

# Digital Signal Processing

## ESD-5 & IV-5 (elektro), E24

### 6. Finite word length effects – a brief intro

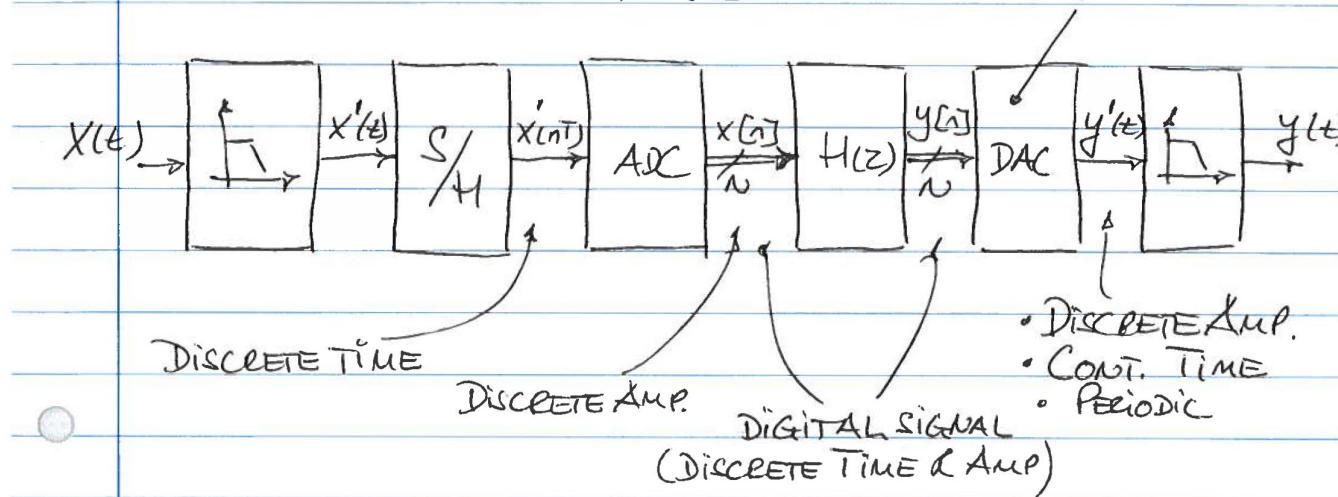
(Today, partly handwritten slides...)

Assoc. Prof. Peter Koch, AAU

# REALISATION OF DIGITAL FILTERS

## THE SIGNAL CHAIN

ZERO ORDER HOLD



THROUGH THE SIGNAL CHAIN, SEVERAL ERRORS — OR QUANTIZATIONS — ARE INTRODUCED.

S/H : TIME QUANTIZATION

ADC : VARIABLE QUANTIZATION

$H(z)$  (COMPUTER) : VARIABLE + COEFFICIENT QUANT.

DAC : TIME QUANTIZATION.

# Time Quantization

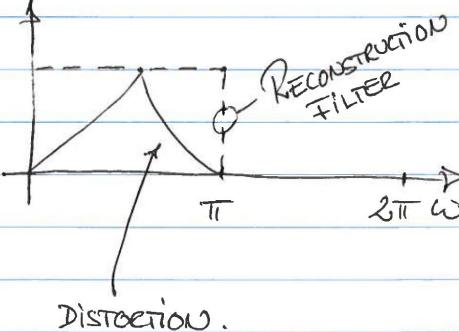
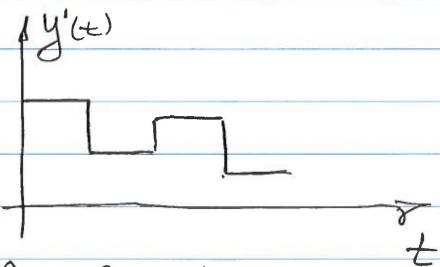
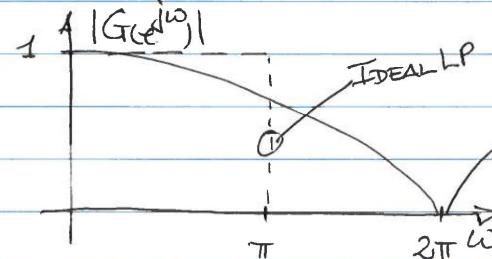
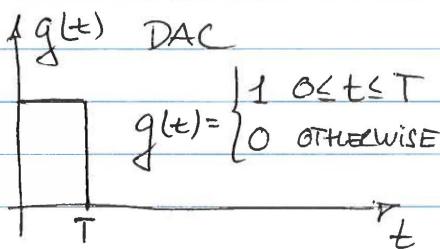
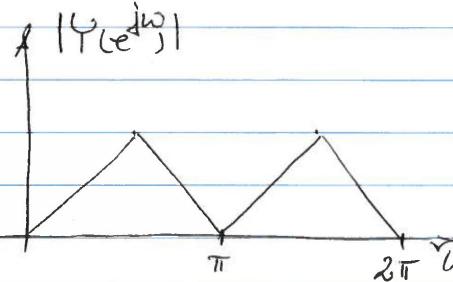
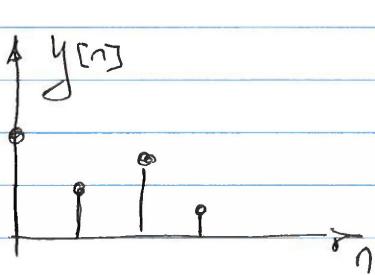


AALBORG UNIVERSITY  
DENMARK

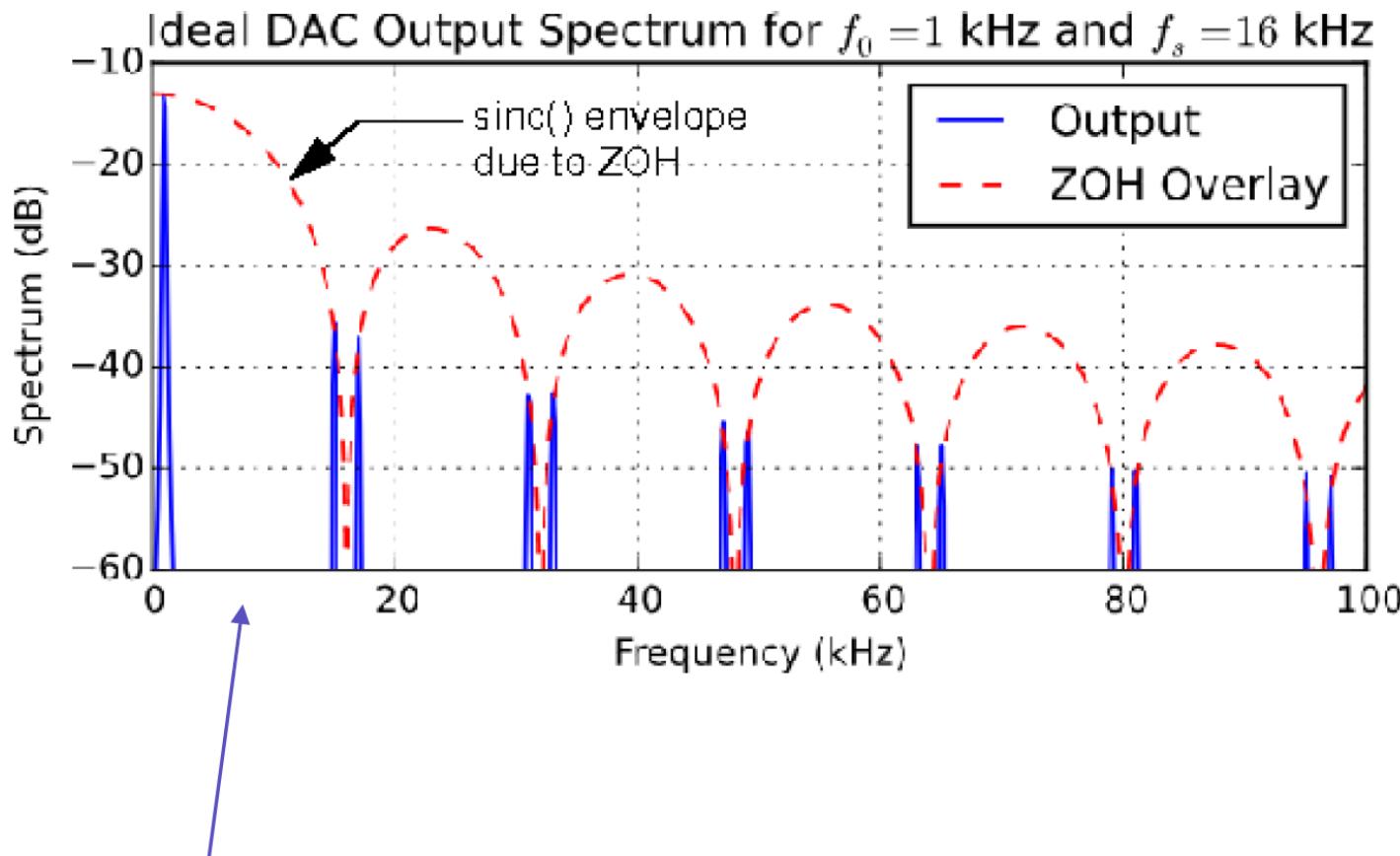
## TIME QUANTIZATION

- $f_{\text{sample}} > 2 \cdot f_{\text{max, analog}}$  (Reconstruction).
- ALWAYS APPLY AN ANTI-AIUSING FILTER.
- DISCRETE TIME  $\Rightarrow$  PERIODIC FREQUENCY.

