## Panel 1 (handwritten)

EE222 Lecture 15   Feb. 26, 2019
Low Power Design of Logic Circuits:

- Combinational Circuits
  - Mixed Signal (digital and analog) Circuits
    for Processing-in-Memory (PIM)
- Sequential circuits
- Clock distribution

## Panel 2

### Combinational Vs. Sequential Circuits

Combinational (or Combinatorial) Networks/Circuits
- Circuit without storage
- Outputs depend only on its current inputs
- Examples: NAND gate, look-up table (LUT)

Vs.

Sequential Networks/Circuits
- Circuit with storage elements
- Outputs depend on present inputs and also on history of inputs
- Examples: RAM, finite state machine (FSM)
- Normally, combinational cells AND storage elements

## Panel 3



ISSCC 2018 / SESSION 31 / COMPUTATION IN MEMORY FOR MACHINE LEARNING / 31.1

31.1  Conv-RAM: An Energy-Efficient SRAM with Embedded Convolution Computation for Low-Power CNN-Based Machine Learning Applications

Avishek Biswas, Anantha P. Chandrakasan

Massachusetts Institute of Technology, Cambridge, MA

## Panel 4 (handwritten)

Computation in Memory (CIM)
(Also called Processing in Memory (PIM)

Ref. Convolution-RAM: An Energy Efficient SRAM with Embedded Convolution Computation for Low Power CNN (Convolution Neural Network) - based Machine Learning Applications
ISSCC 2018   A. Biswas, et al. MIT



(conventional)

## Panel 5



Figure 31.1.1: Concept of embedded convolution computation, performed by averaging in SRAM, for binary-weight convolutional neural networks.

Figure 31.1.2: Overall architecture of the Convolution-SRAM (CSRAM) showing local arrays, column-wise DACs and row-wise ADCs to implement convolution as weighted averaging.

## Panel 6



Figure 31.1.3: Schematic and timing diagram for the column-wise GBL_DAC, which converts the convolution digital input to an analog pre-charge voltage for the SRAM.

Figure 31.1.4: Architecture for the row-wise multiply-and-average (MAV₁) and CSH_ADC. Operational waveforms for convolution computation.

Figure 31.1.5: Measured performance of GBL_DAC and CSH_ADC in the CSRAM array. Also shown is the effect of the offset cancelation technique.

Figure 31.1.6: Energy and output distribution measured results for the first two convolutional-layers (C1, C3) of the LeNet-5 CNN for the MNIST dataset. Table showing comparison to prior work.



$$Y = \sum_i W_i \times X_i = \sum_i \alpha \, w_i \times X_i = \frac{M}{N} \sum_i w_i \times X_i$$
$$w_i \in (+1, -1)$$

$$Y_A = \frac{1}{N} \sum W_i \times X_i = ADC\left[ \frac{1}{N} \sum w_i \times DAC(X_i) \right]$$

weight multiply & average

$$V_{P\,average} = \frac{1}{N} \sum_{i=0}^{N-1} |w_i| \times V_{A_i}$$

$w_i \in (+1, -1)$ ← reduces the storage requirement for weights



Local SRAM array (16×64)



Technology, Machine Learning, Comparison mode, Input, Voltage, SRAM Size, Thru put

| Technology | Machine Learning | Comparison mode | Input | Voltage | SRAM Size | Thru put |
|---|---|---|---|---|---|---|
| 65 nm | CNN | Analog | 7b | 1.2 V | 2 KB | 10.7 GOPS |
| 65 nm (=ISCC2016) | CNN | Digital | 16 b | 1.2 V | 36 KB | 64 GOPS |

Energy Efficiency

Figure of Merit (FOM)
$$= \frac{Thruput \times Energy\ Efficiency}{SRAM\ Size}$$

20× improvement { 28.1 TOPS/Watt / 1.42 TOPS/Watt }
$= \frac{150.3}{2.5}$ × 70 improvement



## Sequential Circuits

❑ Sequential circuits are a function of both the current state and the previous state.

❑ In other words, they have *memory*.

