

**ECE4270**  
**Fundamentals of DSP**  
**Lecture 19**

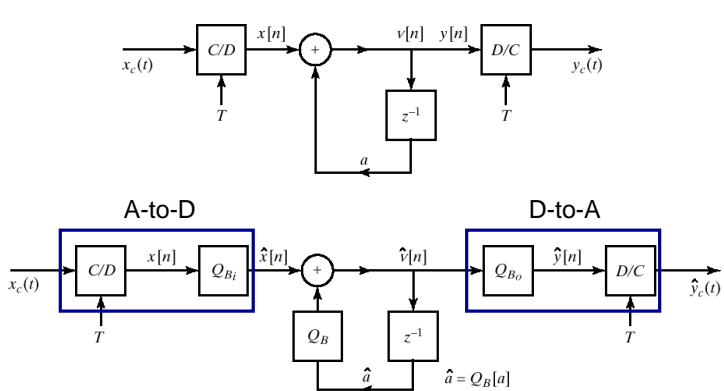
**Fixed-Point Numbers and Arithmetic**

School of Electrical and Computer Engineering  
Center for Signal and Information Processing  
Georgia Institute of Technology

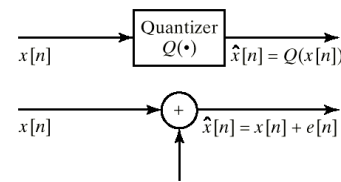
**Overview of Lecture**

- Quantization in LTI Implementation
- Two's-complement arithmetic
  - Integers and fractions
  - Scaling for fixed-point arithmetic
  - Quantizing filter coefficients
  - Addition & Multiplications

**Quantization in LTI Implementation - I**



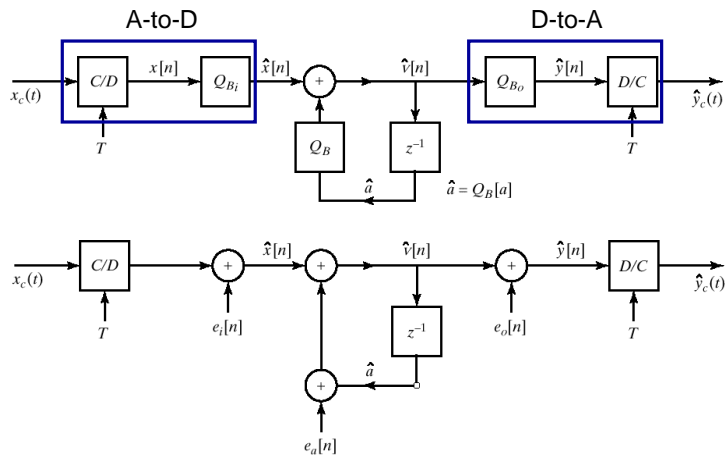
**Linear Noise Model**



- Error is uncorrelated with the input.
- Error is uniformly distributed over the interval  $-(\Delta / 2) < e[n] \leq (\Delta / 2)$ .
- Error is stationary white noise, (i.e. flat spectrum)

$$P_e(\omega) = \sigma_e^2 = \frac{\Delta^2}{12}, \quad |\omega| \leq \pi$$

## Quantization in LTI Implementation - II



ECE4270

Spring 2017

## FIR Filter Coefficients

Coefficient	Unquantized
$h[0] = h[27]$	$1.359657 \times 10^{-3}$
$h[1] = h[26]$	$-1.616993 \times 10^{-3}$
$h[2] = h[25]$	$-7.738032 \times 10^{-3}$
$h[3] = h[24]$	$-2.686841 \times 10^{-3}$
$h[4] = h[23]$	$1.255246 \times 10^{-2}$
$h[5] = h[22]$	$6.591530 \times 10^{-3}$
$h[6] = h[21]$	$-2.217952 \times 10^{-2}$
$h[7] = h[20]$	$-1.524663 \times 10^{-2}$
$h[8] = h[19]$	$3.720668 \times 10^{-2}$
$h[9] = h[18]$	$3.233332 \times 10^{-2}$
$h[10] = h[17]$	$-6.537057 \times 10^{-2}$
$h[11] = h[16]$	$-7.528754 \times 10^{-2}$
$h[12] = h[15]$	$1.560970 \times 10^{-1}$
$h[13] = h[14]$	$4.394094 \times 10^{-1}$