- THE DAC IS NOT AN IDEAL LP.FILTER



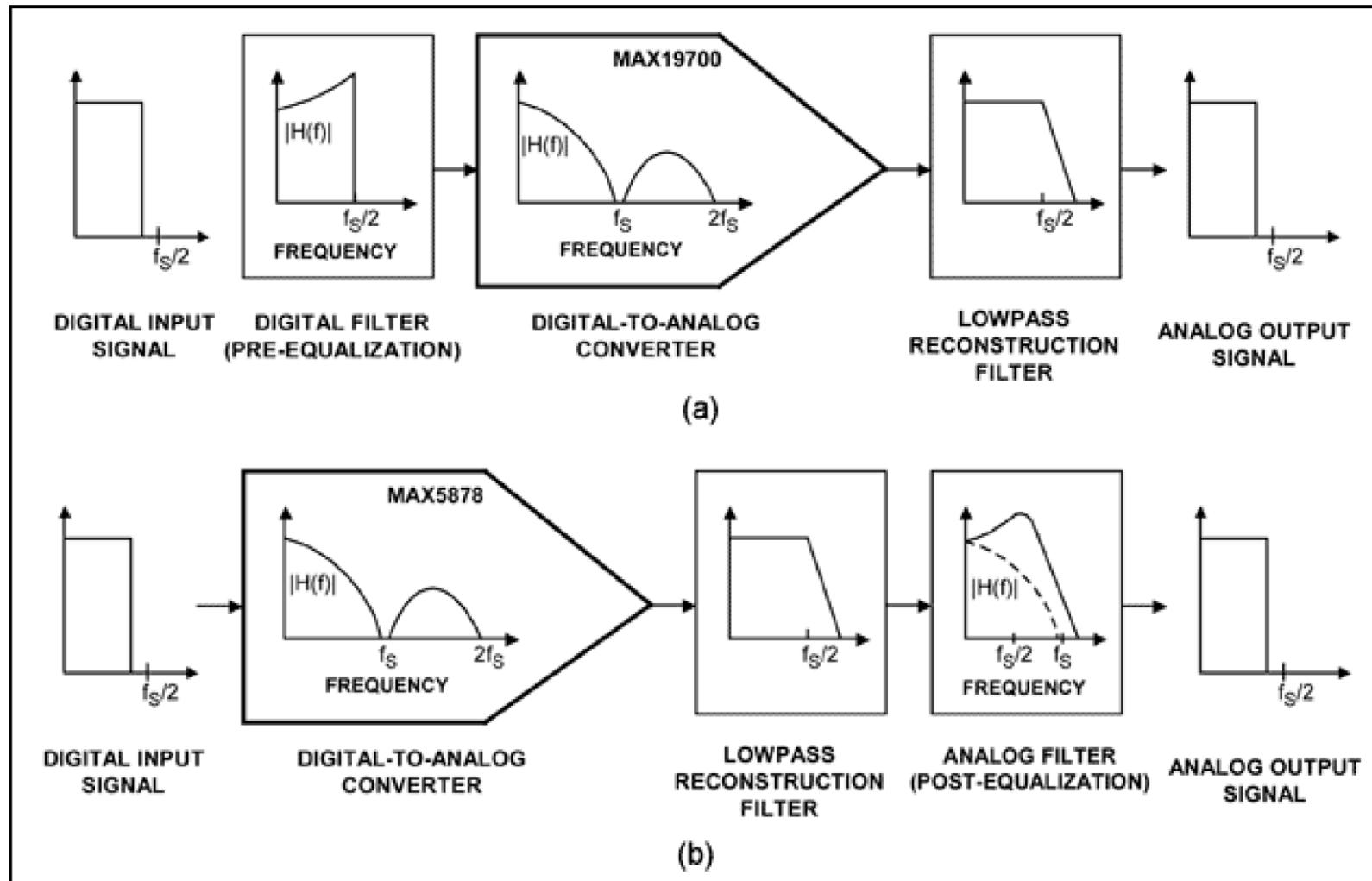
The Sinc shaped ZOH Overlay maintains all the  $2\pi$  periodic images from the output signal  $y[n]$ .



Here we see clearly the need for a reconstruction LP-filter with a cut-off frequency at  $f_s/2$

For frequency components close to the Nyquist frequency, the ZOH introduces distortion.

To compensate for this ZOH distortion we may apply either pre- or post-equalization



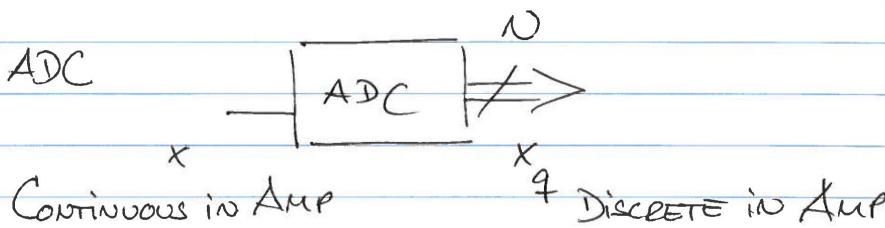
# Variable Quantization



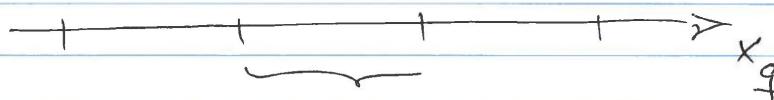
AALBORG UNIVERSITY  
DENMARK

## VARIABLE QUANTIZATION

• ADC



$N$  bit ADC  $\Rightarrow 2^N$  POSSIBLE VALUES FOR  $x_q$



QUANTIZATION STEP

$$\Delta = 1/2^N$$

So, AT THE ADC OUTPUT WE HAVE A FINITE SIGNAL-TO-NOISE RATIO (SNR)

$$SNR \triangleq 20 \cdot \log \frac{\text{SIGNAL (rms)}}{\text{NOISE (rms)}}$$

WITHOUT PROOF;

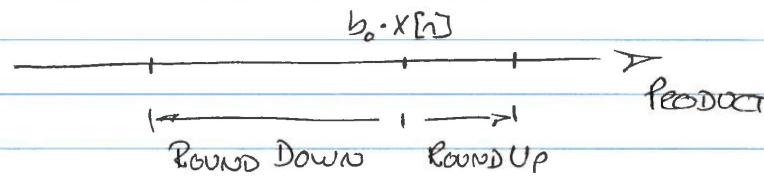
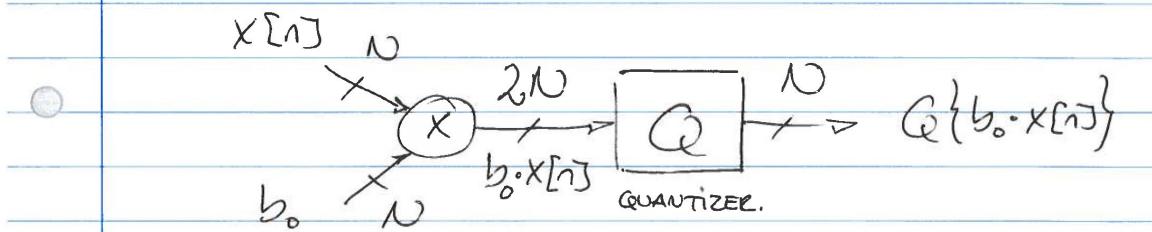
THE SNR ON THE OUTPUT OF THE ADC IS INCREASED WITH 6dB WHEN N IS INCREASED WITH 1.

$$SNR = 6.02N + 1.76 \text{ dB}$$

•  $H(z)$  (Computer)

MULTIPLICATION BETWEEN SIGNAL AND COEF.

$$y[n] = \sum_{k=1}^N a_k \cdot y[n-k] + \sum_{l=0}^M b_l x[n-l]$$



• So, THE QUANTIZER which TAKES US BACK FROM  $2N$  BIT TO  $N$  BIT, INTRODUCES AN ERROR ;

$$e[n] = b_0 \cdot x[n] - Q\{b_0 \cdot x[n]\}$$

$2N$  BIT       $N$  BIT

Q

• THE COMPUTATION ITSELF MAKES NOISE →

## A FEW WORDS ON NUMBER REPRESENTATION.

### • FLOATING POINT

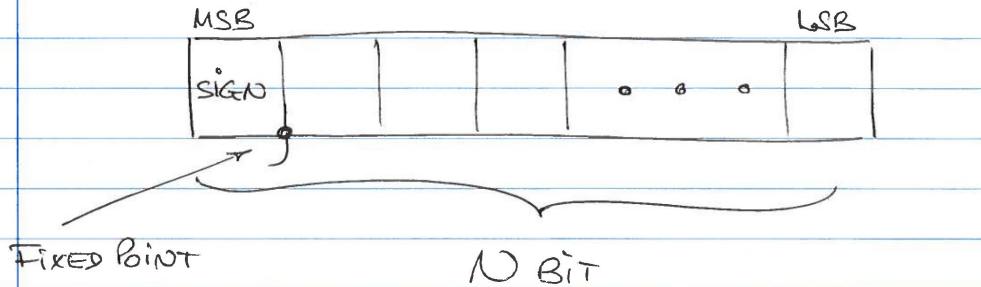
LARGE DYNAMIC RANGE ✓  
TIME, AREA, POWER %