$I_n$ Inputs → COMBINATIONAL LOGIC → $O_n$ Outputs

Current State = register output
$Q_n = D_{n-1}$

$Q_n = D_{n-1}$

Registers  Q  D
$D_n$ for Next state

↑CLK

VLSI



Edge-triggered Flip-flop

• Symbol of edge-triggered D flip-flop

Positive-edge triggered        Negative-edge triggered

# Flip-Flop Timing

- Set-up time: $t_s$
  - Input needs to be stable before trigger
- Hold time: $t_h$
  - Input needs to be stable after trigger
- Propagation delay: $t_p$
  - Some delay from trigger to output change
- Example: Negative edge triggered flip-flip

Clock

$t_s$ : $t_h$

$t_p$

# Sequential Circuit Description



Next state    Present state

X

input

D

C

A

$\overline{A}$

D

C

B

$\overline{B}$

$B$

$A+B$

Clock

$(A+B)\ \overline{X}$

$A$

$\overline{X}$

Y

$X$

At the clock trigger, the next state will be read and transferred to the present state

# Input Equations

$$A_{next} = A_{present}X + B_{present}X$$

$$B_{next} = A'_{present}X$$

Next state in terms of input and present state

$$Y = (A_{present} + B_{present})\overline{X'}$$

Output in terms of input and present state

# State Table

| Present State | | Input | Next State | | Output |
|---|---|---|---|---|---|
| A | B | X | A | B | Y |
| 0 | 0 | 0 | 0 | 0 | 0 |
| 0 | 0 | 1 | 0 | 1 | 0 |
| 0 | 1 | 0 | 0 | 0 | 1 |
| 0 | 1 | 1 | 1 | 1 | 0 |
| 1 | 0 | 0 | 0 | 0 | 1 |
| 1 | 0 | 1 | 1 | 0 | 0 |
| 1 | 1 | 0 | 0 | 0 | 1 |
| 1 | 1 | 1 | 1 | 0 | 0 |

# State Diagram



0/0

1/0

00

01

0/1

0/1

0/1

1/0

10

11

1/0

1/0

# Example of Moore Model



X
Y

D    $Q$

C    $\overline{Q}$

A

Z

Clock

$$A_{next} = A_{present} + XY$$

$$Z = A_{present}$$

| X | Y | $A_{present}$ | $A_{next}$ |
|---|---|---|---|
| 0 | 0 | 0 | 0 |
| 0 | 0 | 1 | 1 |
| 0 | 1 | 0 | 0 |
| 0 | 1 | 1 | 1 |
| 1 | 0 | 0 | 0 |
| 1 | 0 | 1 | 1 |
| 1 | 1 | 0 | 1 |
| 1 | 1 | 1 | 1 |

00,01,10

11

0/0

1/1

00,01,10

A True single-phase clock edge triggered flip-flop w/Reset (R)





FIGURE 1: (a) C²MOS logic and (b) an example of single-phase clocking.

FIGURE 2: Domino logic.

FIGURE 3: NORA logic.

FIGURE 4: (a) A single-phase frequency divider reported by Oguey and Vittoz in 1973 and (b) a latch filed for patent by Piguet in 1974.



One can put some logic functions into N- or P- unit.

FIGURE 5: Yuan's original drawing of a TSPC circuit.

FIGURE 6: (a) A dynamic latch with a single clocked device, (b) cascaded TSPC stages, and (c) a three-stage master-slave flip-flop operating as a frequency divider.



FIGURE 9: The problem of race in TSPC stages.

## References

[1] J. Yuan and C. Svensson, "High-speed CMOS circuit technique," *IEEE J. Solid-State Circuits*, vol. 24, pp. 62–70, Feb. 1989.

[2] Y. Suzuki, K. Odagawa, and T. Abe, "Clocked CMOS calculator circuitry," *IEEE J. Solid-State Circuits*, vol. 8, pp. 462–469, Dec. 1973.

[3] R. H. Krambeck, C. M. Lee, and H. S. Law, "High-speed compact circuits with CMOS," *IEEE J. Solid-State Circuits*, vol. 17, pp. 614–619, June 1982.

