

Computer Architecture

Yuqiao Meng

2022-1-28

Contents

1	Overview	4
1.1	Why the prices of two processors differs?	4
2	Digital Design	5
2.1	Abstraction	5
2.2	The Digital Abstraction	5
2.3	Digital Discipline: Binary Values	6
2.4	Decimal to Binary Conversion	6
2.5	Signed Binary Numbers	6
2.5.1	Sign/Magnitude Numbers	6
2.5.2	Two's Complement Numbers	7
2.6	Logic Gates	7
2.7	Logic levels	7
2.7.1	Noise Margins	8
2.7.2	Transfer Characteristics	8
2.8	MOS Transistor	9
2.8.1	nMOS	10
2.8.2	pMOS	10
2.8.3	Why do we need two types of MOS?	11
2.9	cMos NAND Gate	12
2.10	Pseudo-nMOS Gates	12
2.11	Power Consumption	13
2.12	Circuits	14
2.12.1	Composition	14
2.12.2	Types	14
2.12.3	Definitions	15
2.12.4	SOP Forms	15
2.12.5	POS Forms	16
2.12.6	SOP & POS Forms	16
2.13	Boolean	17
2.13.1	Axiom	17
2.13.2	Theorem	17
2.13.3	Simplification methods	18
2.14	K-Maps	19
2.15	Combinational Building Blocks	20
2.15.1	Multiplexers	20

2.15.2	Decoders	21
2.16	Timing	22
2.17	Critical paths	24
2.18	Glitches	24
2.19	Sequential Logic	24
2.19.1	Definitions	24
2.20	State elements	25
2.20.1	Bistable Circuit	25
2.20.2	SR Latch	25
2.20.3	D Latch	26
2.20.4	D Flip-Flop	26
2.21	Timing	27
2.21.1	Input timing constraint	27
2.21.2	Output timing constraint	27
2.21.3	Setup Time Constraint	27
2.21.4	Hold Time Constraint	28
2.22	Clock Skew	28
2.22.1	Setup Time Constraint with Skew	29
2.22.2	Hold Time Constraint with Skew	29
3	Pipeline Overview	30
3.1	Simple Processor Pipeline	30
3.2	Five-Stage Pipeline for a RISC processor	30
3.3	CPU Block Diagram	31
3.4	Hurdles in Pipelining	31
3.4.1	Structural Hazard: Resources	31
3.4.2	Data Hazard: Data dependency	32
3.4.3	Control Hazard: Branch	32
3.5	Data Dependency	33
3.6	Out-of-order pipeline	33
3.7	Memory Wall	34
3.8	Memory Hierarchy Design	34

1 Overview

Intel Xeon Processor



\$13012



\$400

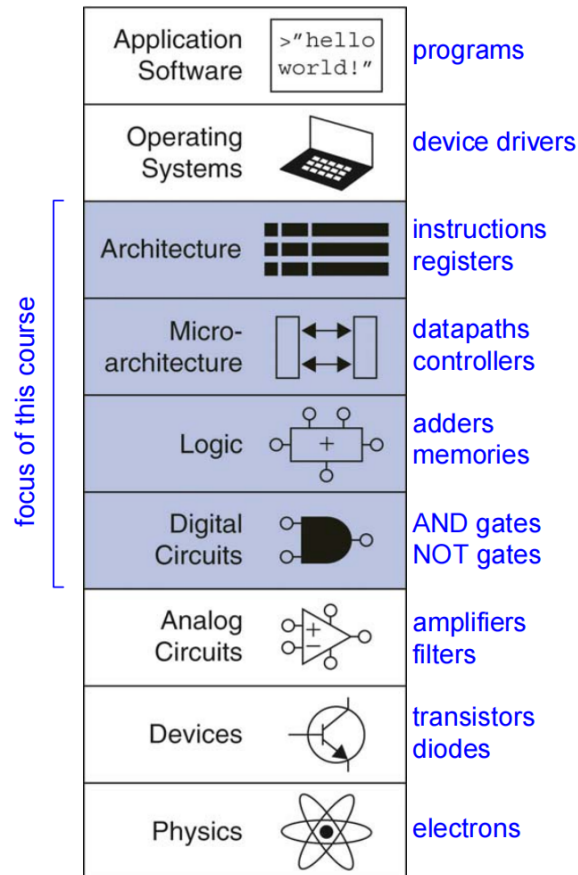
- | | |
|---|---|
| <ul style="list-style-type: none">• Xeon Platinum 8380HL- 28: Cores- 56: Threads- Base Frequency: 2.9 GHz- Max Turbo Frequency: 4.3 GHz- Cache: 38.5MB- Max Memory Size: 4.5 TB- Memory Channels: 6- Max. Memory Speed: 3200MHz- Package Size: 77.5 x 56.5mm | <ul style="list-style-type: none">• Intel i9-10850K- 10: Cores- 20: Threads- Base Frequency: 3.6 GHz- Max Turbo Frequency: 5.2 GHz- Cache: 20MB- Max Memory Size: 128GB- Memory Channels: 2- Max. Memory Speed: 3200MHz- Package Size: 37.5 x 37.5mm |
|---|---|

1.1 Why the prices of two processors differs?

1. The prices increase exponentially as the number of Cores and Threads increase. That's because every core has a bad possibility, so the difficult of making a processor with many cores much harder.
2. Significant: The memory Chaneels of XEON is three times by i9, which means it has three times as many pins as i9 has.

2 Digital Design

2.1 Abstraction



2.2 The Digital Abstraction

- Most physical variables are continuous
 - Voltage on a wire
 - Frequency of an oscillation
 - Position of a mass
- Digital abstraction considers discrete subset of values

2.3 Digital Discipline: Binary Values

- Two discrete Values:
 - 1's and 0's
 - 1, true, high
 - 0, false, low
- 1 and 0: voltage levels, rotating gears, fluid levels
- Digital circuits use voltage levels to represent 1 and 0

2.4 Decimal to Binary Conversion

Method: repeatedly divided by 2, remainders goes in next most significant bit

$$\begin{array}{rcl} 53_{10} = & 53/2 = 26 \text{ R}1 & \\ & 26/2 = 13 \text{ R}0 & \\ & 13/2 = 6 \text{ R}1 & \\ & 6/2 = 3 \text{ R}0 & \\ & 3/2 = 1 \text{ R}1 & \\ & 1/2 = 0 \text{ R}1 & \end{array} \quad = 110101_2$$

