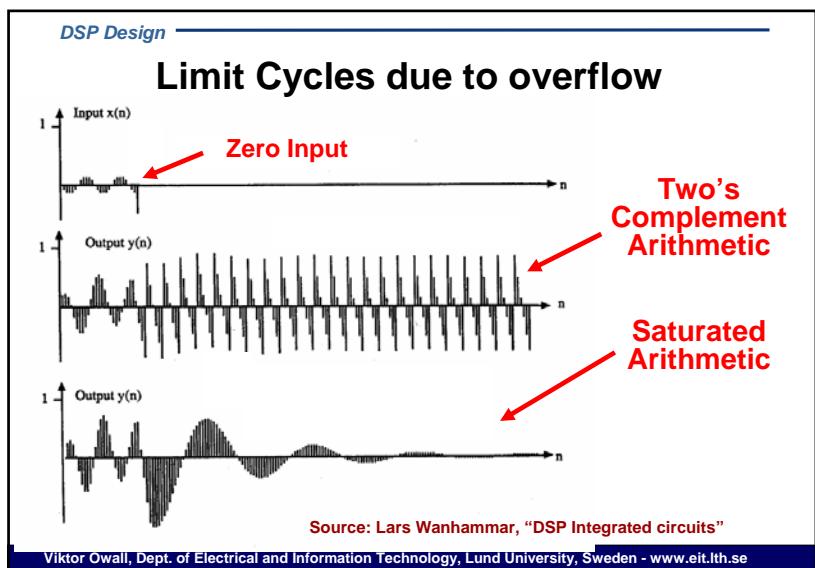
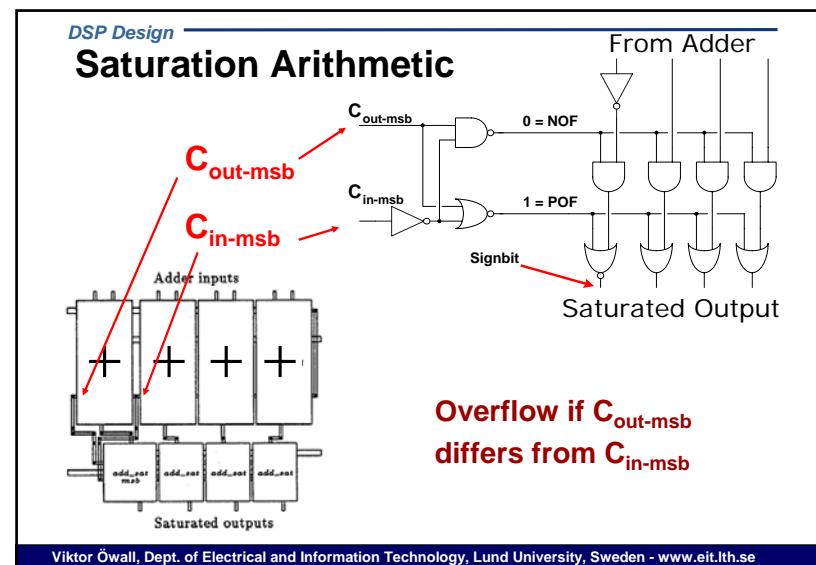
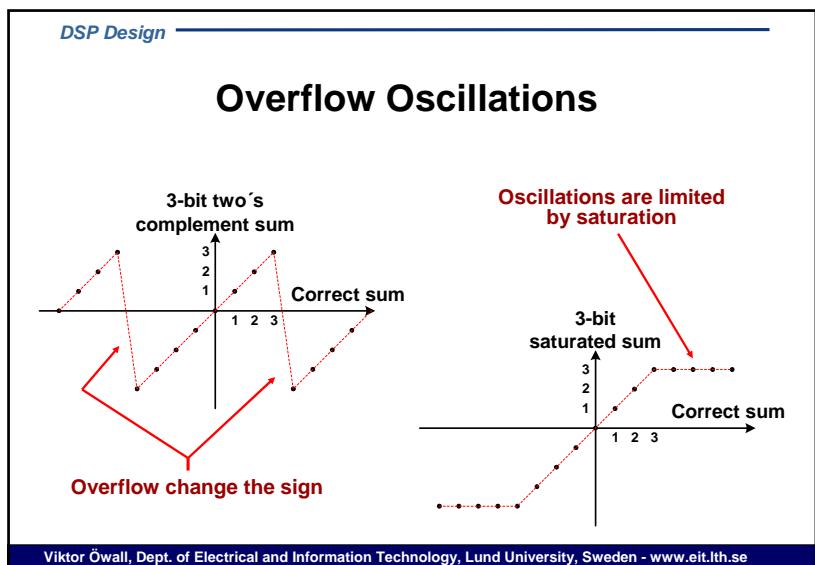
Poles close to the unity circle
Matlab: zplane(1,[1 -1.9375 0.9375])