The condition for linear phase is

$$h[M-n] = \pm h[n]$$

In this case,

$$\max \{h[k]\} \leq 0.5$$

ECE4270

Spring 2017

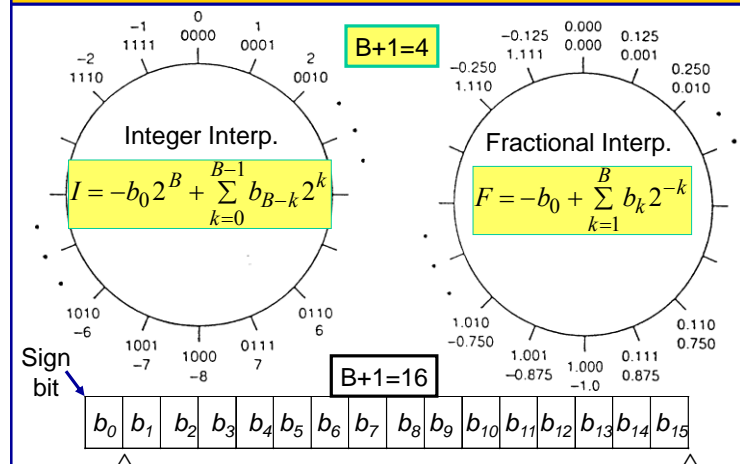
## Fixed-Point Arithmetic in DSP Chips

- Numbers in fixed-point DSPs are represented as **two's complement** numbers. (16 bits in TI chips)
- Although the processor deals implicitly with signed two's complement integers, filter coefficients have fractional parts. Representation of these fractions must be built into the program.
- Proper scaling of the signals and coefficients is required to maintain precision while avoiding overflow. Therefore, the programmer must constantly worry about the following: **scaling**, **quantization** (roundoff noise), and **overflow**.

ECE4270

Spring 2017

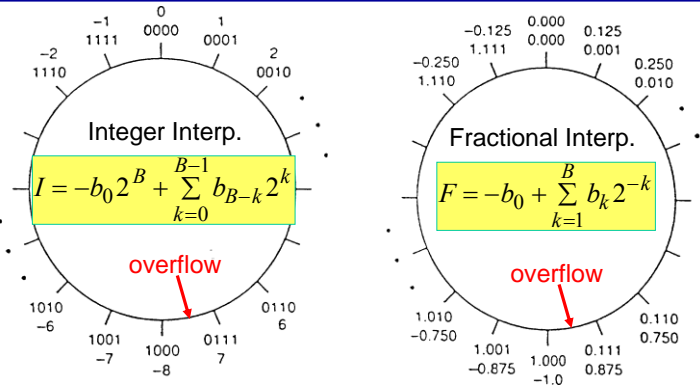
## Two's Complement Numbers



ECE4270

Spring 2017

## Two's Complement Addition



Adding moves clockwise and subtracting moves counter clockwise around the circle. Overflow wraps around.

ECE4270

Spring 2017

## Fixed-Point Scaling - I

- 16-bit two's complement integers range in size from  $-32768_{10}$  to  $+32767_{10}$ .
- Q notation is a convenient way for a programmer to keep track of the binary point when representing fractional numbers by integers. Consider a fractional number  $a$  such that  $-1 \leq a < 1$ , then we can represent  $a$  by a Q15 integer  $A$  such that  $-32768 \leq A \leq 32767$ . The relationship between  $a$  and  $A$  is simply

$$a = A \times 2^{-15} \quad \text{or} \quad A = a \times 2^{15}$$

ECE4270

Spring 2017

## Fixed-Point Scaling - II

- Consider a fractional number  $a$  such that  $-1 \leq a < 1$ , then we can represent  $a$  by a Q15 integer  $A$  such that  $-32768 \leq A \leq 32767$ . The relationship between  $a$  and  $A$  is

$$a = A \times 2^{-15} \quad \text{or} \quad A = a \times 2^{15}$$

- For example

$$a = 0.75 \Leftrightarrow A = 24576_{10} \text{ Q15}$$

0  $\wedge$  110000000000000  
15 bits

Q15 means 15 bits to the right of the binary point.

ECE4270

Spring 2017

## Fixed-Point Scaling - III

- Consider a mixed number  $a$  such that  $-4 \leq a < 4$ . Then we can represent  $a$  by a Q13 integer  $A$  such that  $-32768 \leq A \leq 32767$ . The relationship between  $a$  and  $A$  is

$$a = A \times 2^{-13} \quad \text{or} \quad A = a \times 2^{13}$$

- For example

$$a = 3.5 \Leftrightarrow A = 28672_{10} \text{ Q13}$$

011  $\wedge$  1000000000000  
13 bits

ECE4270

Spring 2017

## Fixed-Point Scaling - IV

- The smallest number that can be represented by a Q15 number is

$$\Delta = \frac{\text{range}}{\text{number of possibilities}} = \frac{2}{2^{16}} = \frac{1}{2^{15}}$$