### • FIXED POINT

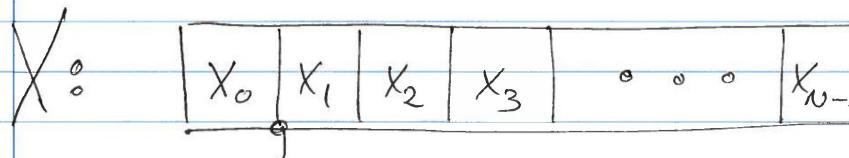
LARGE DYNAMIC RANGE %  
TIME, AREA, POWER ✓

SIGNIFICANTLY USED IN SIGNAL PROCESSING IN

### • TERMS OF 2<sup>s</sup> COMPLEMENT NUMBER REPRESENTATION.



$$\text{SIGN} = \begin{cases} 0 & \rightarrow \text{POS. NUMBER} \\ 1 & \rightarrow \text{NEG. NUMBER.} \end{cases}$$



THE NUMERICAL VALUE OF X;

$$X = -x_0 + \sum_{i=1}^{n-1} x_i \cdot 2^{-i}$$

# MULTIPLYING TWO N BIT NUMBERS

$$P = X \cdot Y = X \cdot \left( -y_0 + \sum_{i=1}^{N-1} y_i \cdot 2^{-i} \right)$$

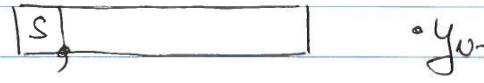
N bit X



+



+



+

..

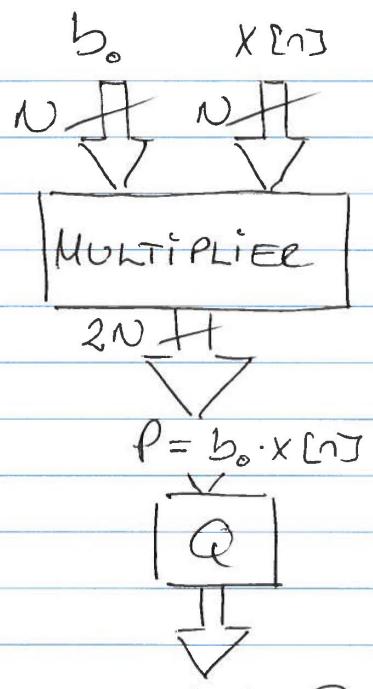
-



2N bit

- ACTUALLY, THE PRODUCT P ONLY HAS  $2N-1$  BIT WHICH IS COMPENSATED FOR VIA LEFT SHIFT (0 INTO LSB POSITION)

- FURTHERMORE, SIGN EXTENSION IS NEED FOR THE ADDITION.

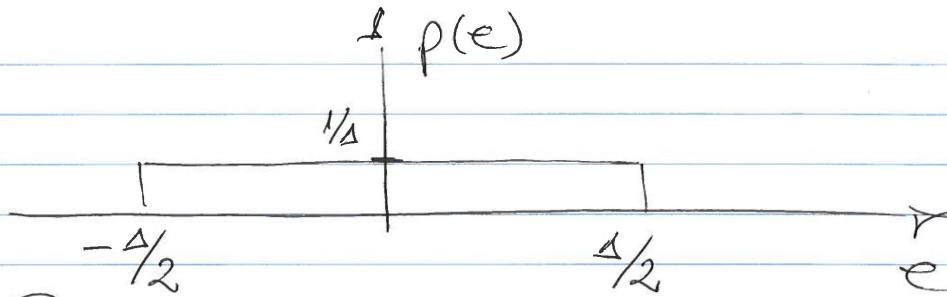


CAN WE SAY  
ANYTHING ABOUT  
THIS ERROR ??

DIFFICULT, BUT  
WE CAN MAKE  
SOME ASSUMPTIONS

- ①  $e[n]$  IS A STATIONARY WHITE NOISE PROCESS
- ② THE AMPLITUDE OF  $e[n]$  IS EVENLY DISTRIBUTED OVER ONE QUANTIZATION STEP.
- ③  $e[n]$  AND  $x[n]$  ARE UN-CORRELATED.

USING THESE ASSUMPTIONS WE CAN SAY SOMETHING DEFINITELY ABOUT THE STATISTICS OF  $e[n]$



### Probability Density Function

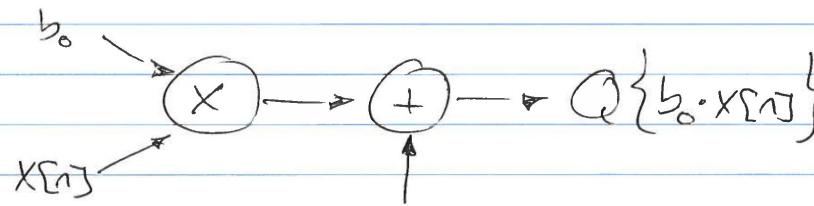
#### MEAN

$$\mu_e = \int_{-\infty}^{\infty} e \cdot p(e) \cdot de = \frac{1}{\Delta} \left[ \frac{e^2}{2} \right]_{-\Delta/2}^{\Delta/2} = 0$$

#### VARIANCE

$$\sigma_e^2 = \int_{-\infty}^{\infty} e^2 \cdot p(e) \cdot de = \frac{1}{\Delta} \left[ \frac{e^3}{3} \right]_{-\Delta/2}^{\Delta/2} = \frac{\Delta^2}{12}$$

USING THESE STATISTICAL VALUES, WE CAN NOW DERIVE A MODEL FOR THE QUANTIZATION NOISE



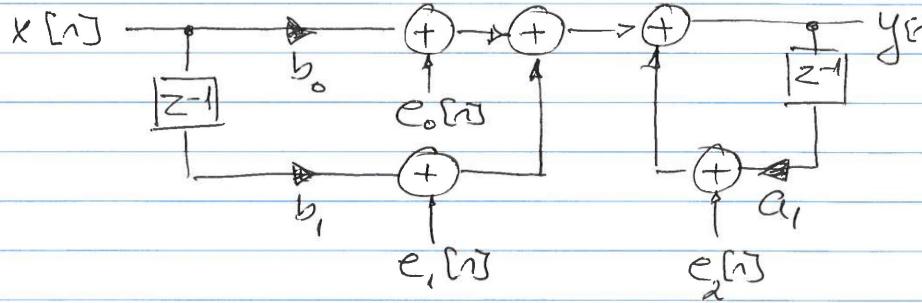
$e[n]$

WHERE  $e[n]$  IS A WHITE SIGNAL  $(\mu_e, \sigma_e^2)$

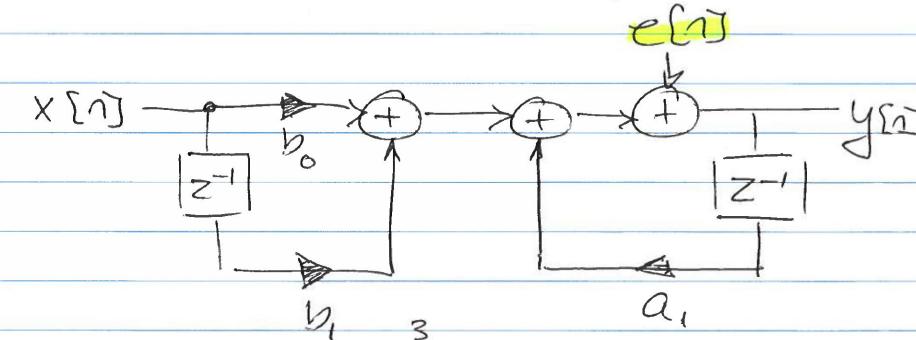
EXAMPLE : 1<sup>ST</sup> ORDER IIR FILTER.

$$y[n] = a_1 \cdot y[n-1] + b_0 \cdot x[n] + b_1 \cdot x[n-1]$$

GRAPHICAL REPRESENTATION ;



SINCE THESE ARE ALL LINEAR OPERATIONS, WE CAN RE-ORGANIZE THE GRAPH ;



WHERE  $e[n] = \sum_{i=1}^3 e_i[n]$ , AND THUS ;

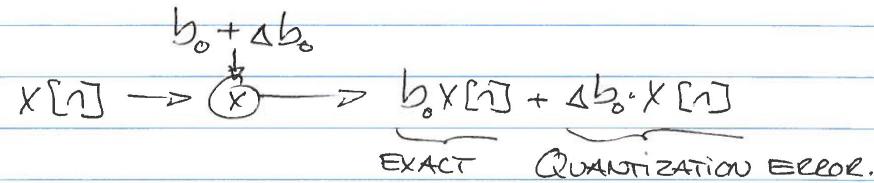
- ① THE TOTAL NOISE IMPACTS DIRECTLY THE OUTPUT.
- ② AT THE OUTPUT THE SNR IS NOT INFINITE

# Coefficient Quantization



AALBORG UNIVERSITY  
DENMARK

## COEFFICIENT QUANTIZATION (in $H(z)$ )



WHERE  $b_0$  IS THE EXACT COEFFICIENT (CALCULATED WITH MANY DIGITS) AND  $\Delta b_0$  IS THE ERROR INTRODUCED ON THE COEFFICIENT WHEN IT IS REPRESENTED IN  $N$  BIT

$$\hat{H}(z) = \frac{\sum_{k=0}^M \hat{b}_k z^{-k}}{\sum_{l=1}^N \hat{a}_l z^{-l}} = \frac{\hat{B}(z)}{\hat{A}(z)}$$

SINCE, GENERALLY,  $\hat{b}_k \neq b_k$  AND  $\hat{a}_l \neq a_l$

WE WILL SEE SOME CHANGES IN THE POLE / ZERO LOCATIONS.