2.5 Signed Binary Numbers

2.5.1 Sign/Magnitude Numbers

Problems

- Has two 0 values, positive 0 and negative 0
- Addition of a negative and a positive will fail

2.5.2 Two's Complement Numbers

Conversion from positive to negative:

- Invert every bit
- Add 1

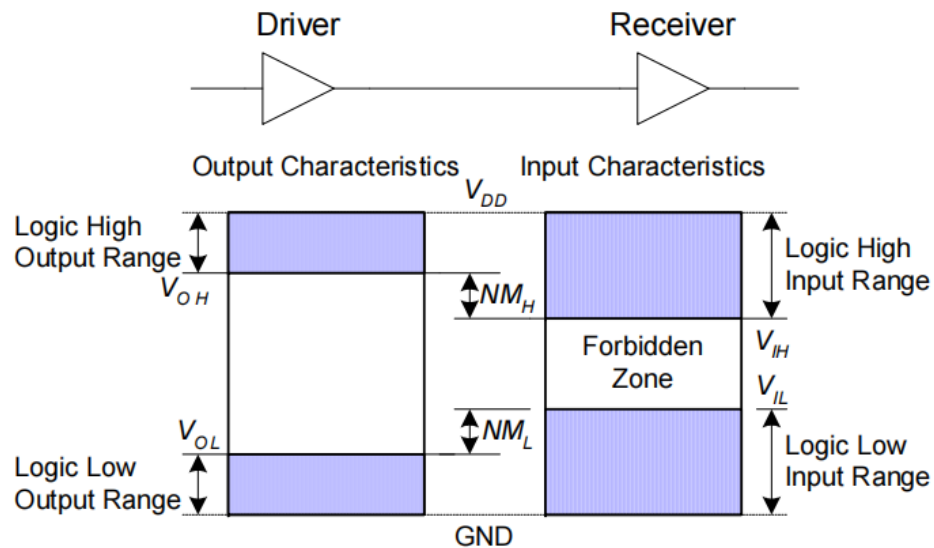
2.6 Logic Gates

- AND
- OR
- XOR
- NAND
- NOR
- XNOR

2.7 Logic levels

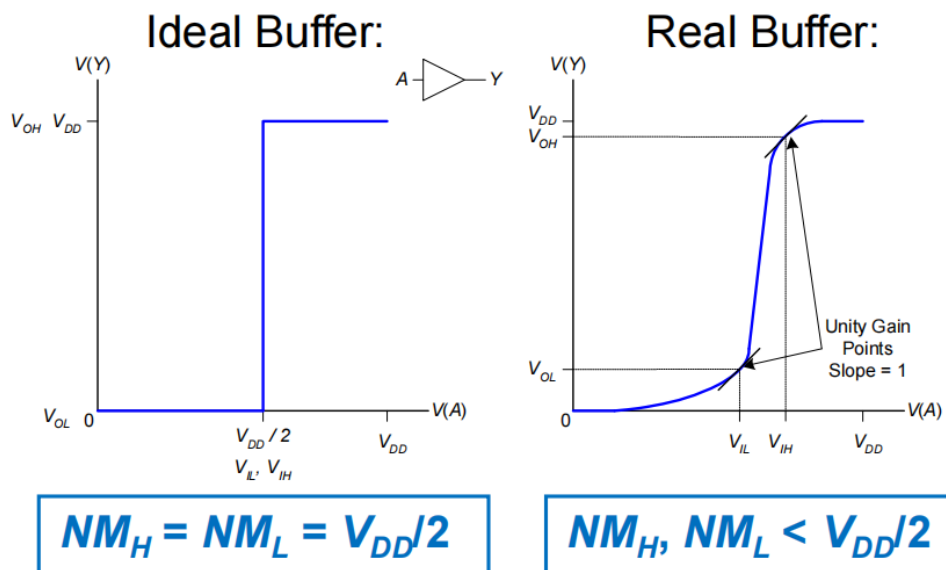
Discrete voltages represent 1 and 0, but there is noise that will degrade the signal

2.7.1 Noise Margins



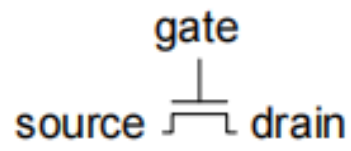
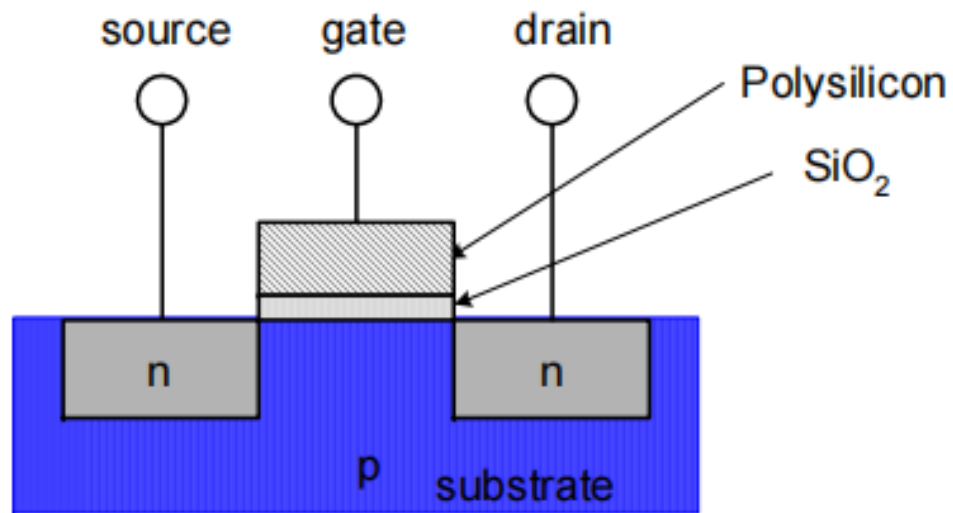
- High Noise Margin: $NM_H = V_{OH} - V_{IH}$
- Low Noise Margin: $NM_L = V_{OL} - V_{IL}$

2.7.2 Transfer Characteristics



2.8 MOS Transistor

- Polysilicon Gate: control the transistor
- Oxide insulator: keep electrons to form a path for electricity
- Doped Silicon

