

CS2318

Real Number Representing

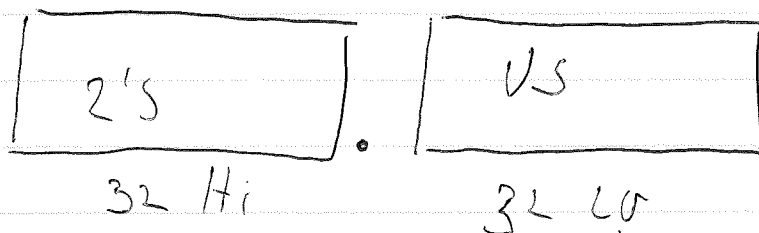
11/10/16

Imaginary Fixed point

Inside 1 to the left of a token

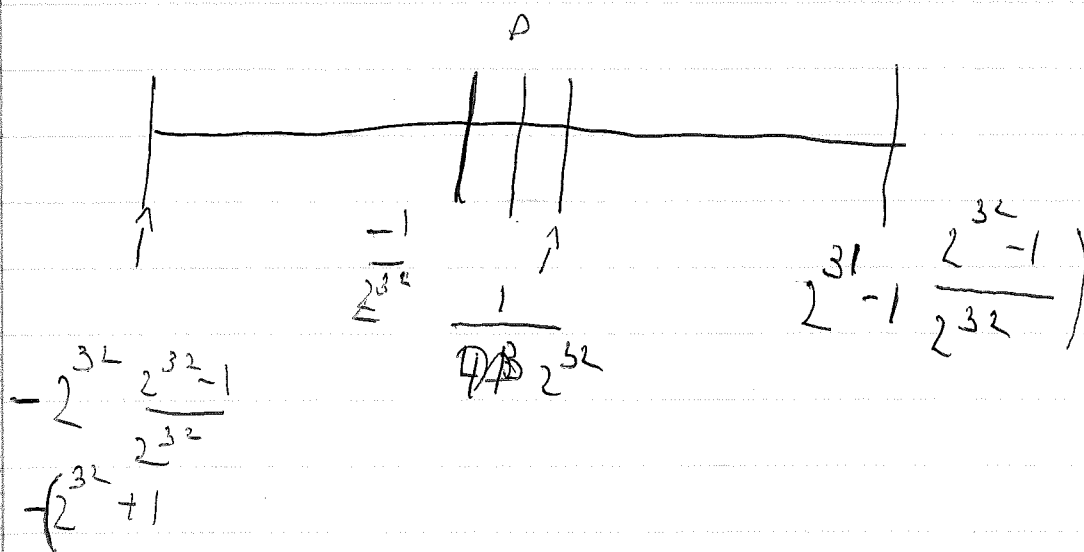
Division in MIPS

HI 2's
LO unsigned

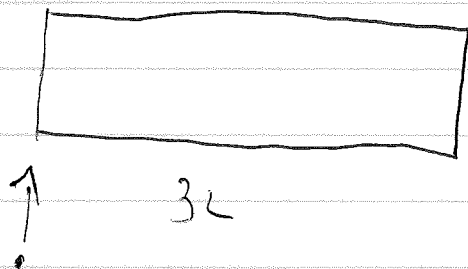


Int

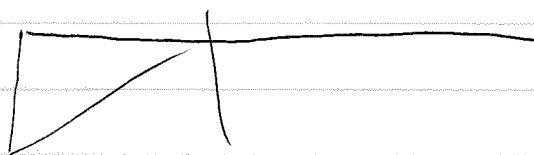
Remainder / mod



2



$$[-2^{31}, 2^{31}-1]$$



$$\left[-\frac{2^{31}}{2^{31}}, \frac{2^{31}-1}{2^{31}} \right] \quad (-1, 1)$$

~~As~~ -1×-1 does not work

with the Int unit

$\times, \div, +, -$ of of these Fixed

Points is done with the Integer

Unit

(3)

IEEE 754 Floating Point

Standard

~~32~~ n bits code \Rightarrow real number

Real number \Rightarrow code

3.55×10^4 F.p.

