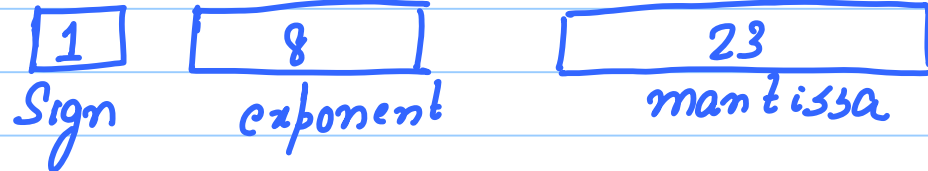


Sept 13th

Note Title

13-09-2011

Floating Point



0	($E=0, m=0$)	NaN	($255(E), m \neq 0$)
$+\infty$	($255(E), m=0$)	Denormal	
$-\infty$	($255(E), m=0$)		($E=0$) ($m \neq 0$)

$$\begin{array}{r}
 7.834 \\
 3.212 \\
 \hline
 11.046
 \end{array}
 \quad (1.1046 \times 10^1)$$

$$(1.105 \times 10^1)$$

Ex 1

$$+ 2^{-50}$$

$$x + y = x$$

Two Floating Point Numbers (A+B)
(A > 0 2 B > 0)

- 1) Take a look at sign bits
- 2) Compare A and B (w.l.g. A > B)
- 3) Align the smaller number, B, by right shifting.
- 4) Add mantissas. (1.x + 1.y = 2.u or 3.v)

5) Normalize

6) Rounding

7) Re-normalize.

Examples

1.011

+

1.11 $\times 2^{-2}$

precision
ε 3 bits

1.011

0.0111 \leftarrow LSB

1.1101

1) Truncate

2) Round to $+\infty$

3) Round to $-\infty$

4) Round:

Rounding an FP Num

FP \rightarrow Int

1. 1.x (0.5 \leq x \leq 0.5) \rightarrow nothing

$\left\{ \begin{array}{l} 0 \leq x < 0.5 \\ x_0 = 0 \end{array} \right\}$
nothing.

$\frac{x}{}$
x₀ x₁ x₂

(0.5 < x < 1) \rightarrow increment

x₀ \neq 0

\rightarrow x = 0.5 (nothing) (s = 0)
 \rightarrow x > 0.5 (s = 1)

Errata

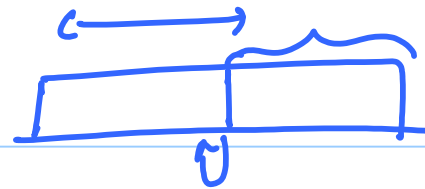
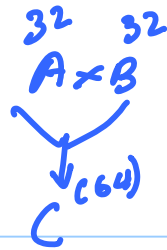
x₀ \rightarrow round bit
x₁ | x₂ | . . . x_n
= s \rightarrow sticky bit.
($x \geq s$)
increment

Round to 1. 23 bits

$\gamma \rightarrow$ 24th bit

s \rightarrow 25 . . . n \leftarrow OR.

Multiplication



[Notes consider a guard & round bit (conceptually same)]

Truncate → Round down (+)
Round up (-).

Round up & Round down have the same connotation.

} Floating Point Assembly
Not covered}
Please Read.

