

CSE4210  
Architecture and Hardware  
for DSP  
Lecture 1  
Introduction &  
Number systems

Administrative Stuff

- CSE4210 Architecture and Hardware for DSP
- Text: VLSI Digital Signal Processing Systems: Design and Implementation.  
K. Parhi. Wiley Interscience
- Posted articles

## Administrative Stuff

- Office hours: Monday 1-2pm TR 3-4pm
- Room 2026 CSEB x40607
- HW 0%
- Quizes 15%
- Midterm 25%
- Project 15%
- Final 45%

## Topics

- Number systems
- Fast arithmetic
- Algorithm representation
- Transformation (retiming, unfolding, folding)
- Systolic arrays and mapping algorithms into hardware
- Low power design

## Introduction

- Introduction to DSP algorithms
- Non-terminating programs in real time.
- Speed depends on applications (audio, video, 2-D, 3-D, ...)
- Need to design families of architectures for specified algorithm complexity and speed constraints

## Typical DSP Programs



3-Dimensional optimization: Area, Speed, Power)

Usually, speed is a requirement, area-power tradeoff

$$P=C V^2 F$$

## Examples

- FIR filter,  $x(n)$  is the input,  $y(n)$  output

$$y(n) = \sum_{j=0}^{J-1} h(j)x(n-j)$$

- IIR filter

$$y(n) = \sum_{i=1}^P a(i)y(n-i) + \sum_{k=0}^Q b(k)x(n-k)$$

## Examples

- Convolution  $y(n) = \sum_{i=0}^{M-1} x(i)h(n-i) = \sum_{j=0}^{N-1} h(j)x(n-j)$

```

For n=1 to M+N-2,
    y(n)=0,
    For i=0:M-1,
        y(n)=y(n)+x(i)*h(n-i)
        end
    end

```

MAC operation

## More Complex Examples Motion Estimation

- Image (frame) is divided into macroblocks
- Each macroblock is compared to a macroblock in the reference frame using some error measure.
- The search is conducted over a predetermined search area.
- A vector denoting the displacement of the (motion) is sent.

## More Complex Examples Motion Estimation

- Many measures of errors could be used.
- The displaced block difference  $s(m,n)$  using MAD (Mean Absolute Difference) is defined as

$$s(m,n) = \sum_{i=0}^{N-1} \sum_{j=0}^{N-1} |x(i,j) - y(i+m, j+n)|$$

$m, n$  are in the search area,  $N$  is the macroblock size.  
The one with the minimum error is chosen

## More Complex Examples

### Vector Quantization

- Used in compression
- A group of samples (vector) are quantized together
- For example consider  $k$  pixels, with  $W$  bits.
- That vector is compared to a group of  $N$  *codewords*, choose the one with the min. distortion.
- We transmit the index of that codeword

## More Complex Examples

### Vector Quantization

- Compression ratio =  $KW/\log_2 N$
- Euclidean distance is used as a measure of distortion.

$$\begin{aligned}
 d(\mathbf{x}, \mathbf{c}_j) &= \|\mathbf{x} - \mathbf{c}_j\|^2 = \sum_{i=0}^{k-1} (x_i - c_{ji})^2 \\
 &= \|\mathbf{x}\|^2 - 2(\mathbf{x} \cdot \mathbf{c}_j + e_j), \quad e_j = -\frac{1}{2} \|\mathbf{c}_j\|^2 = -\frac{1}{2} \sum_{i=0}^{k-1} c_{ji}^2
 \end{aligned}$$

## Discrete Cosine Transform

- The 1-D DCT is defined as

$$X(K) = e(k) \sum_{n=0}^{N-1} x(n) \cos\left[\frac{(2n+1)k\pi}{2N}\right] , \quad k = 0, 1, 2, \dots, N-1$$

$$e(k) = \begin{cases} \frac{1}{\sqrt{2}} & \text{if } k = 0 \\ 1 & \text{otherwise} \end{cases}$$

## More Complex Examples

- Viterbi Decoding
- FFT
- Wavelets and Filter banks
- See the book for details

## Requirements

- Consider block matching algorithm, the computational requirement is as follows
- $3(2p+1)^2NMF$
- $3*(2*7+1)*288*352*30=2GOP$
- Much higher for higher resolution and bigger frames
- How to achieve these requirements?