Pseudo IEEE 16 bit (10) Used by HP

Float 32 - bit

Double 64 - bit

Quad 128 - bit

$35.5 \times 10^3 = \boxed{3.55 \times 10^4}$ Normal
~~35.5~~ $\boxed{0.355 \times 10^5}$ Standard

only ~~for~~ 1 Representation of

a real #

④

$$A \quad 3.55 \times 10^4$$

$$B \quad + 305 \times 10^3$$

Align exponent

move the decimal point to

the left is equivalent to

shifting the number

point to the right \Rightarrow shift left

$$305.0 \times 10^3 = 30.5 \times 10^4 \quad B'$$

$$A \quad 3.55 \times 10^4 \quad 35.5 \times 10^3$$

$$\text{if } A > B \quad E_A > E_B$$

Align to A \Rightarrow Right shift

(point to left) Lose LSB

5

$$\begin{array}{rcl}
 A & \pm M_A \times 10^{E_A} & 3.55 \times 10^4 \\
 + \\
 B & \pm M_B \times 10^{E_B} & 305.0 \times 10^3
 \end{array}$$

1) Normalize

$$\Rightarrow d.xxx \times 10^E \quad d \neq 0 \quad x \text{ any digit}$$

$$\Rightarrow d.dxxx \times 10^E \quad d \neq 0 \quad x \text{ " " "}$$

$$A' \quad 0.355 \times 10^5$$

$$B' \quad 0.305 \times 10^6$$

2 Align Exponent

$$E_{B'} > E_{A'} \quad \text{align to } B$$

$$\begin{array}{rcl}
 A'' + 0.0355 \times 10^6 & S_A M_A E_A = & \text{~~(0.355)}~~ \pm M_A \times 10^{E_A} \\
 + \\
 B + 0.305 \times 10^6 & S_B M_B E_B = & \pm M_B \times 10^{E_B}
 \end{array}$$

$$\begin{array}{rcl}
 \hline
 C & E_C = & S_C M_C E_C
 \end{array}$$

$$3 \quad E_C = E_{B'}$$

6 $M_C = M_{A''} + M_{B'}$

$\pm (M_{A''} + M_{B'}) \times 10^{E_{B'}}$ result

7 normalize / round \hookrightarrow potentially twice

8 check for overflow

F.P. multiplication

$$\begin{array}{r} 3.55 \times 10^4 \\ \times 305.0 \times 10^3 \\ \hline \end{array} \quad \begin{array}{r} +0.355 \times 10^5 \\ \times +0.305 \times 10^6 \\ \hline \end{array}$$

$$\begin{array}{r} S_A M_A E_A \\ \times S_B M_B E_B \\ \hline \end{array} = S_C M_C E_C$$

$$S_C = S_A (+) S_B$$

1) normalize

2) $S_C = S_A (+) S_B$

3) $E_C = E_A + E_B$ (check for overflow)

7

4) $M_C = M_A * M_B$ (2 ints)

5) Normalize Round

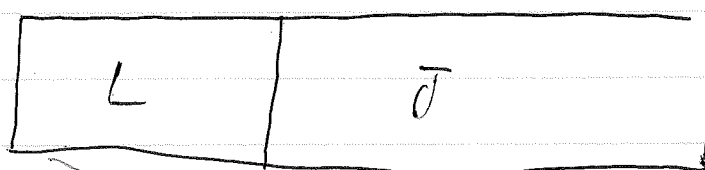
6) check for ~~ovf~~ OVF

1) Naive Representation $\times 10^{-3}$

\Rightarrow improve Naive

a. Efficient coding

b. Reuse (of existing units)