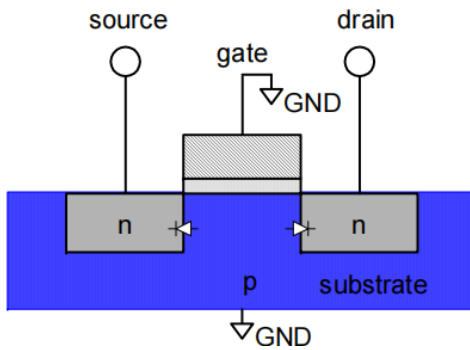


nMOS

2.8.1 nMOS

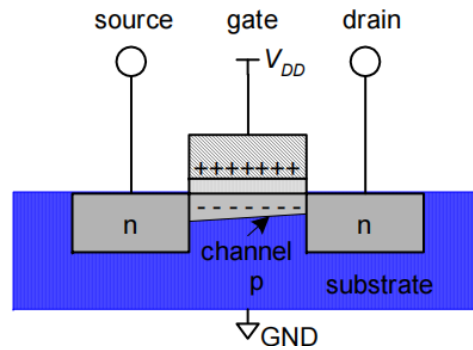
Gate = 0

OFF (no connection between source and drain)



Gate = 1

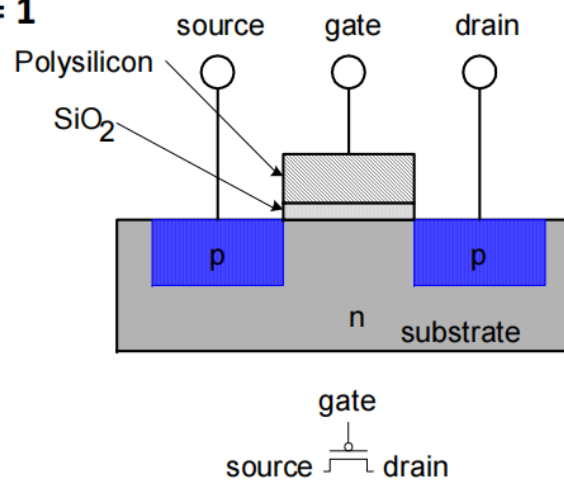
ON (channel between source and drain)



2.8.2 pMOS

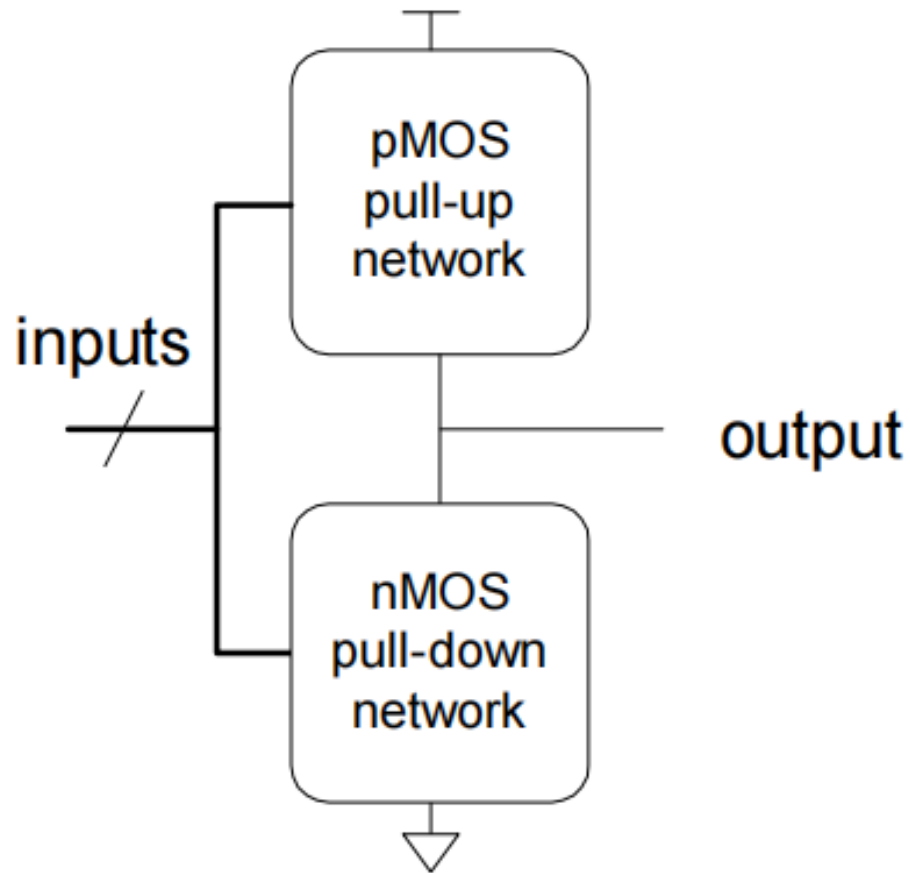
pMOS transistor is opposite

- **ON** when **Gate = 0**
- **OFF** when **Gate = 1**



2.8.3 Why do we need two types of MOS?

- nMOS: pass good 0's, so connect source to GND
- pMOS: pass good 1's, so connect source to VDD