## hardware

- Microprocessors
- Microprocessors with DSP extension
- DSP
- FPGA
- ASIC

## Number System

- Numbers and their representation
- Binary numbers
- Negative numbers
- Unconventional numbers

## Binary Numbers

- An ordered sequence  
 $(x_{n-1}, x_{n-2}, \dots, x_1, x_0) \quad x_n \subset \{0,1\}$
- The value of the number is
- $$X = x_{n-1}2^{n-1} + x_{n-2}2^{n-2} + \dots + x_12 + x_0 = \sum_{i=0}^{n-1} x_i 2^i$$
- The range  $[X_{\min}, X_{\max}]$  is the range of the numbers to be represented, in the previous case  $[0, 2^n - 1]$

## Binary numbers

- The previous representation is non-redundant and weighted ( $w_i$ )
- The n-digit number can be partitioned into a fraction part (n-k bits) and an integral part(k bits)

$$\underbrace{(x_{k-1}x_{k-2}\cdots x_1x_0)}_{\text{integral part}} \cdot \underbrace{(x_{-1}x_{-2}\cdots x_{-m})_r}_{\text{fractional part}}$$

$$X = x_{k-1}r^{k-1} + x_{k-2}r^{k-2} + \cdots + x_1r + x_0 + x_{-1}r^{-1} + \cdots + x_{-m}r^{-m}$$

## Binary numbers

- Given the length of the operand, n, the weight  $r^{-m}$  of the least significant digit indicates the position of the radix point.
- Unit in Last Position ulp= $r^{-m}$
- Simplifies the discussion and there is no need to partition the number into fractional and integral parts.

## Conversion

- Convert 36.4375 into binary

36/2	Division by 2		Multiplication by 2	
	Quotient	Remainder	Integer	Fraction
18	0		0	.875
9	0		1	.75
4	1		1	.5
2	0		1	0.5
1	0		1	0.0
0	1			

100100.0111

## Negative Numbers

- Signed magnitude
- Complement
  - Diminished radix complement (1's complement for binary)
  - Radix complement (2's complement in binary)

## Signed magnitude

- The  $n^{\text{th}}$  bit (digit) is the sign
- $n-1$  digits for magnitude ( $k-1$  integral and  $m$  fractional).
- Largest value  $011\dots11 = X_{\max} = r^{k-1}-\text{ulp}$
- Smallest negative value  $-(r^{k-1}-\text{ulp}) 11..1$
- Two representation for zero

## Signed magnitude

- Operations may be more complicated than using complement.
- For example, adding 2 numbers, a positive number  $X$ , and a -ve number  $Y$ , the result depends on if  $X > Y$  or not
  - If  $X > Y$  the result is  $X+(-Y)$
  - If  $Y > X$  switch the 2 numbers, subtract, attach minus sign  $-(Y-X)$

## Complement representation

- Positive numbers are represented just as signed-magnitude
- Negative numbers are represented as R-number, where R is a constant
- Note that  $-(-y) = R - (R - Y) = Y$
- The choice of R must satisfy 2 conditions
  - Calculating the complement is easy
  - Simplifying or eliminating correction

## Complement representation

- For radix complement  $R = r^k$
- For the diminished radix complement  
 $R = r^k - ulp$

$$r=2, k=n=4, m=0, \text{ulp}=2^0 = 1$$

Sequence	Two's complement	One's complement	Signed-magnitude
0111	7	7	7
0110	6	6	6
0101	5	5	5
0100	4	4	4
0011	3	3	3
0010	2	2	2
0001	1	1	1
0000	0	0	0
1111	-1	-0	-7
1110	-2	-1	-6
1101	-3	-2	-5
1100	-4	-3	-4
1011	-5	-4	-3
1010	-6	-5	-2
1001	-7	-6	-1
1000	-8	-7	-0

## Example

2's complement

$$6 = 0110 \quad -6 = 1010 \quad 4 = 0100 \quad -4 = 1100$$

$$\begin{array}{cccc} 6-4 & 4-6 & 6+4 & -4-6 \\ 0110 & 0100 & 0110 & 1100 \\ 1100 & 1010 & 0100 & 1010 \\ \hline 1\ 0010 & 1110 & 1010 & 1\ 0110 \end{array}$$