[4] M. Shoji, *CMOS Digital Circuit Technology.* Englewood Cliffs, NJ: Prentice Hall, 1988.

[5] N. P. Goncalves and H. J. de Man, "NORA: A racefree dynamic CMOS technique for pipelined logic structures," *IEEE J. Solid-State Circuits*, vol. 18, pp. 261–268, June 1983.

*1st paper on True single phase clock circuit*

*Domino paper*

*NORA*

---

## Popular Clocking Disciplines with FF & Latches

❑ Options for timing of sequential circuits:

» **Using Flip-Flops** (1-phase single-edge-tr.)

*Same behavior*

» **Using Transparent Latches** (2-phase clocking)

» **Using Pulsed Latches** (1-phase single-edge-tr.)
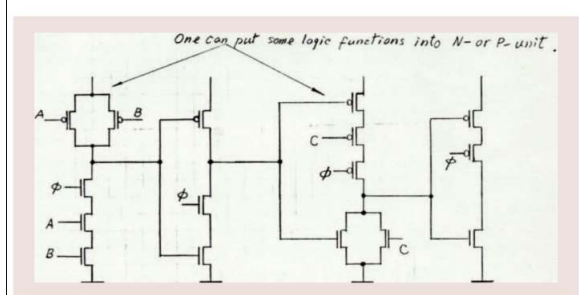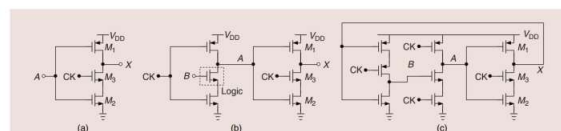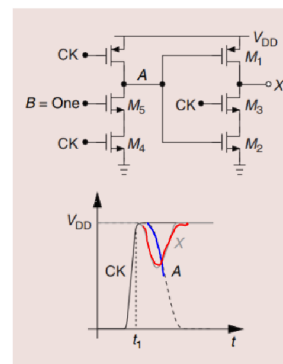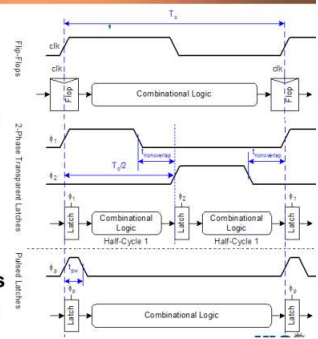


---

## Synchronous Clocking

Definitions:

▶ A (sub)circuit is said to operate synchronously iff all state transitions are restricted to occur periodically at precise moments of time that are determined by a special signal referred to as the clock.

▶ A clock domain is a (sub)circuit where all clock signals maintain fixed frequency and phase relationships because they are derived from a common source.



Toy example of synchronously clocked circuit

---

## Synchronous Clocking – Design Rules

**Design rule**

Consistently dissociate signals into

▶ asynchronous reset signals (when to enter the start state)
▶ clock signals (when to move from one state to the next)
▶ information signals (what state to enter, what output to produce)

**Design rule**

Make the clock period long enough so that all transient effects have died out before the next active clock edge instructs registers (and other storage devices) to accept new data!

---

## Delays - Abstraction



---

## Timing Parameters of Combinational & Sequential Circuits



Combinational circuit with three input signals and three output signals.

$t_{pd}$ (**Propagation delay**): The time required to process new input from applying a stable logic value at a data or clock input until the output has settled on its final value.

$t_{cd}$ (**Contamination delay**): The time from altering the logic value at a data or clock input until a the output value starts to change.

**Combinational and sequential circuits**

*Propagation delay tpd*
- New stable input (data or clock) – output settled on final value
- Example NAND: A-to-Z, and B-to-Z
- Example latch: D-to-Q, and/or CLK-to-Q

*Contamination delay (or retain delay) tcd*
- Altering input (data or clock) – first change of value at output
- By definition: 0≤tcd ≤tpd