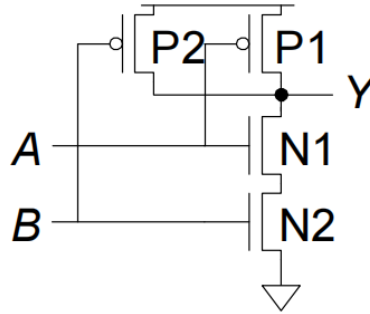
2.9 cMos NAND Gate

NAND



$$Y = \overline{AB}$$

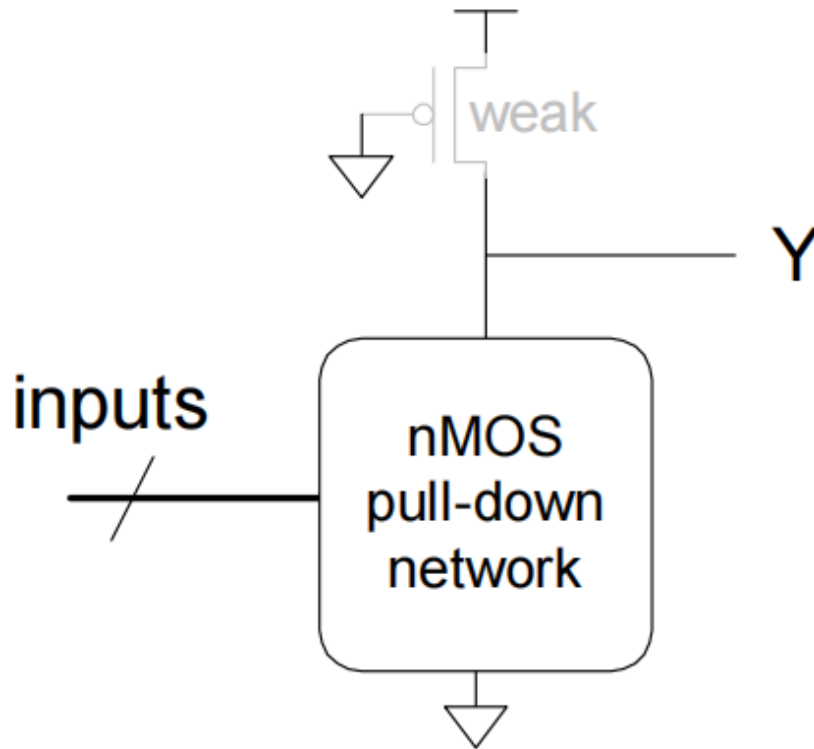
A	B	Y
0	0	1
0	1	1
1	0	1
1	1	0



A	B	P1	P2	N1	N2	Y
0	0	ON	ON	OFF	OFF	1
0	1	ON	OFF	OFF	ON	1
1	0	OFF	ON	ON	OFF	1
1	1	OFF	OFF	ON	ON	0

2.10 Pseudo-nMOS Gates

- Replace pull-up network with weak pMOS transistor that is always on
- Tradeoff: need more energy



2.11 Power Consumption

- Dynamic power consumption: Power to charge transistor gate capacitances
 - $P_{Dynamic} = \frac{1}{2}CV_{DD}^2f$
 - f: Frequency
 - CV_{DD}^2f : energy to charge a capacitance
- Static power consumption
 - power consumed when no gates are switching
 - caused by quiescent supply current I_{DD} (there is more and more current leakage for transistor is smaller and smaller, and the distance between source and drain is shorter and shorter, so there is more and more quiescents leaked)

- $P_{static} = I_{DD}V_{DD}$
- $P = \frac{1}{2}CV_{DD}^2f + I_{DD}V_{DD}$

2.12 Circuits

2.12.1 Composition

- inputs
- outputs
- functional specification
- timing specification

2.12.2 Types

- Combinational
 - Memoryless
 - output \rightarrow current value of inputs
- Sequential
 - has memory
 - output \rightarrow previous and current value of inputs

2.12.3 Definitions

- **Complement:** variable with a bar over it
 $\bar{A}, \bar{B}, \bar{C}$
- **Literal:** variable or its complement
 $A, \bar{A}, B, \bar{B}, C, \bar{C}$
- **Implicant:** product of literals
 $ABC, \bar{A}C, BC$
- **Minterm:** product that includes all input variables
 $\bar{A}\bar{B}\bar{C}, \bar{A}B\bar{C}, A\bar{B}\bar{C}$
- **Maxterm:** sum that includes all input variables
 $(A+\bar{B}+C), (\bar{A}+B+\bar{C}), (\bar{A}+\bar{B}+C)$

2.12.4 SOP Forms

- Each row has a minterm(products of literals)

A	B	Y	minterm	minterm name
0	0	0	$\bar{A} \bar{B}$	m_0
0	1	1	$\bar{A} B$	m_1
1	0	0	$A \bar{B}$	m_2
1	1	1	$A B$	m_3

$$Y = F(A, B) = \bar{A}B + AB = \Sigma(1, 3)$$

2.12.5 POS Forms

- Each row has a maxterm(sum of literals)

A	B	Y	maxterm	maxterm name
0	0	0	$A + B$	M_0
0	1	1	$A + \overline{B}$	M_1
1	0	0	$\overline{A} + B$	M_2
1	1	1	$\overline{A} + \overline{B}$	M_3

$$Y = F(A, B) = (A + B)(\overline{A} + \overline{B}) = \Pi(0, 2)$$

2.12.6 SOP & POS Forms

SOP – sum-of-products

O	C	E	minterm
0	0	0	$\overline{O} \overline{C}$
0	1	0	$\overline{O} C$
1	0	1	$O \overline{C}$
1	1	0	$O C$

$$E = O\overline{C}$$

$$= \Sigma(2)$$

POS – product-of-sums

O	C	E	maxterm
0	0	0	$O + C$
0	1	0	$O + \overline{C}$
1	0	1	$\overline{O} + C$
1	1	0	$\overline{O} + \overline{C}$

$$E = (O + C)(O + \overline{C})(\overline{O} + \overline{C})$$

$$= \Pi(0, 1, 3)$$

2.13 Boolean

2.13.1 Axiom

Number	Axiom	Dual	Name
A1	$B = 0 \text{ if } B \neq 1$	$B = 1 \text{ if } B \neq 0$	Binary Field
A2	$\overline{0} = 1$	$\overline{1} = 0$	NOT
A3	$0 \bullet 0 = 0$	$1 + 1 = 1$	AND/OR
A4	$1 \bullet 1 = 1$	$0 + 0 = 0$	AND/OR
A5	$0 \bullet 1 = 1 \bullet 0 = 0$	$1 + 0 = 0 + 1 = 1$	AND/OR

Dual: Replace: \bullet with $+$
 0 with 1

2.13.2 Theorem

#	Theorem	Dual	Name
T6	$B \bullet C = C \bullet B$	$B + C = C + B$	Commutativity
T7	$(B \bullet C) \bullet D = B \bullet (C \bullet D)$	$(B + C) + D = B + (C + D)$	Associativity
T8	$B \bullet (C + D) = (B \bullet C) + (B \bullet D)$	$B + (C \bullet D) = (B + C) (B + D)$	Distributivity
T9	$B \bullet (B + C) = B$	$B + (B \bullet C) = B$	Covering
T10	$(B \bullet C) + (B \bullet \overline{C}) = B$	$(B + C) \bullet (B + \overline{C}) = B$	Combining
T11	$(B \bullet C) + (\overline{B} \bullet D) + (C \bullet D) = (B \bullet C) + (\overline{B} \bullet D)$	$(B + C) \bullet (\overline{B} + D) \bullet (C + D) = (B + C) \bullet (\overline{B} + D)$	Consensus

Warning: T8' differs from traditional algebra:
OR (+) distributes over AND (\bullet)