M_A
Signed

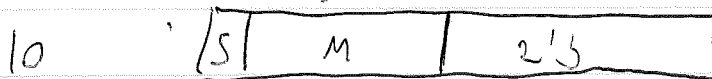
k

E_A

Assume $S+M$ 2's

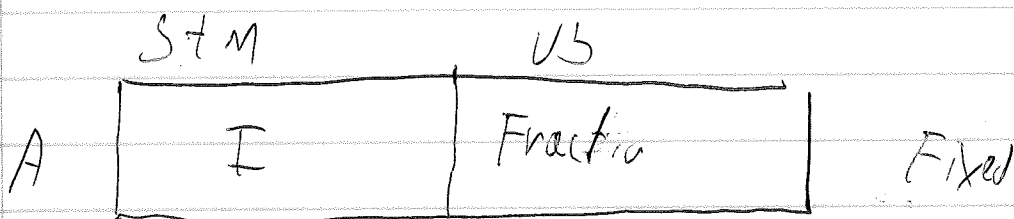
5

5

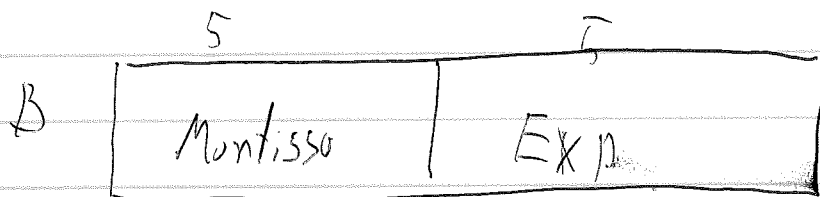


M_A

8



5 ↑ 5



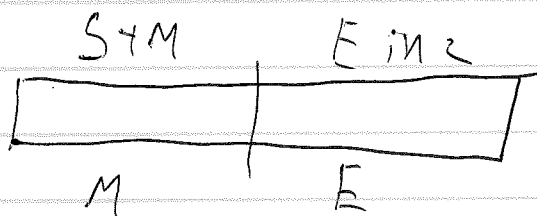
5+5 2's

STM 2's
11110 | 10111

$$A = -14 \frac{23}{32}$$

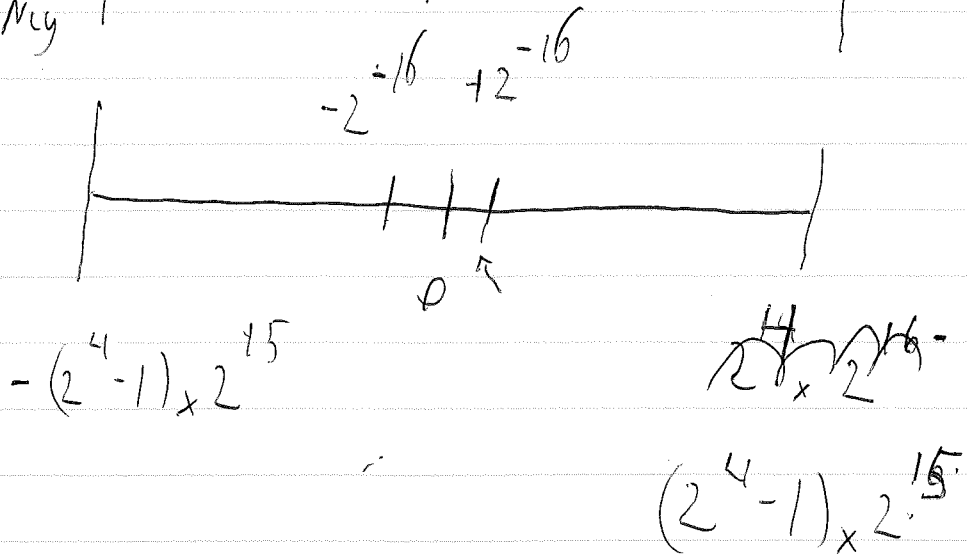
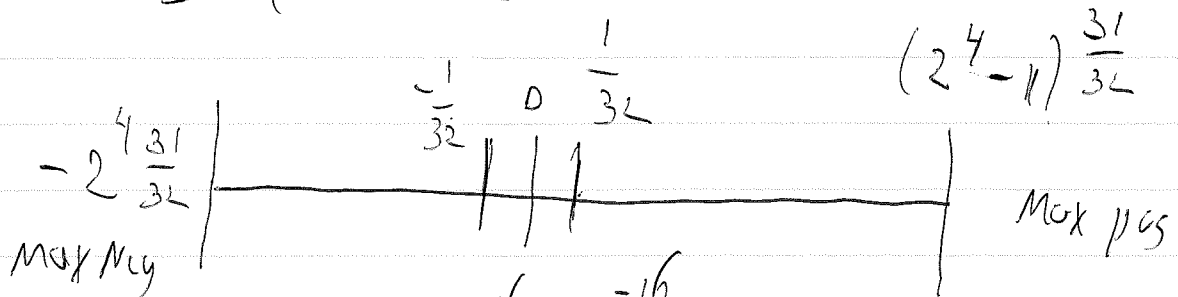
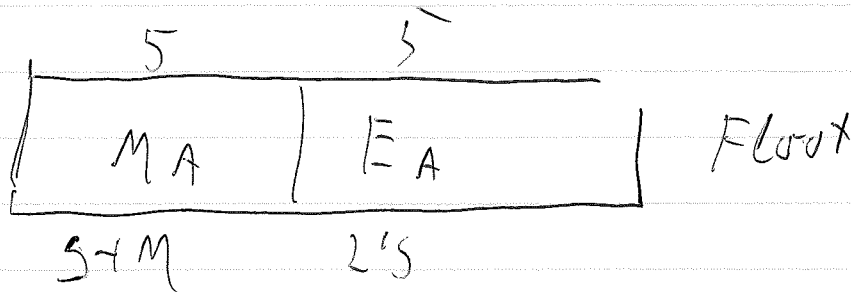
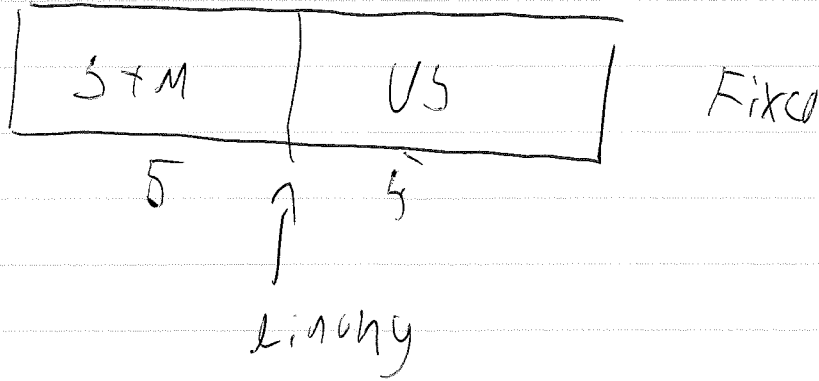
$$B = -14 \times 2^{\cancel{23}} = -9$$