---

## Synchronous Timing



Synchronous timing:
All registers synchronized with same CLK

---

## Timing Constraints (Setup & Hold)

- There are two main problems that can possibly arise in synchronous logic:
  - » Max Delay: The data doesn't have enough time to pass from one register to the next before the next clock edge.
  - » Min Delay: The data path is so short that it causes a hold violation in capturing register.
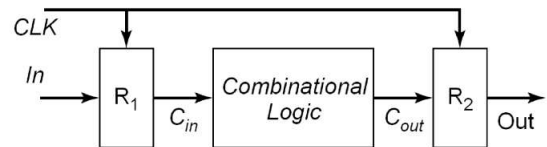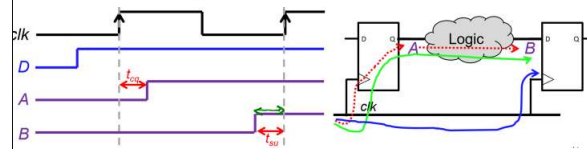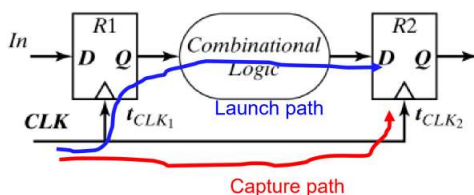
- Max delay violations are a result of a slow (long) data paths, including the register's $t_{su}$, therefore they are often called "Setup violations".

- Min delay violations are a result of a fast (short) data path, causing the data to change before the $t_{hold}$ of the reg has passed, therefore they are often called the "Hold violations".

---

## Setup (Max) Constraint

- Let's see what makes up our clock cycle:
  - » After the clock rises, it takes $t_{cq}$ for the data to propagate to point *A*.
  - » Then the data goes through the delay of the logic to get to point *B*.
  - » The data has to arrive at point *B* $t_{su}$ before the next clock edge.
- In general, our timing path is a race:
  - » Between the Data Arrival, starting with the *launching* clock edge.
  - » And the Data Capture, one clock period later.



---

## Setup (Max) Constraint



Launch path

Capture path

$$T > t_{CQ} + t_{\text{logic}} + t_{SU}$$

---

## Hold (Min) Constraint

- Hold problems occur due to the logic changing before $t_{hold}$ has passed.
- This is not a function of cycle time – it is relative to a single clock edge!
- Example of meeting the hold constraint:
  - » The clock rises and the data at A changes after $t_{cq}$.
  - » The data at B changes $t_{pd}(logic)$ later.
  - » Since the data at B had to stay stable for $t_{hold}$ after the clock (for the second register), the change at B has to be at least $t_{hold}$ **after** the clock edge.



B changes after $t_{hold}$, no hold violation!

## Hold (Min) Constraint



$$t_{CQ} + t_{\text{logic}} > t_{hold}$$

## Summary

- For *Setup* constraints, the clock period has to be longer than the data path delay:
  - » This sets our maximum frequency. $\boxed{T > t_{CQ} + t_{\text{logic}} + t_{SU}}$
  - » If we have setup failures, we can always just **slow down the clock**.

- For *Hold* constrains, the data path delay has to be longer than the hold time:
  - » This is **independent of clock period**. $\boxed{t_{CQ} + t_{\text{logic}} > t_{hold}}$
  - » If there is a hold failure, your chip will never work!
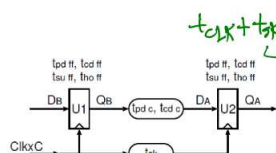
## Fundamental Timing Conditions – Reminder



The formulas to the right apply to the data input of register U2.

**Setup Condition**
$$t_{clk} \geqslant t_{pd_{ff}} + t_{pd_c} + t_{su_{ff}}(-t_{sk})$$

**Hold Condition**
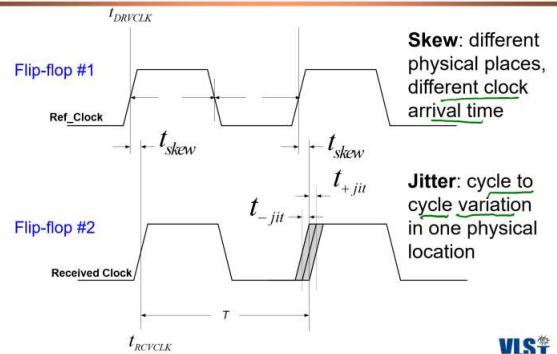$$t_{ho_{ff}} \leqslant t_{cd_{ff}} + t_{cd_c}(-t_{sk})$$

**Clock Skew**
Positive clock skew $t_{sk}$ relaxes the setup condition, while having detrimental effects on the hold condition.