ONE POSSIBLE METHOD TO INVESTIGATE HOW SEVERE THIS PROBLEM IS, IS TO MAKE AN ANALYSIS OF THE POLE-SENSITIVITY

Sensitivity Analysis (without proof)

$$A(z) = 1 - \sum_{k=1}^N a_k z^{-k} = \prod_{k=1}^N (1 - d_k z^{-1})$$

WHERE  $d_k$  ARE THE ROOTS IN  $A(z)$  (POLES).

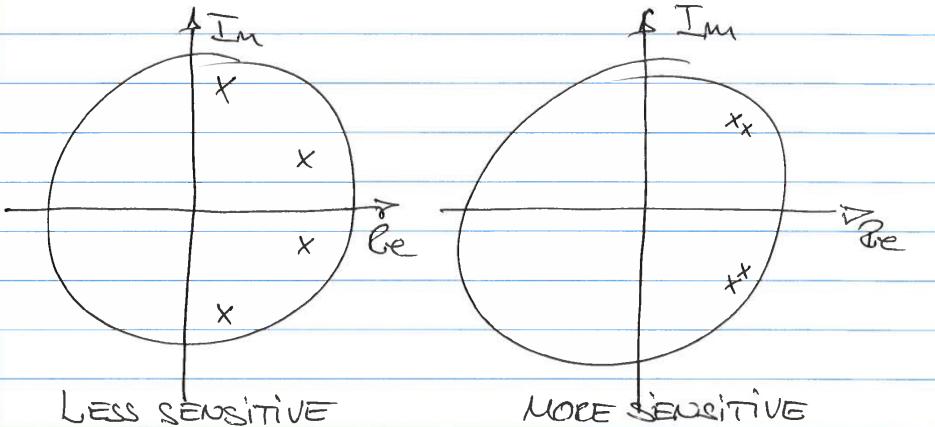
BASICALLY, WE WANT TO KNOW „THE CHANGE IN THE POLES AS RELATED TO THE CHANGE IN THE FILTER COEFFICIENTS“;

$$\frac{\partial d_i}{\partial a_k} = \frac{d_i^{N-k}}{\prod_{\substack{j=1 \\ j \neq i}}^N (d_i - d_j)}$$

(THIS EQUATION IS NOT SHOW IN THE BOOK)

SO, WHAT WE SEE IS THAT A CHANGE IN THE POLE LOCATION AS RELATED TO A CHANGE IN THE COEFFICIENTS (DUE TO QUANTIZATION) IS INVERSE PROPORTIONAL TO THE DISTANCE BETWEEN THE POLES

So, the closer together the poles are located, the more sensitive the pole location is due to a change in the filter coefficients



CAN WE DO SOMETHING ABOUT THIS PROBLEM ?

YES WE CAN .

WE MAY RE-ORGANIZE  $H(z)$  IN SUCH A WAY THAT CLOSELY SPACES POLES ARE DE-COUPLED INTO INDIVIDUAL FILTER SECTIONS

TAKE THE EXAMPLE ABOVE ;

$$\begin{aligned}
 H(z) &= \frac{B(z)}{A(z)} = \frac{B(z)}{1 - \sum_{k=1}^4 a_k z^{-k}} = \frac{B(z)}{\left(1 - \sum_{l=1}^3 a_l z^{-l}\right)\left(1 - \sum_{m=1}^2 a_m z^{-m}\right)} \\
 &= \frac{B(z)}{A_1(z) \cdot A_2(z)}
 \end{aligned}$$

- So, the trick is to re-organize  $H(z)$  into a cascade of 2<sup>nd</sup> order sections (evt. also one 1<sup>st</sup> order section).

$$-\boxed{H_1(z)} - \boxed{H_2(z)} - \boxed{H_3(z)} - \dots - \boxed{H_{N/2}(z)}$$

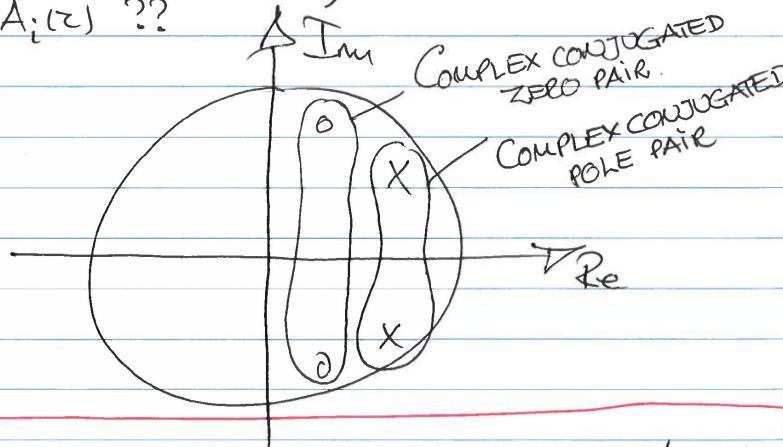
↓

$$H(z)$$

$$H(z) = \prod_{i=0}^{N/2} H_i(z) \quad (\text{Factorization})$$

$$\text{WHERE } H_i(z) = \frac{B_i(z)}{A_i(z)}$$

Now, question is ; Which zeros from  $H(z)$  (i.e., roots from  $B(z)$ ) should be paired with  $A_i(z)$  ??



The idea is that closely spaced pole/zero pairs are merged into one section  $H_i(z)$

NEXT QUESTION ;

IS  $\boxed{H_1(z)}$  —  $\boxed{H_2(z)}$  IDENTICAL TO

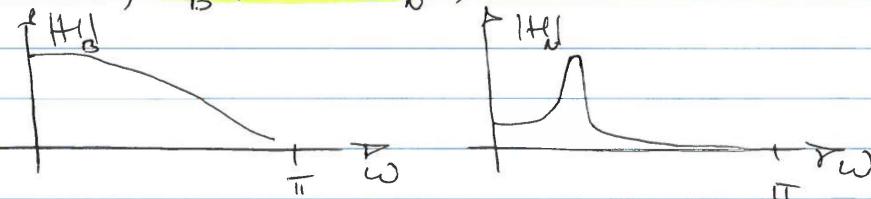
$\boxed{H_2(z)}$  —  $\boxed{H_1(z)}$  ?

IF WE DO THE COMPUTATION IN INFINITE WORDLENGTH (FLOATING POINT), THEN THE ANSWER IS "YES".

HOWEVER, FOR FIXED POINT COMPUTATION, THE ANSWER MOST LIKELY IS "NO".

HENCE, WHICH ONE TO CHOOSE ?

TO ANSWER THIS QUESTION, WE DISTINGUISH BETWEEN BROAD-BAND AND NARROW-BAND SECTIONS;  $H_B(z)$  AND  $H_N(z)$



1)  $H_N(z)$  FIRST — REDUCTION OF NOISE SPECTRUM THROUGH THE CASCADE

2)  $H_N(z)$  LAST — THE AMPLITUDE LEVEL IS REDUCED ONLY AT THE OUTPUT OF THE FILTER

RECOMMENDED

# Pole Locations in a $2^{nd}$ Order Section



AALBORG UNIVERSITY  
DENMARK

LET'S HAVE A LOOK AT A 2<sup>nd</sup> ORDER SYSTEM;

$$H(z) = \frac{1}{A(z)} = \frac{1}{\sum_{l=0}^{\infty} a_l z^{-l}}$$

$$= \frac{1}{1 + a_1 z^{-1} + a_2 z^{-2}}$$

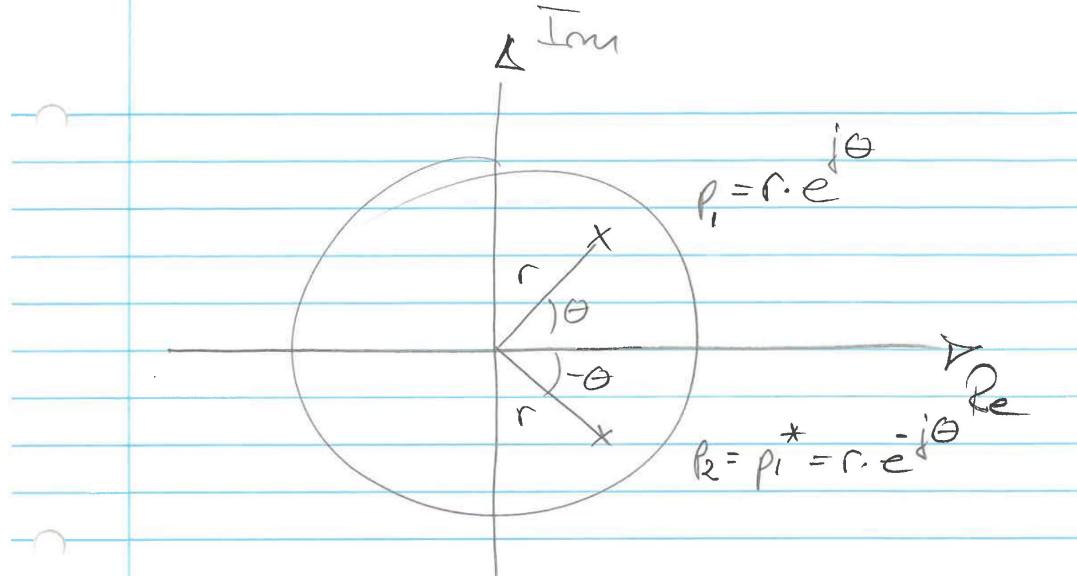
NOTE THE  
NORMALIZED  
POLYNOMIUM,  
 $a_0 = 1.$

$$= \frac{z^2}{z^2 + a_1 z + a_2}$$

$$= \frac{z^2}{(z - p_1)(z - p_2)}$$

$$= \frac{1}{(1 - p_1 z^{-1})(1 - p_2 z^{-2})}$$

WHERE  $p_2 = p_1^*$  (OTHERWISE WE GET  
COMPLEX COEFFICIENTS  $a_i$ )