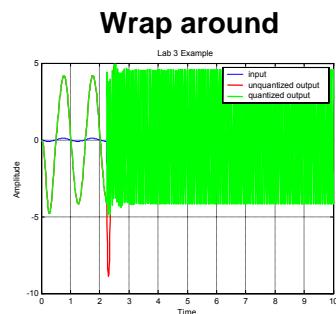
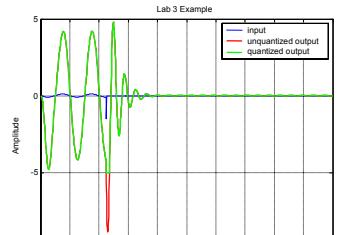
zplane(1,[1 -1.9375 0.9375])



Viktor Öwall, Dept. of Electrical and Information Technology, Lund University, Sweden - www.eit.lth.se



Limit cycles due to overflow



Viktor Öwall, Dept. of Electrical and Information Technology, Lund University, Sweden - www.eit.lth.se

Simple Noise Analysis

Viktor Öwall, Dept. of Electrical and Information Technology, Lund University, Sweden - www.eit.lth.se

Scaling and White Noise Input

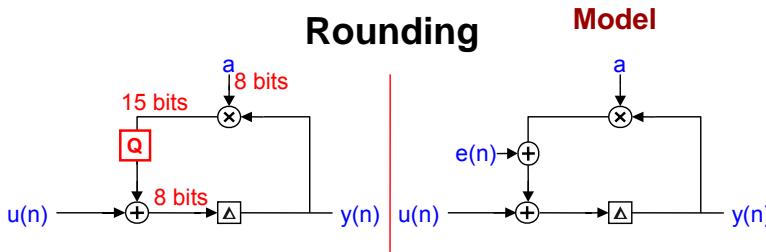
$$\beta = \sum_{i=0}^{\infty} |f(i)|, \text{Safe-scaling}$$