+2                  -2  
 Carry in = carry out  
 Ignore carry out

-6                  +6  
 Carry in ≠ carry out  
 Overflow

## Example

1's complement

$$6 = 0110 \quad -6 = 1001 \quad 4 = 0100 \quad -4 = 1011$$

$$\begin{array}{r} 6-4 \\ 4-6 \\ \hline \end{array}$$

$$\begin{array}{r} 0110 \\ 0100 \\ \hline 1011 \end{array}$$

$$\begin{array}{r} 1011 \\ 1001 \\ \hline 0100 \end{array}$$

$$\begin{array}{r} 10001 \\ 1101 \\ \hline 1010 \end{array}$$

$$+2$$

$$-2$$

$$-5$$

$$+4$$

$\underbrace{\hspace{2cm}}$  Carry in = carry out

Add to LSB

$\underbrace{\hspace{2cm}}$  Carry in  $\neq$  carry out

Overflow

## Arithmetic Shift

- Consider the number  $\{x_{n-2}, x_{n-3}, \dots, x_0\}$
- Finite extension of signed magnitude  
is  $\dots 0, 0 \{x_{n-2}, x_{n-3}, \dots, x_0\} 0, 0, \dots$
- 2's complement  $\dots x_{n-1}, x_{n-1} \{x_{n-2}, x_{n-3}, \dots, x_0\} 0, 0, \dots$
- 1's complement  $\dots x_{n-1}, x_{n-1} \{x_{n-2}, x_{n-3}, \dots, x_0\} x_{n-1}, x_{n-1}$

## Arithmetic Shift

2's Comp

-6	1010
-12	10100
-24	101000
-6	1010
-3	1101
-2	1110
-1	1111

1's Comp

-6	1001
-12	10011
-24	100111
-6	1000
-3	1100
-1	1110
-0	1111

## Unconventional Number System

- Negative radix number system
- A general class of fixed-radix number system
- Signed-digit number system
- Residue number system

## Negative radix Number System

- The radix could be negative,  $r=-\beta$ ,  $\beta$  is a positive number.
- Digit set  $0, 1, \dots, \beta-1$
- Value of  $(x_{n-1}, x_{n-2}, \dots, x_0)$
- $$X = \sum_{i=0}^{n-1} x_i (-\beta)^i$$

$387_{-10} = 300 - 80 + 7 = 227$   
 $\text{Range} = [090, 909]_{-10}$ , or  $[-90, 909]_{10}$

$1010_{-2} = -8 - 2 = -10$   
 $\text{Range} = [1010, 0101]_{-2} = [-10, 5]_{10}$

## Negative radix Number System

- Algorithms do exist for basic operations.
- Not better than 2's complement systems
- Carry could be -ve or positive

## General Class of Fixed Radix Number System

- Characterized by  $(n, \beta, \Lambda)$ ,  $\beta$  is a positive radix, digit set  $0, 1, \dots, \beta-1$ , and a vector of length  $n$   $\Lambda = (\lambda_{n-1}, \lambda_{n-2}, \dots, \lambda_0)$   
 $\lambda_i = \{-1, 1\}$

$$X = \sum_{i=0}^{n-1} \lambda_i x_i \beta^i$$

- 2's Complement  $\Lambda = \{-1, 1, 1, \dots, 1\}$

## General Class of Fixed Radix Number System

$$P = \{p_{n-1}, p_{n-2}, \dots, p_0\} \quad \text{Max. positive number}$$

$$p_i = \begin{cases} \beta - 1 & \text{if } \lambda_i = +1 \\ 0 & \text{otherwise} \end{cases} = \frac{1}{2}(\lambda_i + 1)(\beta - 1)$$

$$\begin{aligned} P &= \sum_{i=0}^{n-1} 1/2(\lambda_i + 1)(\beta - 1)\beta^i = 1/2 \left[ \sum_{i=0}^{n-1} \lambda_i (\beta - 1)\beta^i + \sum_{i=0}^{n-1} (\beta - 1)\beta^i \right] \\ &= 1/2 [Q + (\beta^n - 1)] \end{aligned}$$

Where Q is the value of the tuple  $(\beta - 1, \beta - 1, \dots, \beta - 1)$

## Signed Digit Number System

- The digits could be positive or negative
- Redundant (more than one representation for the same number)
- For a radix  $\beta$ ,  $x_i \in \{\overline{\beta-1}, \overline{\beta-2}, \dots, \overline{1}, \overline{0}, \overline{1}, \dots, \overline{\beta-1}\}$
- To reduce redundancy,

$$X_i = \left\{ \overline{a}, \overline{a-1}, \dots, \overline{1}, \overline{0}, \overline{1}, \dots, \overline{a} \right\} \text{ where } \left\lceil \frac{r-1}{2} \right\rceil \leq a \leq r-1$$