Using Euler we can rewrite the poles into;

$$p_1 = r(\cos \theta + j \sin \theta)$$

$$p_2 = r(\cos \theta - j \sin \theta)$$

$$\begin{aligned}
 H(z) &= \frac{z^2}{(z-p_1)(z-p_2)} \\
 &= \frac{z^2}{(z-r(\cos \theta + j \sin \theta))(z-r(\cos \theta - j \sin \theta))}
 \end{aligned}$$

$$= \frac{z^2}{z^2 + r^2 \cos^2 \theta - 2zr \cos \theta + r^2 \sin^2 \theta}$$

$$= \frac{z^2}{z^2 + r^2 - 2zr \cos \theta}$$

$$= \frac{1}{1 - 2r \cos \theta \cdot z^{-1} + r^2 z^{-2}} = \frac{Y(z)}{X(z)}$$

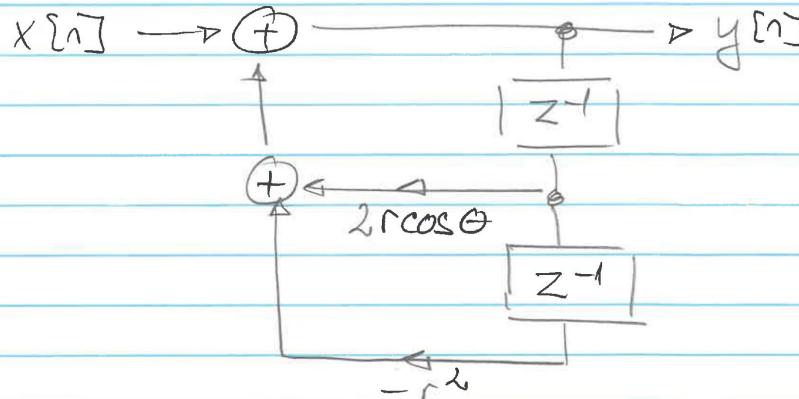
$$Y(z) = X(z) + Y(z)(2r \cos \theta z^{-1} - r^2 z^{-2})$$

$$z^{-1}$$

$$y[n] = x[n] + 2r \cos \theta y[n-1] - r^2 y[n-2]$$

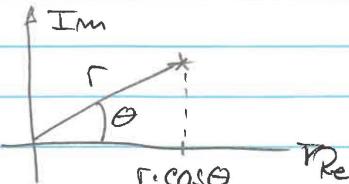
$$y[n] = x[n] + 2r \cos \theta y[n-1] - r^2 y[n-2]$$

LET'S NOW USE A DIRECT FORM :



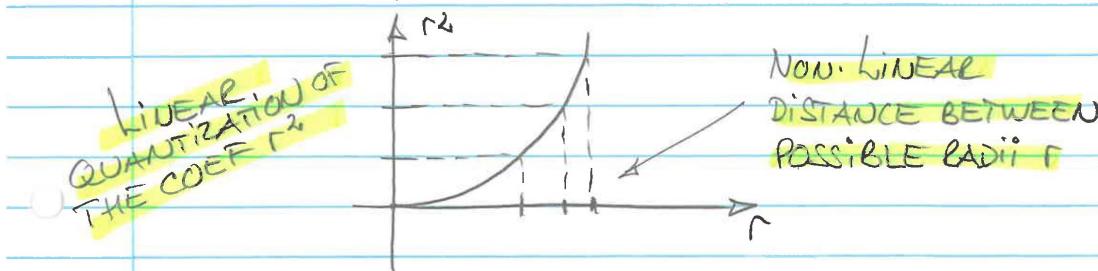
IN THIS STRUCTURE, THE COEFFICIENTS  $a_1$  AND  $a_2$  ARE  $2r \cos \theta$  AND  $-r^2$  WHICH NOW MUST BE QUANTIZED

$$\bullet r \cdot \cos \theta = \operatorname{Re}\{p_1\} = \operatorname{Re}\{p_2\}$$

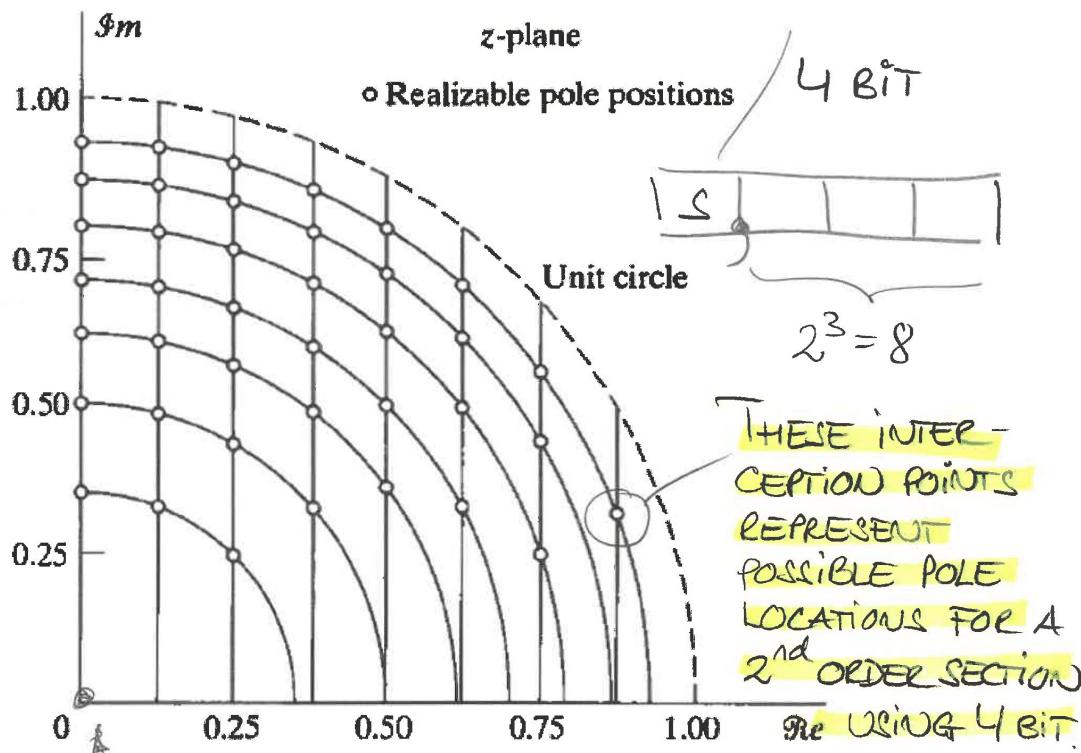
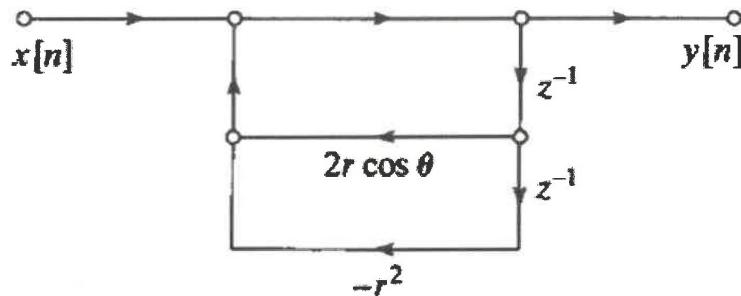