$$\beta = \delta \sqrt{\sum_{i=0}^{\infty} f^2(i)}, \text{possible overflow}$$

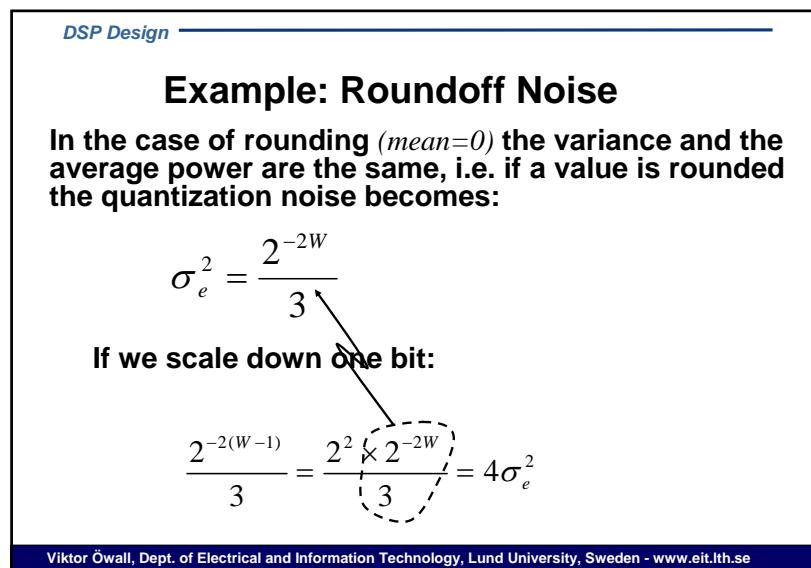
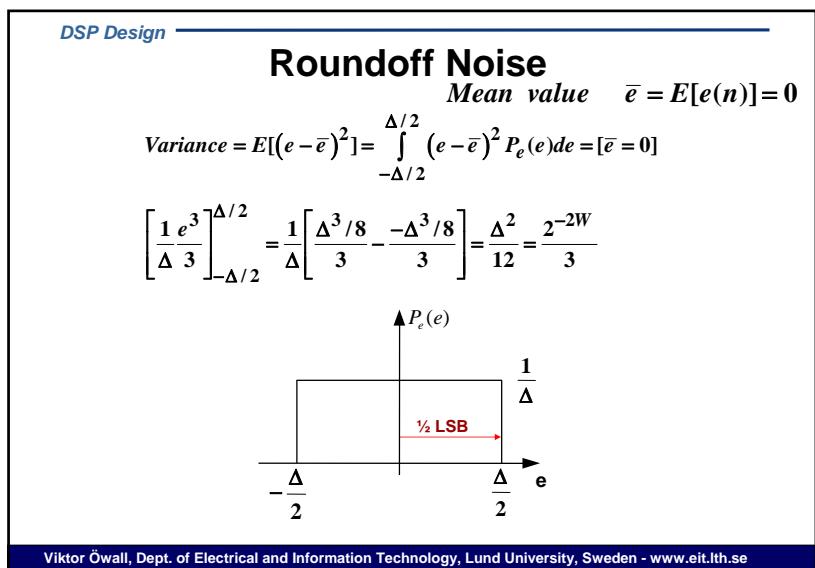
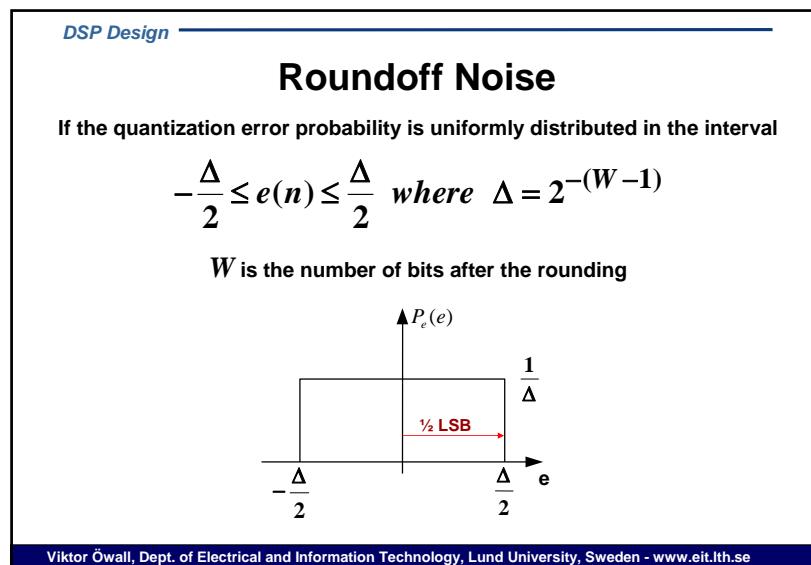
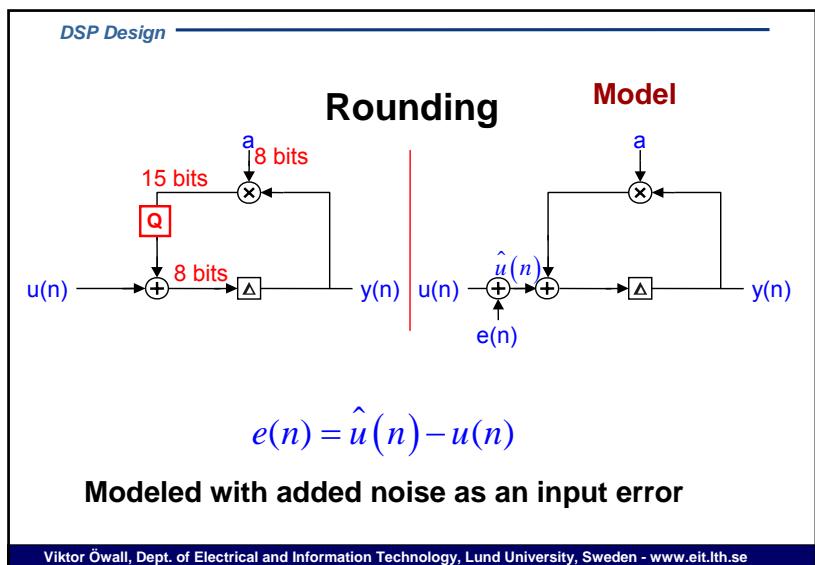
$f(i)$ = unit sample response, $\sum_{i=0}^{\infty} f^2(i)$ = Variance white noise input