- In general, the smallest number representable as a QB number will be

$$\Delta = \frac{1}{2^B}$$

- That is, in a QB number, the least significant bit (LSB) has value

$$\Delta = \frac{1}{2^B}$$

ECE4270

Spring 2017

## Example of Quantizing Coefficients

```
% Multiply the coefficient by 2^15
»a=-.001359657; A=a*2^(15)
A =
-44.55324057600000
% Round (or truncate) the result
»Ahat=round(A)
Ahat =
-45
% The equivalent quantized fraction is therefore
»ahat=Ahat/2^(15)
ahat =
-0.00137329101562
```

ECE4270

Spring 2017

## Quantized FIR Filter Coefficients

Coefficient	Unquantized	Q15	Q12
$h[0] = h[27]$	$1.359657 \times 10^{-3}$	$45 \times 2^{-15}$	$6 \times 2^{-12}$
$h[1] = h[26]$	$-1.616993 \times 10^{-3}$	$-53 \times 2^{-15}$	$-7 \times 2^{-12}$
$h[2] = h[25]$	$-7.738032 \times 10^{-3}$	$-254 \times 2^{-15}$	$-32 \times 2^{-12}$
$h[3] = h[24]$	$-2.686841 \times 10^{-3}$	$-88 \times 2^{-15}$	$-11 \times 2^{-12}$
$h[4] = h[23]$	$1.255246 \times 10^{-2}$	$411 \times 2^{-15}$	$51 \times 2^{-12}$
$h[5] = h[22]$	$6.591530 \times 10^{-3}$	$216 \times 2^{-15}$	$27 \times 2^{-12}$
$h[6] = h[21]$	$-2.217952 \times 10^{-2}$	$-727 \times 2^{-15}$	$-91 \times 2^{-12}$
$h[7] = h[20]$	$-1.524663 \times 10^{-2}$	$-500 \times 2^{-15}$	$-62 \times 2^{-12}$
$h[8] = h[19]$	$3.720668 \times 10^{-2}$	$1219 \times 2^{-15}$	$152 \times 2^{-12}$
$h[9] = h[18]$	$3.233332 \times 10^{-2}$	$1059 \times 2^{-15}$	$132 \times 2^{-12}$
$h[10] = h[17]$	$-6.537057 \times 10^{-2}$	$-2142 \times 2^{-15}$	$-268 \times 2^{-12}$
$h[11] = h[16]$	$-7.528754 \times 10^{-2}$	$-2467 \times 2^{-15}$	$-308 \times 2^{-12}$
$h[12] = h[15]$	$1.560970 \times 10^{-1}$	$5115 \times 2^{-15}$	$639 \times 2^{-12}$
$h[13] = h[14]$	$4.394094 \times 10^{-1}$	$14399 \times 2^{-15}$	$1800 \times 2^{-12}$

ECE4270

Spring 2017

## Quantized Filter Coefficients

- For fixed-point implementation, the filter coefficients generally will be computed by a design algorithm that gives the filter coefficients as floating-point numbers. Therefore, they must be quantized to  $B+1$  bits

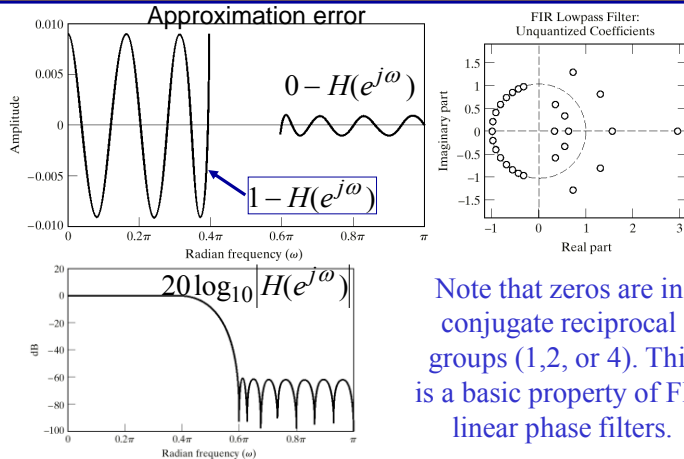
$$\hat{h}[n] = Q_B \{h[n]\} = h[n] + \Delta h[n]$$

$$\begin{aligned} \hat{H}(e^{j\omega}) &= \sum_{n=0}^M (h[n] + \Delta h[n]) e^{-j\omega n} \\ &= H(e^{j\omega}) + \sum_{n=0}^M \Delta h[n] e^{-j\omega n} \end{aligned}$$