## Signed Digit Number System

- Example:  $\beta=10, a=6$ ,

## Addition

- Add two numbers X,Y

$$w_i = x_i + y_i - rt_{i+1}$$

$$t_{i+1} = \begin{cases} 1 & \text{if } (x_i + y_i) \geq a \\ \bar{1} & \text{if } (x_i + y_i) \leq \bar{a} \\ 0 & \text{if } |x_i + y_i| < a \end{cases}$$

$$s_i = w_i + t_i \quad , t_0 = 0$$

Example, a=6,r=10

1634+3366

## Adition

## Converting to SD

- We can use the previous addition algorithm to do conversion
- EX:

## Breaking the carry chain

- For no carry  $|s_i| = |w_i + t_i| \leq a \Rightarrow |w_i| \leq a-1$

Case 1  $x_i + y_i = 2a$  (upper bound)

$$w_i = 2a - r(1) = 2a - r \Rightarrow a \leq r - 1$$

Case 2  $x_i + y_i = a$  (lower bound)

$$w_i = (a) - r(1) = a - r$$

$$|w_i| = r - a$$

$$r - a \leq a - 1 \Rightarrow \left\lceil \frac{r+1}{2} \right\rceil \leq a$$

$$\left\lceil \frac{r+1}{2} \right\rceil \leq a \leq r - 1$$

## Breaking the carry chain

- For binary SD,  $r=2, a=1$  it is impossible for the previous condition to hold.
- We do have carry in SD addition

## Breaking the carry chain

$X_i, Y_i$	00, <u>11</u>	01	01	<u>01</u>	<u>01</u>	11	<u>11</u>
$X_{i-1}, Y_{i-1}$		Neither one is <u>1</u>	At least one is <u>1</u>	Neither one is <u>1</u>	At least one is <u>1</u>		
$t_{i+1}$	0	1	0	0	<u>1</u>	1	<u>1</u>
$w_i$	0	<u>1</u>	1	<u>1</u>	0	0	0

## Breaking the carry chain

- Example

## Breaking the carry chain

- Example

## Residue Number System

- The number systems mentioned so far are positional.
- The residue number system uses positional bases that are relatively prime to each other.
- For example, the residues of (2,3,5) can encode 30 numbers as follows

## Residue Number System

- Converting long (carry) arithmetic into short arithmetic
- Addition, subtraction and multiplication is very fast (no carry)
- Division and comparison is slow

## Chinese Remainder Theorem

- Given a set of relatively prime moduli  $(m_1, m_2, \dots, m_n)$  then for any  $X < M$ , the set of residue  $\{X \bmod m_i \mid 1 \leq i \leq n\}$  is unique, where  $M = \prod_{i=1}^n m_i$
- $m_i = 4$  and  $m_j = 9$ ,  $\gcd(4, 9) = 1$
- Although Neither 4 Nor 9 is Prime, They are Relatively Prime

## Residue Number System

RNS

RNS(8,7,5,3)

## RNS Circuit Structure

## Residue Number System

	Residue to base					Residue to base					Residue to base			
N	5	3	2		N	5	3	2		N	5	3	2	
0	0	0	0		10	0	1	0		20	0	2	0	
1	1	1	1		11	1	2	1		21	1	0	1	
2	2	2	0		12	2	0	0		22	2	1	0	
3	3	0	1		13	3	1	1		23	3	2	1	
4	4	1	0		14	4	2	0		24	4	0	0	
5	0	2	1		15	0	0	1		25	0	1	1	
6	1	0	0		16	1	1	0		26	1	2	0	
7	2	1	1		17	2	2	1		27	2	0	1	
8	3	2	0		18	3	0	0		28	3	1	0	
9	4	0	1		19	4	1	1		29	4	2	1	

## Operations in RNS

- Addition and multiplication are done on per position modulo that position with no carries

$$\begin{array}{rcl} 9 & \rightarrow & [4, 1, 0] \\ +16 & \rightarrow & [1, 1, 0] \\ \hline & & [0, 1, 1] \end{array} \quad \begin{array}{rcl} 7 & \rightarrow & [2, 1, 1] \\ \times 4 & \rightarrow & [4, 1, 0] \\ \hline & & [3, 1, 0] \end{array}$$

## Selection of the Moduli

- Two criteria
  - Efficient in their binary representation
  - Straightforward computational op's
- Moduli of the form  $2^{k_1}, 2^{k_2} - 1, 2^{k_3} - 1, \dots, 2^{k_n} - 1$
- A sufficient condition for  $2^{a-1}$  and  $2^{b-1}$  to be a relatively prime pair is that  $a$  and  $b$  are relatively prime

## Addition using ROM