- “Safe scaling” but not guaranteed
- δ sets the probability for an overflow
- Typically one overflow every 10^6 sample is accepted in audio [Wanhammar]

Viktor Öwall, Dept. of Electrical and Information Technology, Lund University, Sweden - www.eit.lth.se



Viktor Öwall, Dept. of Electrical and Information Technology, Lund University, Sweden - www.eit.lth.se



Signal to Noise Ratio (SNR)

One extra bit reduces quantization error by a factor 4

$$SNR = 10 \log \frac{4\sigma_e^2}{\sigma_e^2} = 6.02 \text{ dB}$$

Good to remember: 6 dB increase in SNR per bit

Signal to Noise Ratio (SNR)

Signal power (variance)

$$SNR = 10 \log \frac{\sigma_x^2}{\sigma_e^2} = 10 \log \frac{3}{2^{-2W}} \sigma_x^2$$

Roundoff error power (variance)

Signal to Noise Ratio (SNR)

Example: Full scale sinus wave rounded to 8 bits

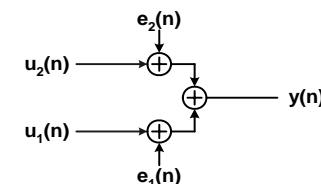
$$SNR = 10 \log \frac{3}{2^{-2 \times 8}} \left(\frac{A}{\sqrt{2}} \right)^2 = 50 \text{ dB}; -1 \leq A \leq 1$$

Roundoff Noise: Addition

$$E[(e_1 + e_2)^2] = E[e_1^2 + 2e_1e_2 + e_2^2] =$$

$$= E[e_1^2] + \underbrace{E[2e_1e_2]}_{\substack{\text{zero if} \\ u_1 \text{ and } u_2 \\ \text{independent}}} + E[e_2^2] =$$

u₁ and u₂ independent



$$= E[e_1^2] + E[e_2^2]$$

Example: Roundoff Noise

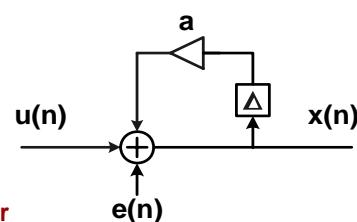
First order IIR-filter, the variance is:

$$\sigma_e^2 \sum_{i=0}^{\infty} f^2(i) = \sigma_e^2 (1 + (a^1)^2 + (a^2)^2 + (a^3)^2 \dots) = \sigma_e^2 \frac{1}{1 - a^2}$$

$$a = 0.100 \Rightarrow 1.01 \sigma_e^2$$

$$a = 0.500 \Rightarrow 1.33 \sigma_e^2$$

$$a = 0.998 \Rightarrow 500 \sigma_e^2$$



Narrow band filter

Viktor Öwall, Dept. of Electrical and Information Technology, Lund University, Sweden - www.eit.lth.se

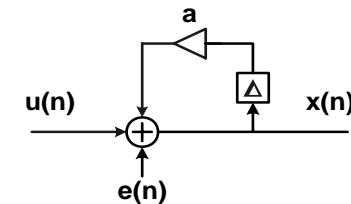
Example: SNR

Example: Full scale sinus,
rounded to 8 bits in IIR

$$\sigma_e^2 \sum_{i=0}^{\infty} f^2(i) = \sigma_e^2 \frac{1}{1 - a^2}$$

No feedback $\Rightarrow SNR = 50\text{dB}$

$$a = 0.998 \Rightarrow 500 \sigma_e^2 \Rightarrow SNR = 23\text{dB}$$



Viktor Öwall, Dept. of Electrical and Information Technology, Lund University, Sweden - www.eit.lth.se