ECE4270

Spring 2017

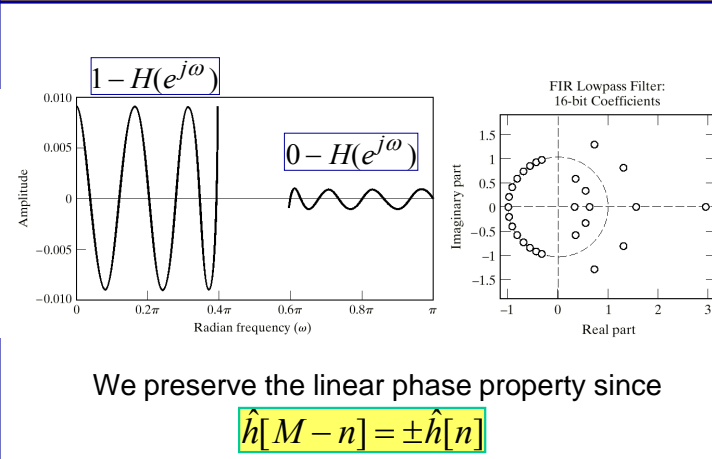
## Unquantized FIR Filter Response (II)



ECE4270

Spring 2017

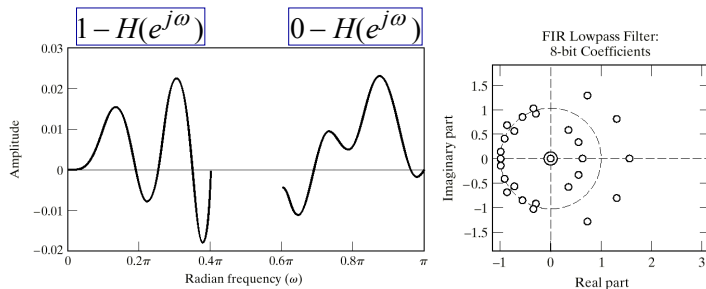
## 16-Bit Quantization of FIR Filter (III)



ECE4270

Spring 2017

## 8-Bit Quantization of FIR Filter (IV)



Zeros stay in reciprocal quads, but shift significantly due to the quantization.

ECE4270

Spring 2017

## 2's Complement Numbers

- To change the sign of a 2's complement number, just complement all the bits, and add 1 to the least significant bit.

$$-6 = -(0110) = 1001 + 1 = 1010$$

- When accumulating 3 or more 2's complement numbers, the intermediate sums can overflow, but the final sum will be correct if it does not exceed the word length of the numbers.

$$\underline{6} + \underline{4} + (-6) = 10 + (-6) = 4$$

$$\underline{0110} + \underline{0100} + 1010 = 1010 + 1010 = 0100$$

Final Sum correct

Partial sum overflows

ECE4270

Spring 2017



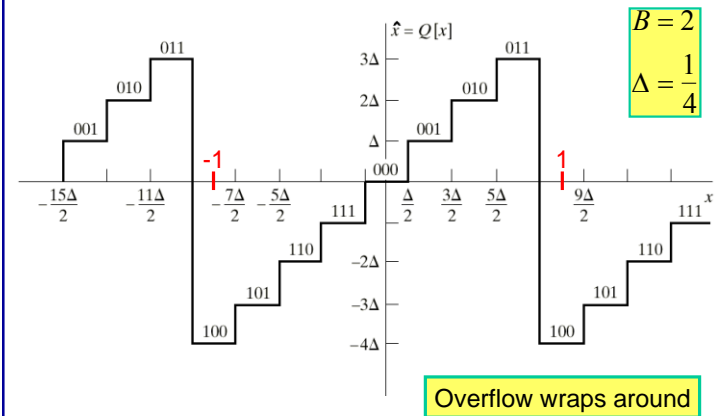
## Overflow

- Overflow can occur in two ways:
  - Accumulator overflows due to additions
  - Taking 16 bits out of the wrong part of a product
- Overflow can be prevented by:
  - More precision - use larger accumulator
  - Saturation arithmetic - clip accumulator at largest value
  - Shift products to the right to discard low-order bits. This results in loss of precision.

ECE4270

Spring 2017

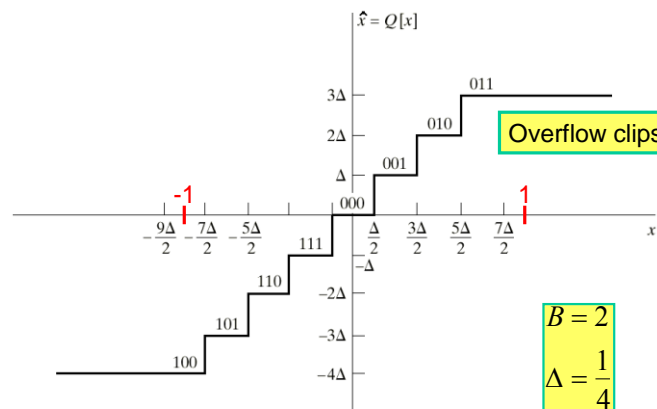
## Two's Complement Quantizer



ECE4270

Spring 2017

## Two's Complement Saturation



ECE4270

Spring 2017