GIVEN  $N$  BIT, WE CAN QUANTIZE THE REAL AXIS (POSITIVE AND NEGATIVE) INTO  $2^N$  STEPS, EACH EQUAL TO  $2^{-N}$

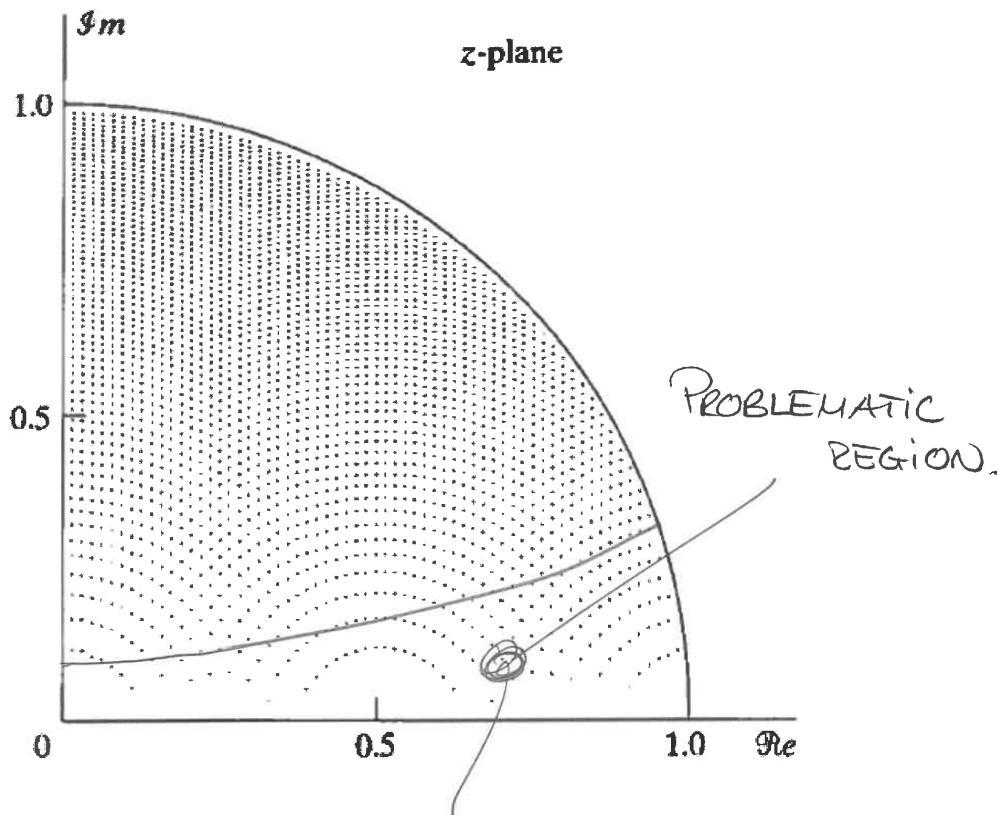
$\bullet r^2$  IS THE POLE RADII SQUARED



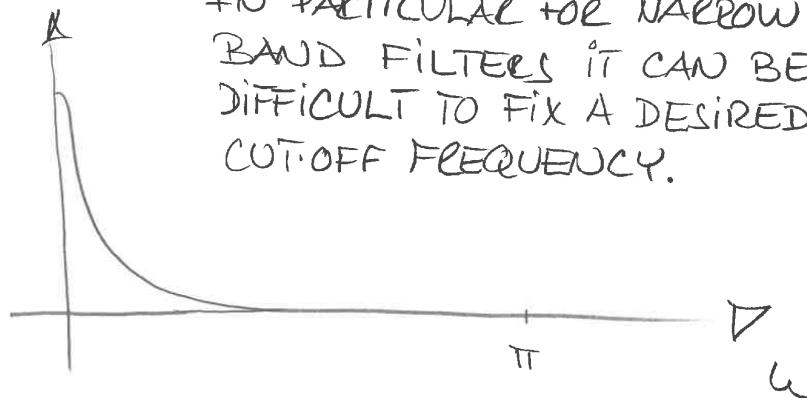
## DIRECT FORM STRUCTURE



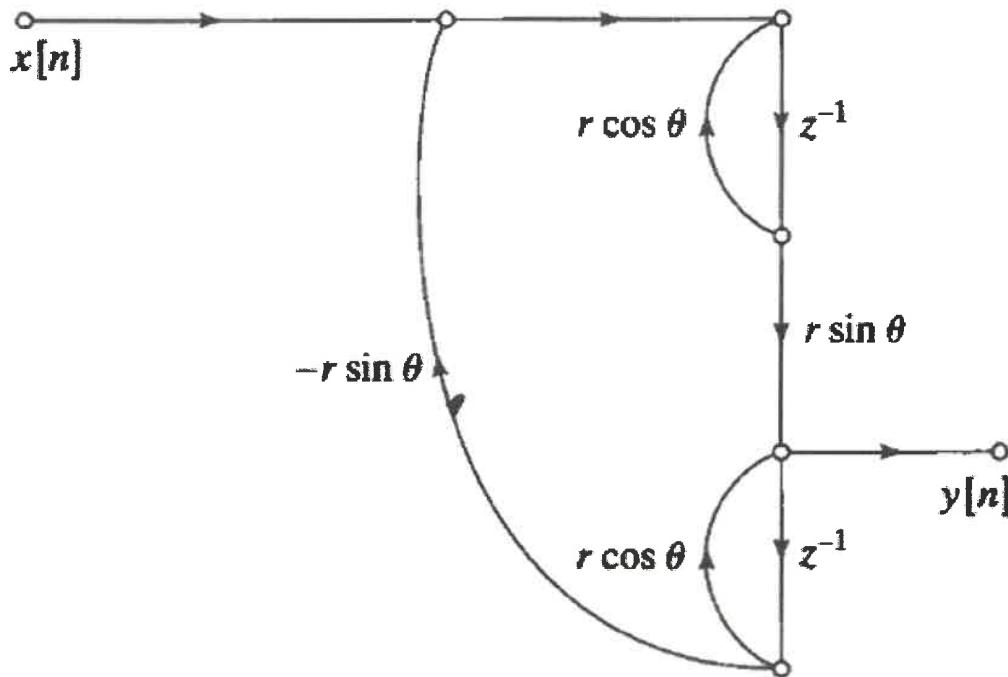
# Possible Pole Locations for 7 bit



IN PARTICULAR FOR NARROW  
BAND FILTERS IT CAN BE  
DIFFICULT TO FIX A DESIRED  
CUT-OFF FREQUENCY.