2.13.3 Simplification methods

- **Distributivity (T8, T8')**
 $B(C+D) = BC + BD$
 $B + CD = (B+C)(B+D)$
- **Covering (T9')**
 $A + AP = A$
- **Combining (T10)**
 $\overline{PA} + PA = P$
- **Expansion**
 $P = \overline{PA} + PA$
 $A = A + AP$
- **Duplication**
 $A = A + A$
- **“Simplification” theorem**
 $\overline{PA} + A = P + A$
 $PA + \overline{A} = P + \overline{A}$

2.14 K-Maps

$$PA + P\bar{A} = P$$

<i>A</i>	<i>B</i>	<i>C</i>	<i>D</i>	<i>Y</i>
0	0	0	0	1
0	0	0	1	0
0	0	1	0	1
0	0	1	1	1
0	1	0	0	0
0	1	0	1	1
0	1	1	0	1
0	1	1	1	1
1	0	0	0	1
1	0	0	1	1
1	0	1	0	1
1	0	1	1	0
1	1	0	0	0
1	1	0	1	0
1	1	1	0	0
1	1	1	1	0

<i>Y</i> <i>CD</i> \ <i>AB</i>		00	01	11	10
		00	01	11	10
	00	1	0	0	1
	01	0	1	0	1
	11	1	1	0	0
	10	1	1	0	1

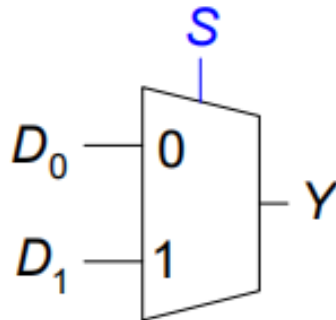
$$Y = \bar{A}\bar{C} + \bar{A}BD + A\bar{B}\bar{C} + \bar{B}\bar{D}$$

2.15 Combinational Building Blocks

2.15.1 Multiplexers

\lg_2^N bit select input

2:1 Mux

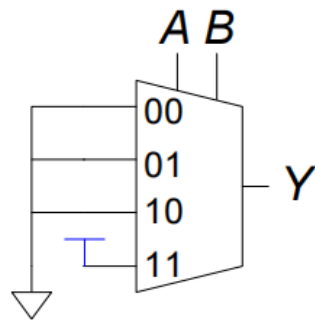


S	D_1	D_0	Y	S	Y
0	0	0	0	0	D_0
0	0	1	1	1	D_1
0	1	0	0		
0	1	1	1		
1	0	0	0		
1	0	1	0		
1	1	0	1		
1	1	1	1		

- Logic use

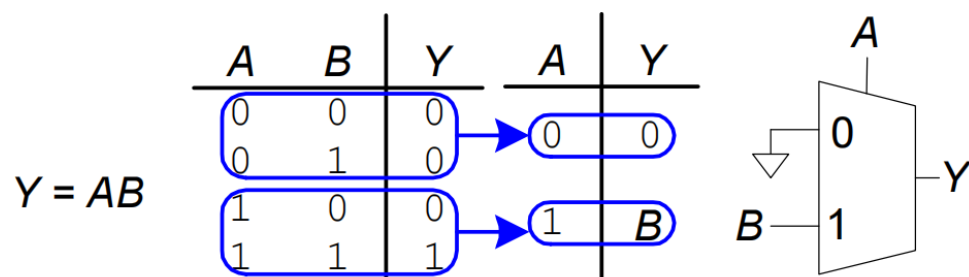
A	B	Y
0	0	0
0	1	0
1	0	0
1	1	1

$$Y = AB$$

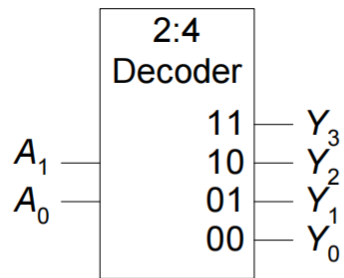


2.15.2 Decoders

N inputs, 2^N outputs



- Logic use

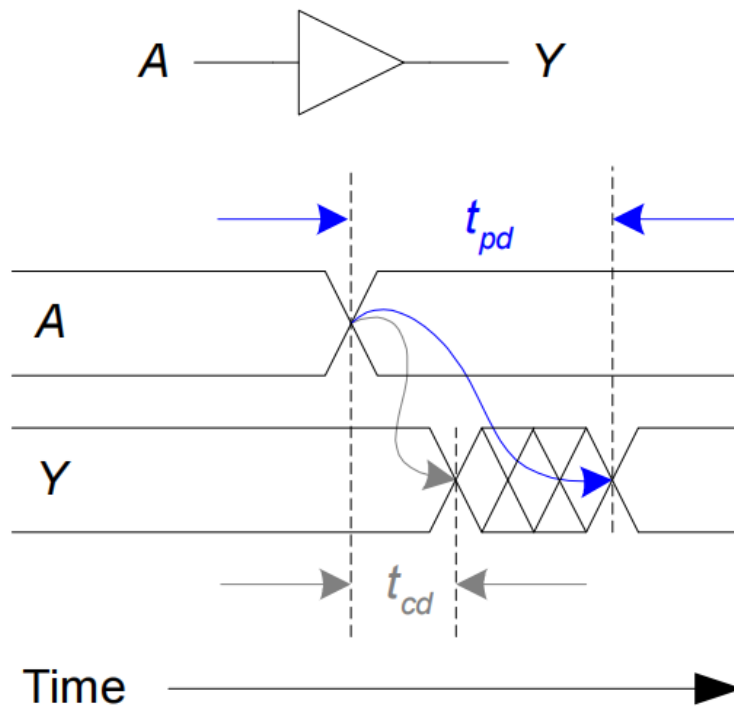


A_1	A_0	Y_3	Y_2	Y_1	Y_0
0	0	0	0	0	1
0	1	0	0	1	0
1	0	0	1	0	0
1	1	1	0	0	0

2.16 Timing

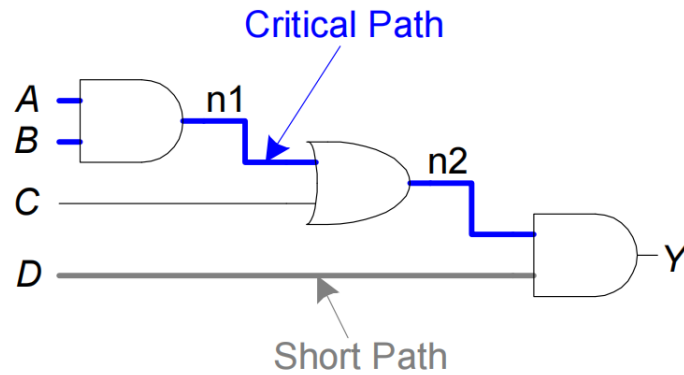
Because of wire resistance and we need time to fully charge the capacity