$$(-5) \times M \times 2^E \quad \text{Navid}$$



10111 0100

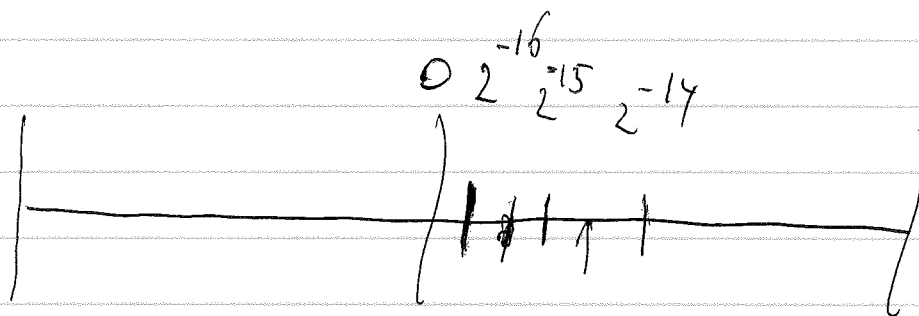
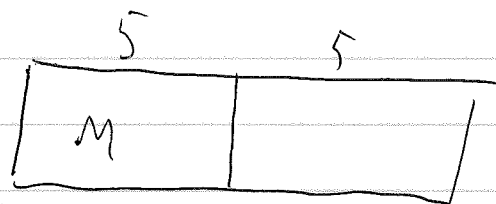
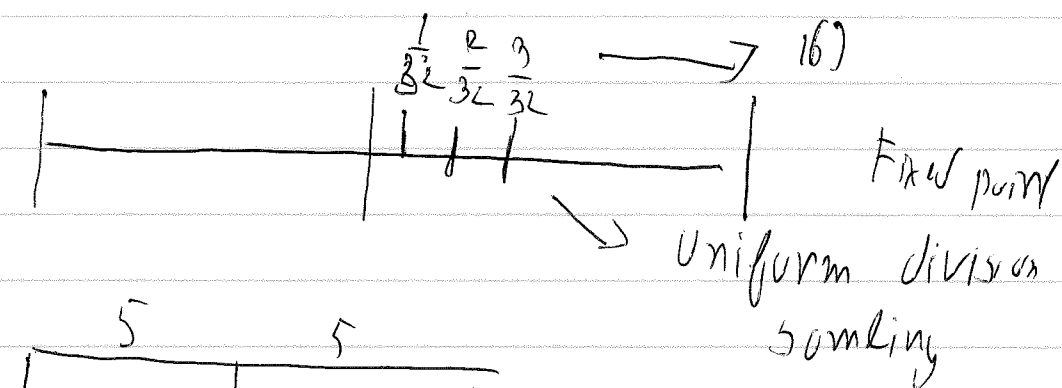
9



Why Fixed point

Why Float Large dynamic range

10

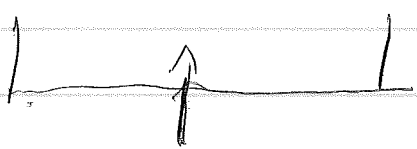


Exponential logarithmic
division

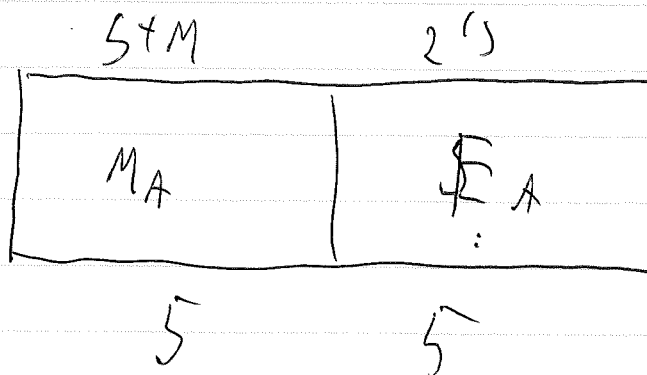


With Fixed error is $\frac{1}{2}$ LSB

5 bits $\frac{1}{2}$ LSB $\frac{1}{64}$ $\frac{1}{2^5} / 2$



14



2^6 Representation

~ 15 Representation

NORMALISE