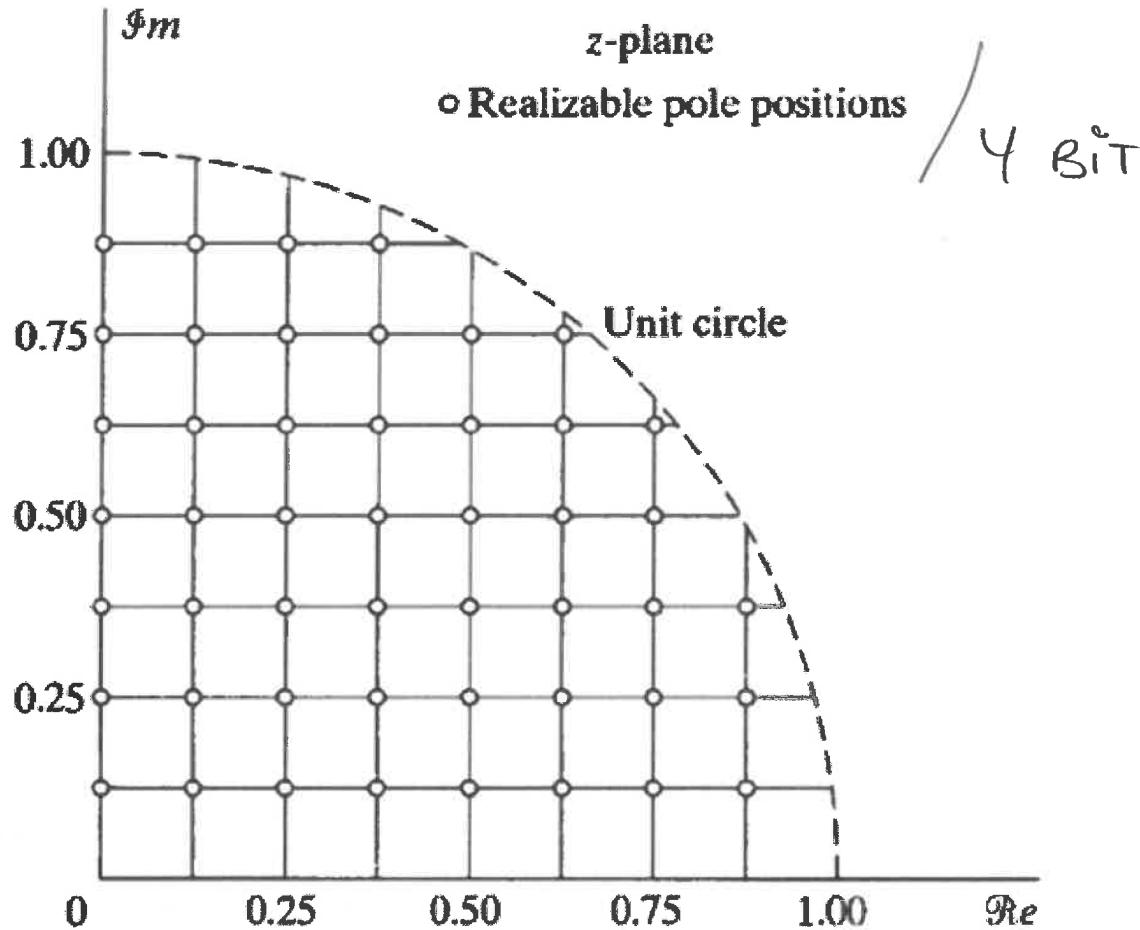
## Coupled Form



## Exercise

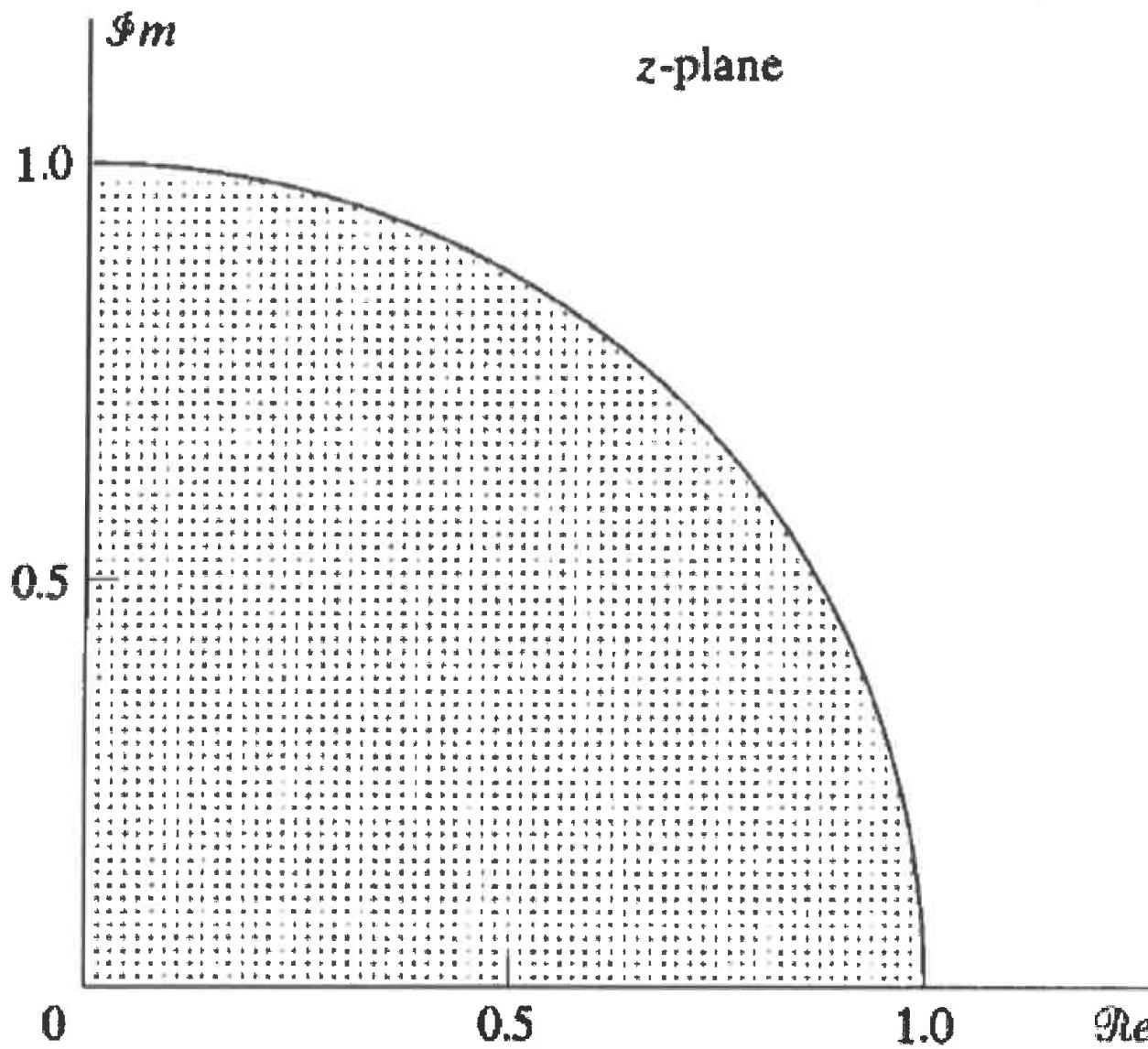
Show THAT THE POLES (FLOATING POINT)  
ARE IDENTICAL IN THE DIRECT FORM AND  
IN THE COUPLED FORM.

# COUPLED FORM



Now POLE LOCATIONS IN A  
REGULAR GRID WITHIN THE  
UNIT CIRCLE

COUPLED FORM / 7-BIT



# Scaling

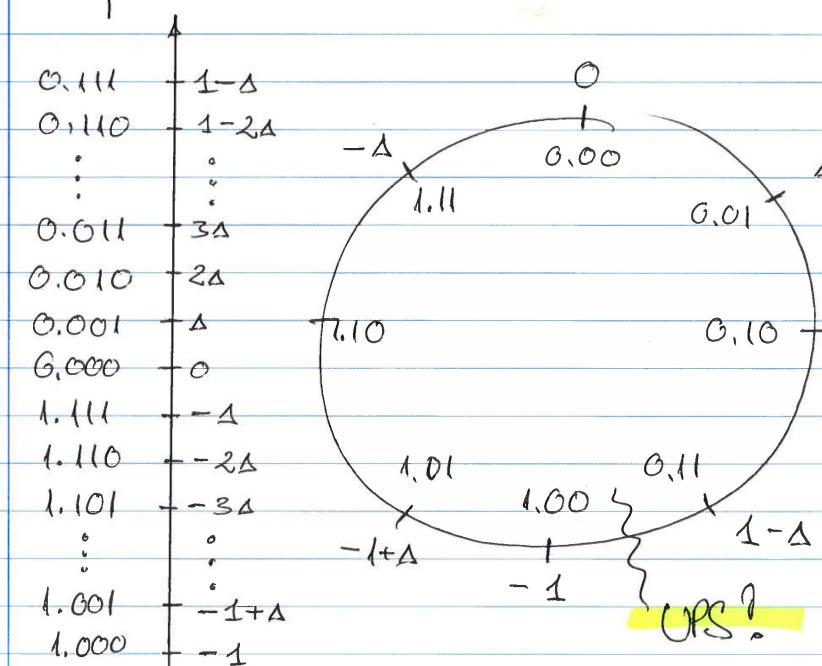


AALBORG UNIVERSITY  
DENMARK

# Scaling

NO MATTER WHETHER WE ARE USING DF-I  
OR DF-II, WE HAVE TO ENSURE THAT  
OUR IMPLEMENTATION OF THE FILTER PROVIDES  
THE LARGEST POSSIBLE SNR, WITHOUT OVERFLOW

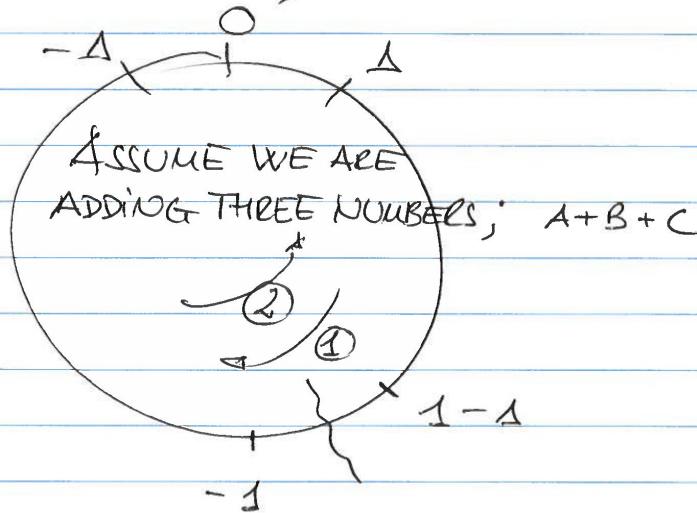
WHAT IS OVERFLOW?



$2^S$  COMPLEMENT.

OVERFLOW IS POSSIBLE ONLY WHEN ADDING  
TWO NUMBERS WITH THE SAME SIGN

THEREFORE, 2<sup>s</sup> COMPLEMENT IS VERY USEFUL  
COMBATING OVERFLOW ;



①  $|A+B| \geq 1 \Rightarrow$  OVERFLOW

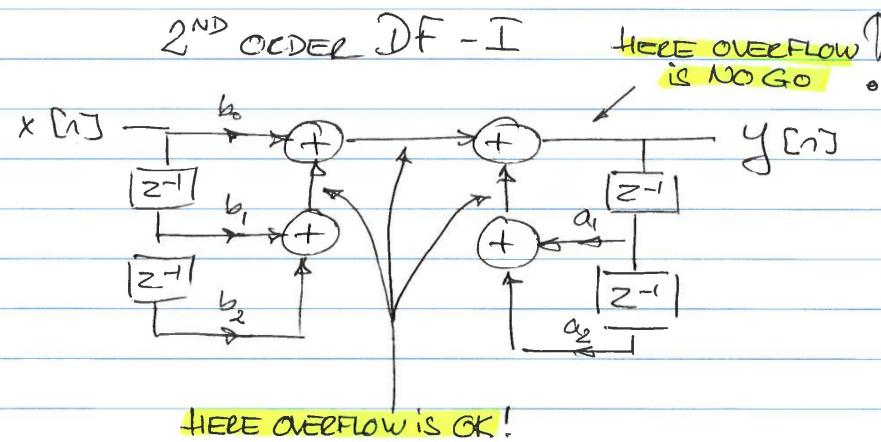
②  $| (A+B) + C | < 1 \Rightarrow$  NO OVERFLOW.