## Clock Skew and Jitter



**Skew**: different physical places, different clock arrival time

**Jitter**: cycle to cycle variation in one physical location

- **Clock skew**
  - » Spatial variation in temporally equivalent clock edges; deterministic (can control it) + random, $t_{SK}$

- **Clock jitter**
  - » Temporal variations in consecutive edges of the clock signal; purely random
  - » Cycle-to-cycle (short-term) $t_{JS}$
  - » Long term $t_{JL}$

Both skew and jitter affect the effective cycle time

- **Variation of the pulse width**
  - » Important for level-sensitive clocking (latches)

## Sources of Clock Uncertainties



*Sources of clock uncertainty*
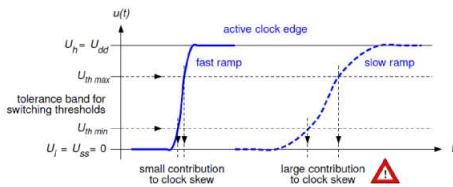
**Sources of skew include:**
- Different buffers & wires between source and leafs
- Temperature gradient

**Sources of jitter include:**
- VDD variations/noise (data-dependent, varies from cycle to cycle)
- Capacitive coupling from nearby signals
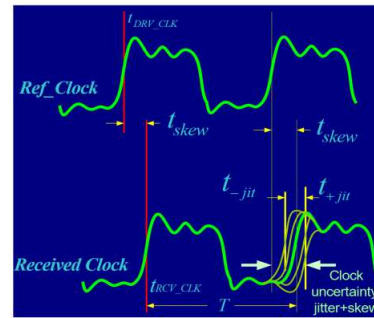
## Sources of Clock Uncertainties



- Disparities of switching thresholds across clocked cells translate ramp times into skew.
- Bistables get characterized for fast clock ramps. Setup and hold times tend to grow when a cell is being clocked with slow ramps.
- Correct behavior and accurate timing are put at risk.
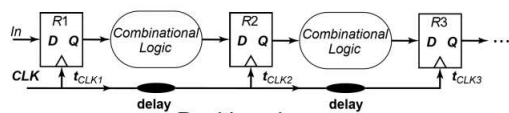  **Avoid slow transitions on clock nets**

VLSI

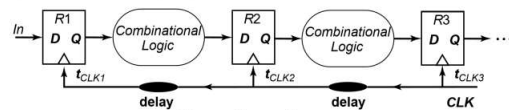## Clock Non-Idealities – Waveform Examples



Digital System Clocking: Oklobdzija, Stojanovic, Markovic, Nedovic VLSI

## Positive and Negative Skew



Positive skew
Data & clock propagate in the same direction. Relaxed setup constraints (longer effective clock cycle).

Negative skew
Opposite directions. Shorter effective clock cycle.

VLSI

## Setup (Max) Constraint

- The *Launch path* (still) consists of:
  - » $t_{cq} + t_{logic} + t_{setup}$    $t_{launch} = t_{CQ} + t_{logic} + t_{SU} + t_{jitter}$
  - » But if jitter makes the launch clock later, we need to add it to the data path delay.
- The *Capture path* consists of:
  - » The clock period ($T$)    $t_{capture} = T + \delta - t_{jitter}$
  - » Positive skew means the capture clock path is longer.
  - » If jitter makes the capture clock earlier, we need to subtract it.
- Our max constraint is:    $t_{capture} \geq t_{launch}$
- So we get:

$$+\delta \text{ Neg. skew}$$
$$T \geq t_{CQ} + t_{logic} + t_{SU} + 2t_{jitter} - \delta \text{ Pos. skew}$$

Jitter: always worse T    VLSI

## Setup (Max) Constraint

- Data has to arrive before *next* active clock edge.