- t_{pd} Propagation delay: the time from input starting to change to output becoming stable
- t_{cd} Contamination delay: the time from input starting to change to output starting to change



- Different rising and falling delays
- multiple inputs and outputs
- circuits slow down when hot

2.17 Critical paths



Critical (Long) Path: $t_{pd} = 2t_{pd_AND} + t_{pd_OR}$

Short Path: $t_{cd} = t_{cd_AND}$

2.18 Glitches

When a single input changes causes an output to change several times

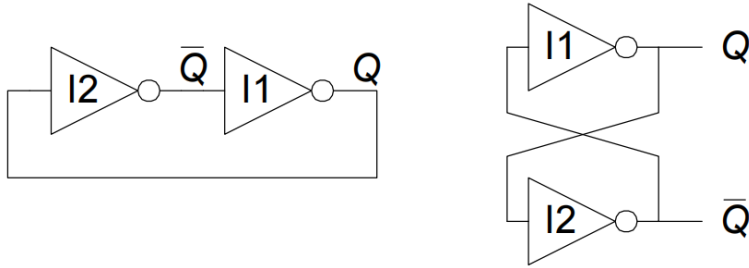
2.19 Sequential Logic

2.19.1 Definitions

- State: influence the future behavior
- Latches and flip-flops
- Synchronous sequential circuits

2.20 State elements

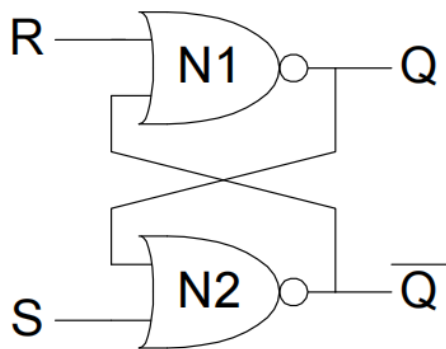
2.20.1 Bistable Circuit



2.20.2 SR Latch

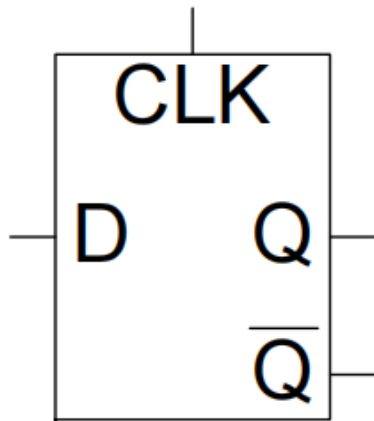
S: set, R: reset

- S: 0, R: 0: previous
- S: 1, R: 0: 1
- S: 0, R: 1: 0
- S: 1, R: 1: invalid



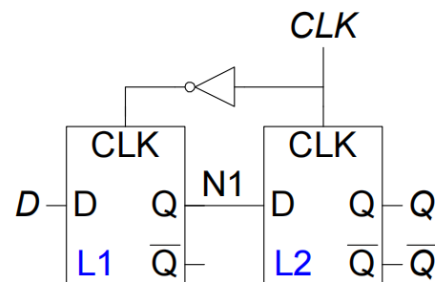
2.20.3 D Latch

D Latch Symbol



- CLK: controls when the output changes
- D: controls what the output changes to

2.20.4 D Flip-Flop

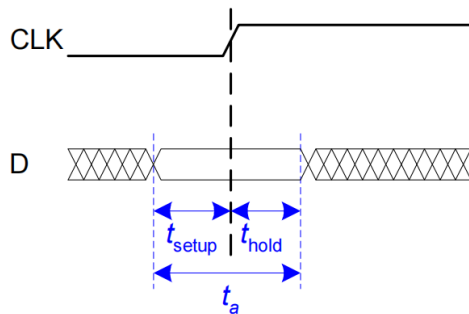


Samples D on rising edge of CLK

2.21 Timing

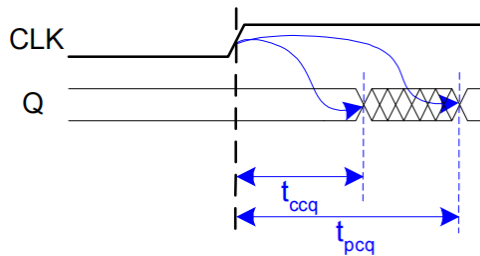
2.21.1 Input timing constraint

- t_{setup} : time before clock edge data must be stable
- t_{hold} : time after clock edge data must be stable
- Aperture time t_a : time around clock edge data must be stable



2.21.2 Output timing constraint

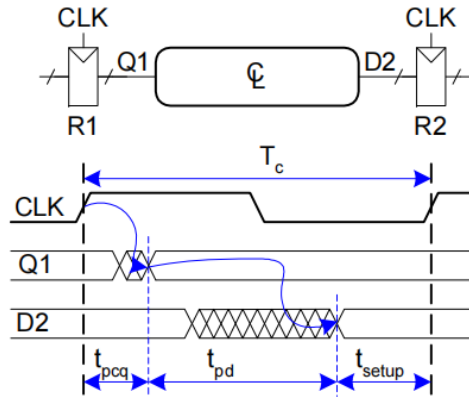
- t_{pcq} : time after clock edge that the output Q is guaranteed to be stable
- t_{ccq} : time after clock edge data might be unstable



2.21.3 Setup Time Constraint

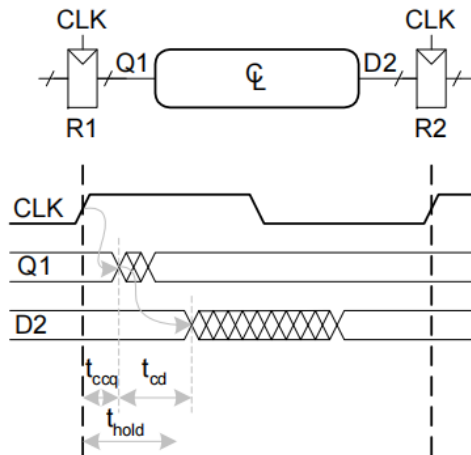
- input to register R2 must be stable at least t_{setup} before clock edge
- clock edge time i Q be stable time + propagation time + set up time
- $T_c \geq t_{pcq} + t_{pd} + t_{setup}$