THUS ;

THERE MAY BE SOME INTERNAL VARIABLES  
IN OUR FILTER STRUCTURE WHERE THE  
PARTIAL SUMS ARE ALLOWED TO  
OVERFLOW, AS LONG AS WE ENSURE  
THAT THE TOTAL SUM DOES NOT OVERFLOW

## AN ALTERNATIVE FORMULATION;

WE WANT TO DESIGN OUR FILTER WITH BEST POSSIBLE UTILIZATION OF THE DYNAMIC RANGE (I.E., BEST SNR) OF THE VARIABLES, AND AT THE SAME TIME REDUCE THE RISK OF OVERFLOW IN THOSE VARIABLES WHERE OVERFLOW IS NOT ALLOWED



"... REDUCE THE RISK OF OVERFLOW..."

OUR PROBLEM IS THAT 0% RISK AND A HIGH SNR (FULL UTILIZATION OF THE DYNAMIC RANGE) ARE CONFLICTING REQUIREMENTS

RSS

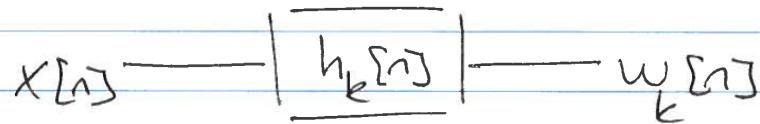
# THREE POSSIBLE SCALING STRATEGIES

PP. 463-465

## 1) MAX-VALUE SCALING

THE IDEA IS TO SCALE THE INPUT SIGNAL SUCH THAT OVERFLOW DOES NOT OCCUR IN THOSE VARIABLE WHERE OVERFLOW IS NOT ALLOWED.

- SUCH A VARIABLE IS DENOTED  $w_k$  AND THE SEQUENCE IN THIS VARIABLE IS  $w_k[n]$ ,



WHERE  $h_k[n]$  IS THE IMPULSE RESPONSE FROM THE  $k^{\text{th}}$  INPUT TO THE VARIABLE  $w_k$ .

$$|w_k[n]| = \left| \sum_{m=-\infty}^{\infty} h_k[m]x[n-m] \right|$$

NOW, ASSUME THAT WE KNOW THE NUMERICAL LARGEST VALUE IN  $x[n] - x_{\max}$

$$w_k[n] = \left| \sum_{m=-\infty}^{\infty} h_k[m] \cdot x[n-m] \right|$$

$$\left| w_k[n] \right| \leq x_{\max} \left| \sum_{m=-\infty}^{\infty} h_k[m] \right| \leq x_{\max} \cdot \sum_{m=-\infty}^{\infty} |h_k[m]|$$

$< 1$

REQUIREMENT  
(NO OVERFLOW)

$$x_{\max} < \frac{1}{\sum_{m=-\infty}^{\infty} |h_k[m]|}$$

WHICH SHOULD BE TRUE FOR ALL VARIABLES

$w_k$  WHERE OVERFLOW IS NOT ALLOWED

IF THIS IS NOT THE CASE, THEN WE CAN  
SCALE THE INPUT SIGNAL (MULTIPLY WITH  $s$ )

$$s \cdot x_{\max} < \frac{1}{\max_k \left\{ \sum_{m=-\infty}^{\infty} |h_k[m]| \right\}}$$

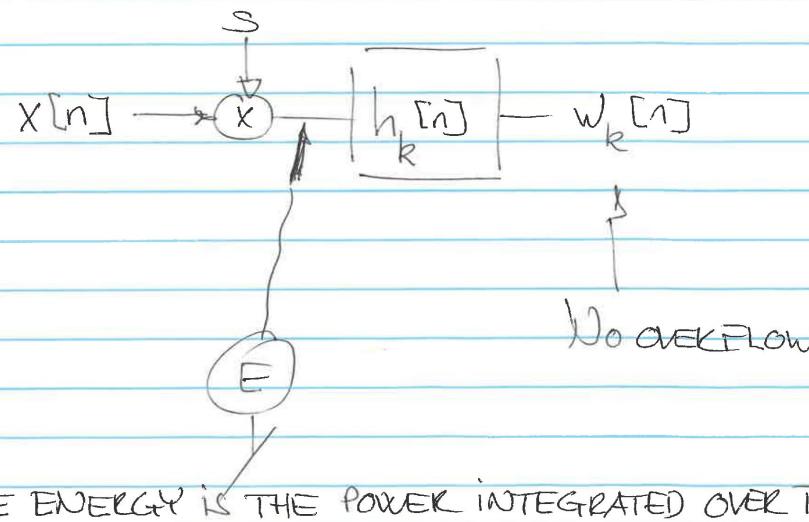
$$x[n] \rightarrow \underbrace{\otimes}_s \rightarrow \underbrace{[h_k[n]]}_{\text{NEVER OVERFLOW}} \rightarrow w_k[n]$$

AN ISSUE WITH THE MAX-VALUE SCALING STRATEGY IS THAT  $S$  MAY BECOME SO SMALL THAT IT COMPROMISES THE SNR OF THE FILTER.

2) SINUSOID SCALING (USED FOR NARROW BAND SIGNALS — SELF STUDY)

3) VARIANCE SCALING

THE IDEA HERE IS TO ADJUST THE ENERGY (VARIANCE) IN THE INPUT SIGNAL SUCH THAT THERE IS NO OVERFLOW IN THE VARIABLE WHERE OVERFLOW IS NOT ALLOWED



$$w_k[n] = h_k[n] * x[n]$$

$$\Downarrow \quad W(e^{j\omega}) = H_k(e^{j\omega}) \cdot X(e^{j\omega})$$

WE CAN NOW USE THE IDTFT ON THE  
RIGHT-HAND SIDE ;

$$w_k[n] = \stackrel{\circ}{\text{IDTFT}} \{W(e^{j\omega})\} = \stackrel{\circ}{\text{IDTFT}} \{H_k(e^{j\omega}) \cdot X(e^{j\omega})\}$$

$$\Downarrow \quad w_k[n] = \frac{1}{2\pi} \int_{-\pi}^{\pi} H_k(e^{j\omega}) \cdot X(e^{j\omega}) e^{j\omega n} d\omega$$

$$\Downarrow \quad |w_k[n]|^2 = \left| \frac{1}{2\pi} \int_{-\pi}^{\pi} H_k(e^{j\omega}) \cdot X(e^{j\omega}) e^{j\omega n} d\omega \right|^2$$

$$\Downarrow \quad |w_k[n]|^2 \leq \left( \frac{1}{2\pi} \int_{-\pi}^{\pi} |H_k(e^{j\omega}) \cdot X(e^{j\omega}) e^{j\omega n}|^2 d\omega \right)^2$$

SCHWARTZ

$$= \left( \frac{1}{2\pi} \int_{-\pi}^{\pi} |H_k(e^{j\omega})|^2 |X(e^{j\omega})|^2 d\omega \right)$$

$$\leq \left( \frac{1}{2\pi} \int_{-\pi}^{\pi} |H_k(e^{j\omega})|^2 d\omega \right) \cdot \left( \frac{1}{2\pi} \int_{-\pi}^{\pi} |X(e^{j\omega})|^2 d\omega \right)$$

(Equ. 125 p. 464 (without proof here))

$$(*) |W_k[n]|^2 \leq \left( \frac{1}{2\pi} \int_{-\pi}^{\pi} |H_k(e^{j\omega})|^2 d\omega \right) \cdot \left( \frac{1}{2\pi} \int_{-\pi}^{\pi} |X(e^{j\omega})|^2 d\omega \right)$$

Now, let's introduce PARSEVAL's THEOREM ;

$$E_x = \sum_{n=-\infty}^{\infty} |x[n]|^2 = \frac{1}{2\pi} \int_{-\pi}^{\pi} |X(e^{j\omega})|^2 d\omega$$

We now use PARSEVAL in (\*)

$$|W_k[n]|^2 \leq \sum_{n=-\infty}^{\infty} |x[n]|^2 \cdot \frac{1}{2\pi} \int_{-\pi}^{\pi} |H_k(e^{j\omega})|^2 d\omega$$

Our REQUIREMENT WAS ;

"ADJUST THE ENERGY IN THE INPUT-SIGNAL  
 $x[n]$  SUCH THAT THERE IS NO OVERFLOW  
 IN  $W_k[n]$ "

$$\Downarrow |W_k[n]|^2 < 1$$

$$\Downarrow \sum_{n=-\infty}^{\infty} |s \cdot x[n]|^2 \cdot \frac{1}{2\pi} \int_{-\pi}^{\pi} |H_k(e^{j\omega})|^2 d\omega < 1$$

$$S^2 \sum_{n=-\infty}^{\infty} |x[n]|^2 \cdot \frac{1}{2\pi} \int_{-\pi}^{\pi} |H_k(e^{j\omega})|^2 d\omega \leq 1$$

$$S^2 \cdot E_x < \frac{1}{2\pi} \int_{-\pi}^{\pi} |H_k(e^{j\omega})|^2 d\omega = \sum_{n=-\infty}^{\infty} |h_k[n]|^2$$

which should hold for all variables where overflow is not allowed.

If we assume that the energy in the input signal  $E_x = 1$  then

$$S < \sqrt{\sum_{n=-\infty}^{\infty} |h_k[n]|^2}$$



$$S = \max_k \left\{ \sqrt{\sum_{n=-\infty}^{\infty} |h_k[n]|^2} \right\}$$