$$T + \delta > t_{CQ} + t_{logic} + t_{SU} + 2t_{jitter}$$

## Hold (Min) Constraint

- The *Launch path* (still) consists of:
  - » $t_{cq} + t_{logic}$    $t_{launch} = t_{CQ} + t_{logic} - t_{jitter}$
  - » But if jitter makes the launch clock *earlier*, we need to subtract it from the data path delay.
- The *Capture path* consists of:    $t_{capture} = \delta + t_{jitter} + t_{hold}$
  - » Skew that makes the clock edge arrive at the capture register later than at the launch register.
  - » Actually, since it is a single clock edge, jitter should effect the capture clock the same as the launch clock.
  - » But as a worst case, we will add it as spacial jitter. Clock non-idealities: More difficult to meet hold constraint, need longer contamination delay
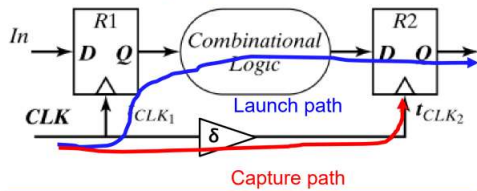- Our min constraint is:    $t_{launch} \geq t_{capture}$
- So we get:

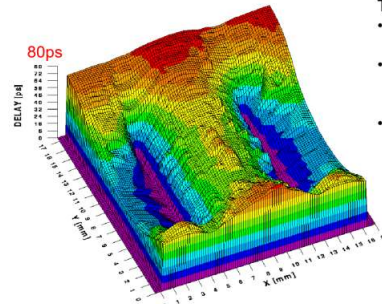$$t_{CQ} + t_{logic} \geq \delta + t_{hold} + 2t_{jitter}$$

VLSI

8

## Hold (Min) Constraint

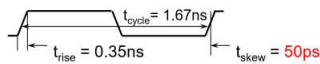❑ Data has to arrive *after* same clock edge has arrived at capturing register



Launch path

Capture path

$$t_{CQ} + t_{\text{logic}} > t_{hold} + \delta + 2t_{jitter}$$

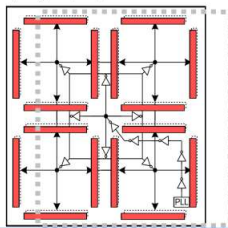---

## Clock Skew in Alpha Processor



80ps

Terrible clock skew
• Zero skew by clock drivers
• Data path in middle has intermediate skew
• Peripherals for I/O has high skew

---

## EV6 (Alpha 21264) Clocking

$t_{cycle}$ = 1.67ns

$t_{rise}$ = 0.35ns    $t_{skew}$ = 50ps
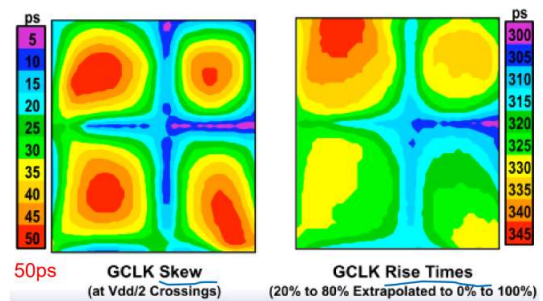
**Global clock waveform** (target spec)



Next version of Alpha processor had less skew: target >50ps

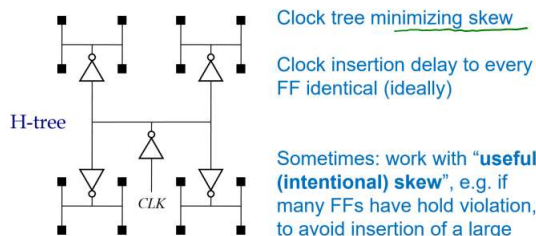Insert clock from PLL into center, then H-tree, and more local final clock buffers

Also, implemented more clock gating for power reduction

---

## EV6 Clock Skew



50ps

**GCLK Skew** (at Vdd/2 Crossings)

**GCLK Rise Times** (20% to 80% Extrapolated to 0% to 100%)

---

## Clock Distribution – H-Tree



H-tree

CLK

Clock tree minimizing skew

Clock insertion delay to every FF identical (ideally)