- $t_{pd} \leq T_c - (t_{pcq} + t_{setup})$



2.21.4 Hold Time Constraint

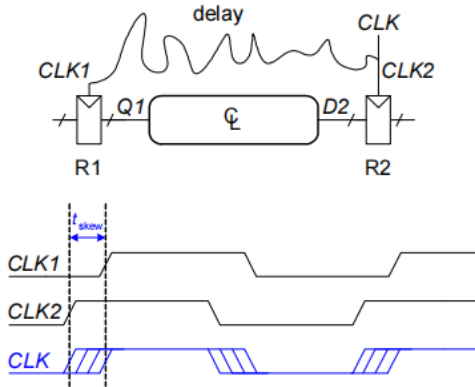
- input to register R2 must be stable at least t_{hold} after clock edge
- hold time : Q1 starts to change time + D2 starts to change time
- $t_{hold} < t_{ccq} + t_{cd}$
- $t_{cd} > t_{hold} - t_{ccq}$



2.22 Clock Skew

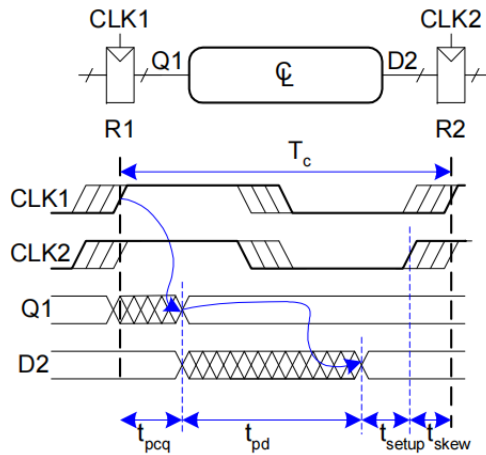
- Skew: difference between two clock edges

- Perform worst case analysis to guarantee dynamic discipline is not violated for any register



2.22.1 Setup Time Constraint with Skew

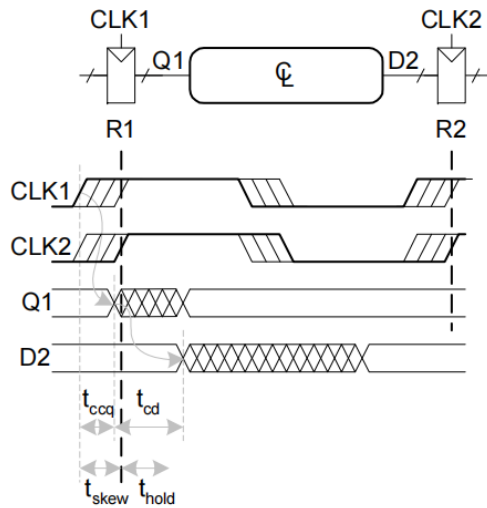
- in worst case, CLK2 is earlier than CLK1
- $T_c \geq t_{pcq} + t_{pd} + t_{setup} + t_{skew}$
- $t_{pd} \leq T_c - (t_{pcq} + t_{setup} + t_{skew})$



2.22.2 Hold Time Constraint with Skew

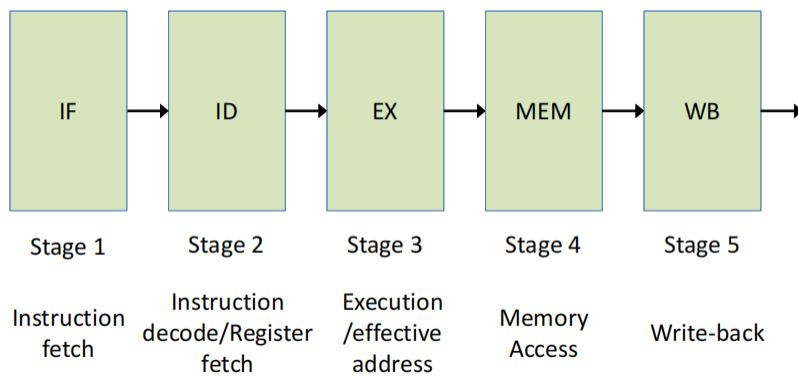
- in worst case, CLK2 is later than CLK1

- $t_{hold} + t_{skew} < t_{ccq} + t_{cd}$
- $t_{cd} > t_{hold} + t_{skew} - t_{ccq}$



3 Pipeline Overview

3.1 Simple Processor Pipeline



3.2 Five-Stage Pipeline for a RISC processor

- IF(Instruction Fetch)

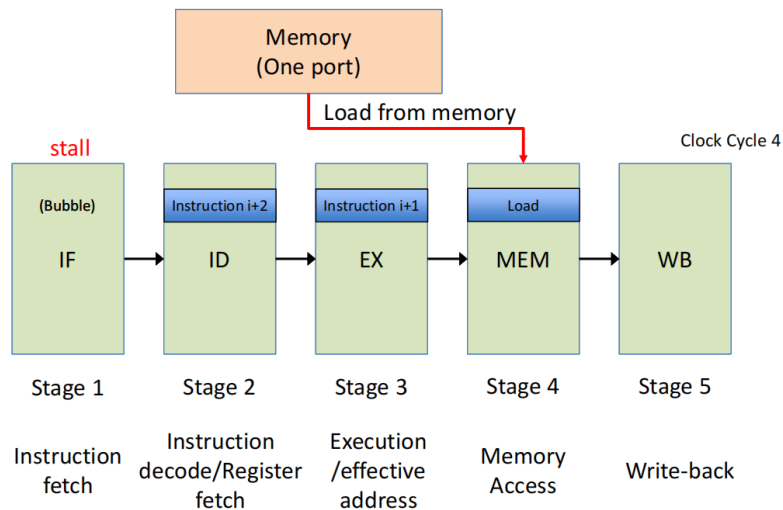
- ID(Instruction Decode/Register Fetch)
- EX(Execution/Effective address)
- MEM(Memory access): store instruction
- WB(Write-back): write the result back into the REGISTER FILE

3.3 CPU Block Diagram

3.4 Hurdles in Pipelining

3.4.1 Structural Hazard: Resources

- All the overlapped instructions in pipeline requires pipelining of functional units and duplication of resources
- Resource conflict can happen when some combination of instructions cannot be accommodated.