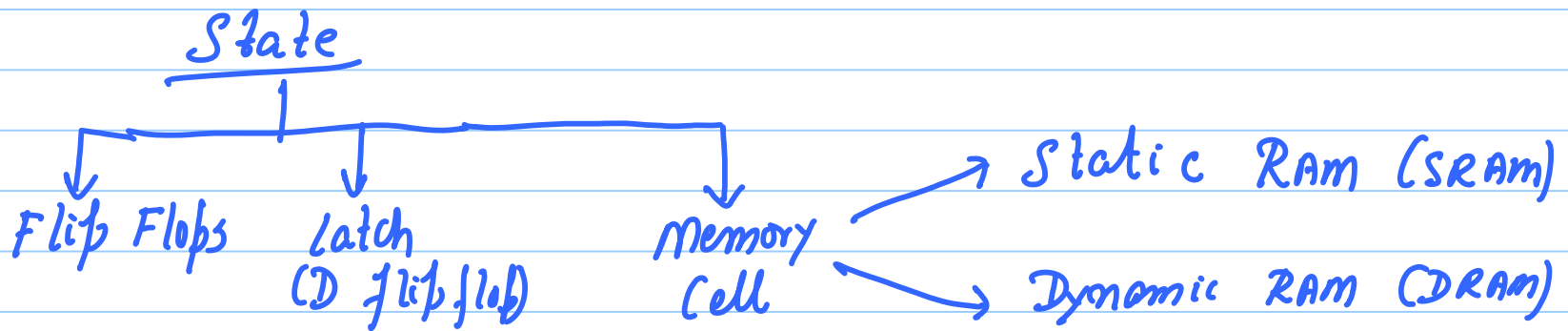
Chapter 4

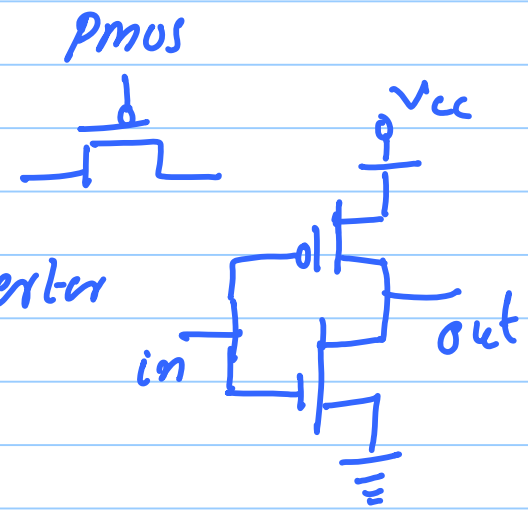
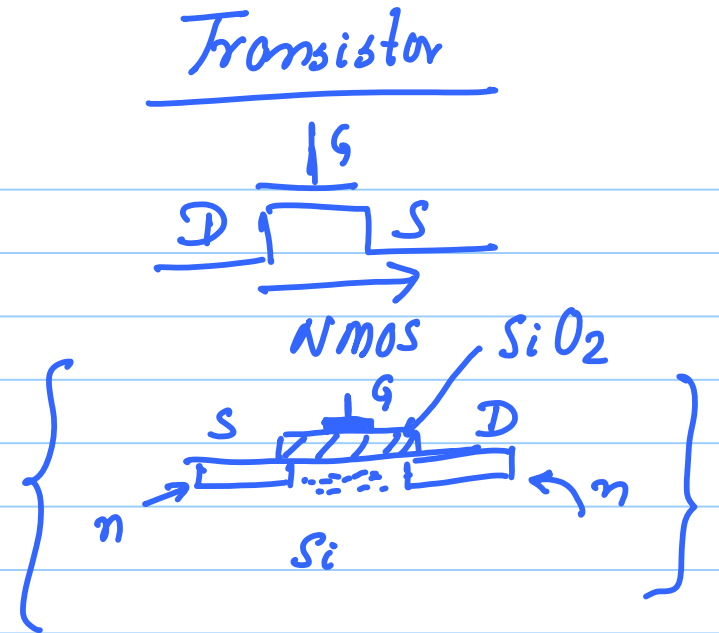
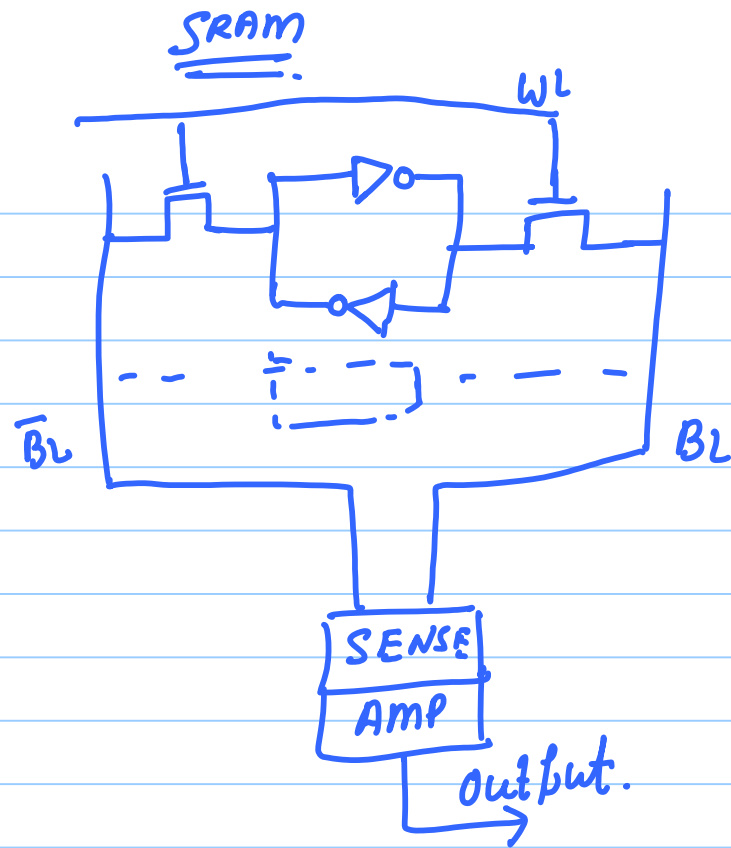
1 2 3 Architecture

→ State (Registers, Mem.)

→ Logic. (Adden Sub,
Divide, ----)

{
4
5
6
} Organization.





Once I enable the wordline.
 BL or \overline{BL} will swing to 1
 other will swing to 0

1) Precharge: B_L and $\overline{B_L}$ to 0.5V $V_{dd} = 1V$

2) Enable WL

3) B_L and $\overline{B_L}$ would start swinging in opp. directions.

4) $450mV$ $550mV$