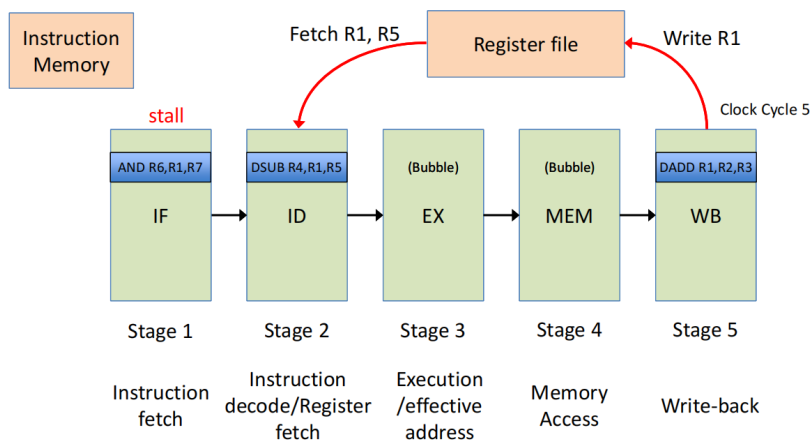


3.4.2 Data Hazard: Data dependency

```

DADD R1, R2, R3 // R2+R3 = R1
DSUB R4, R1, R5 // R1-R5 = R4
AND R6, R1, R7 // R1 & R7 = R6
OR R8, R1, R9 // R1 | R9 = R8
XOR R10, R1, R11 // R1 ^ R11 = R10
  
```

- All the instructions after the DADD use the result of DADD
- The value of R1 is ready at WB stage of DADD, but DSUB needs it at its ID stage



3.4.3 Control Hazard: Branch

```

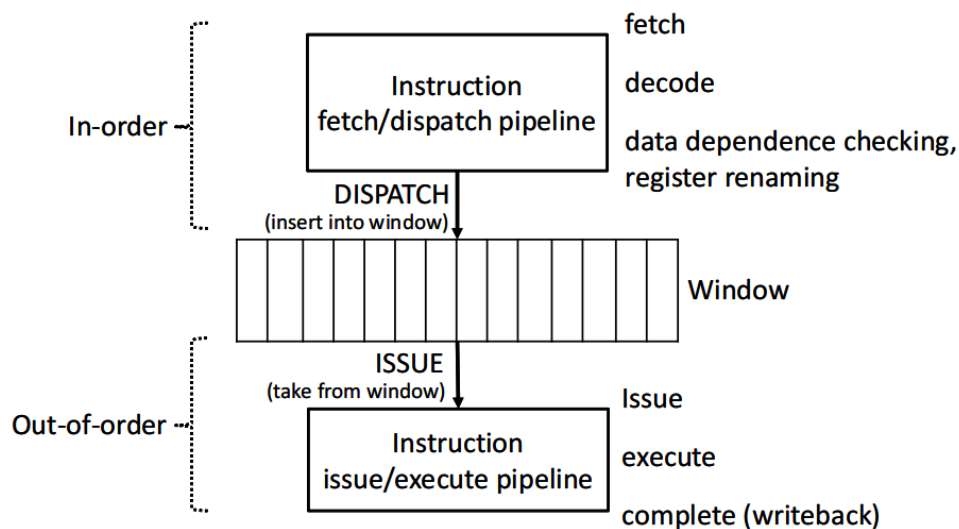
Not taken → BEQZ T, R1
            → DADD R2, R3, R4
            ...
Taken      → T: OR R6, R7, R8
            → XOR R9, R10, R11
            ...
  
```


- Branch may change or not change the PC
- Which instruction to fetch after the branch?
- IF stalls until the branch condition is evaluated.

3.5 Data Dependency

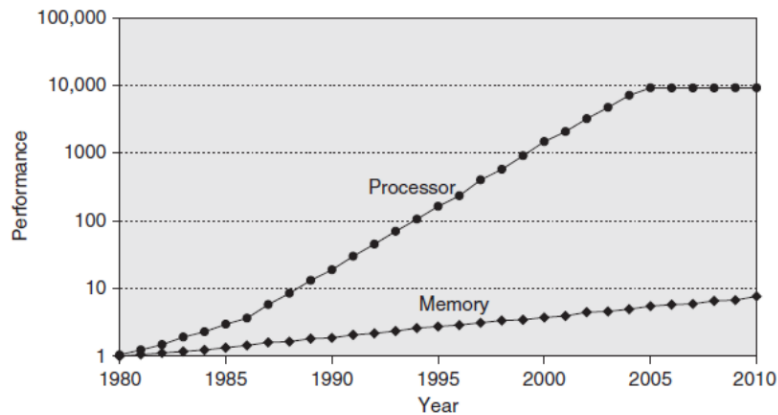
- True-dependence(pure-dependence, flow-dependence): an instruction depends on the result of a previous instruction
- Anti-dependence: an instruction requires a value that is later updated(because we are not executing instructions in order)
- Output-dependence: the ordering of instructions affects the final output result

3.6 Out-of-order pipeline



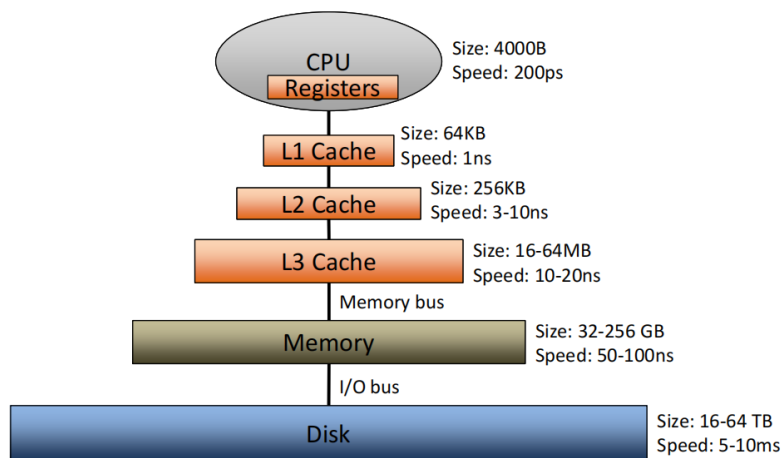
- Instruction fetching do its job in order.
- Then put these instructions into a window
- When a instruction in the window gets ready, execute pipeline executes it

3.7 Memory Wall



- Reason: The growing disparity of speed between CPU and main memory
- Given this trend, memory would become a performance bottleneck

3.8 Memory Hierarchy Design



- Large memory is slow
- processor wants large and fast memory
- Solution: take advantage of locality