Sometimes: work with "**useful (intentional) skew**", e.g. if many FFs have hold violation, to avoid insertion of a large amount of buffers

Equal wire length/number of buffers to get to every location

---

## More realistic H-tree



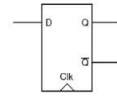True H-tree only works for regular structures like FPGAs

But can still have same delay on all paths from source to end points

[Restle98]
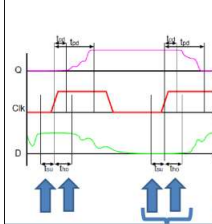
---

9

## Dealing with Skew and Jitter

- Balance clock paths using regular distribution network, such as H-tree
- Use local clock GRIDS (increased cap and power)
- Route data and clock in opposite directions to improve hold at the expense of setup.
- Shield clock wires (minimize capacitive coupling)
- Use dummy metal density fillers for regular wires
- Use decoupling capacitors (for stable VDD)
- Time borrowing (or cycle stealing): long path borrows time from subsequent short path, accomplished using latches
- "Useful skew" to avoid expensive buffering for hold fix

---

## Additional Timing Parameters of Sequential Circuits

**Rising-edge-triggered D-type flip-flop**
The value at the input D becomes available at the output Q at the rising edge of the clock

$t_{su}$ **(Set-up time):** The lapse of time before the active clock edge during which an input is required to assume a fixed logic value of either 0 or 1 at the input of a clocked circuit.

$t_{ho}$ **(Hold time):** The lapse of time after the active clock edge during which data are required to remain logically unchanged at the input of a clocked circuit.

Data call window: tsu + tho

---

## Timing Quantities (Cont'd) – Seq. Circuits Only
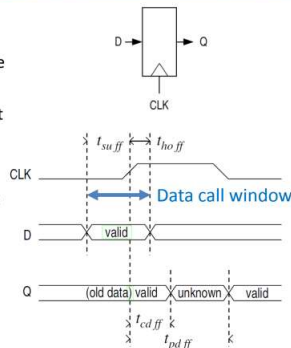
**Sequential circuits only**

*Setup time tsu*
- Time lapse before active clock edge during which input of bistable element (flip-flop, latch, SRAM, ROM, …) needs to settle for correct sampling
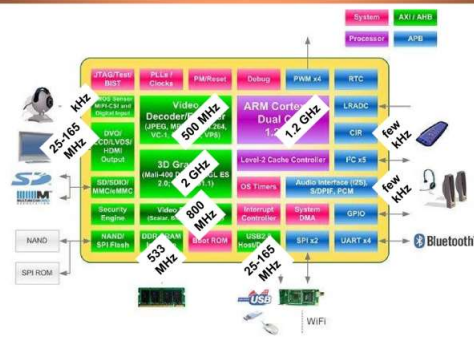
*Hold time tho:*
- Time lapse after active clock edge during which input data of bistable element needs to remain unchanged

***Data call window***
- Sum of *tsu* and *tho*
- Typically centered around active (positive) clock edge
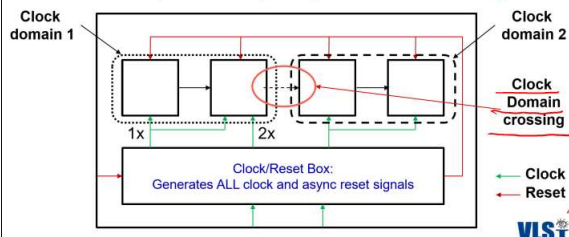
---

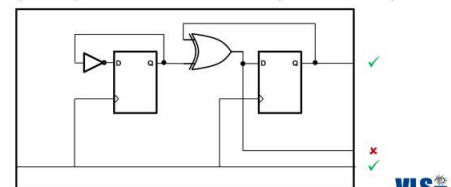## A Modern SoC has Many Clock Domains



---

## Designs with Multiple Clocks

- **Clock domain**: (Sub)circuits in which all clock signals maintain fixed frequency and phase relationship
  - » A single clock domain may use multiple divided clocks
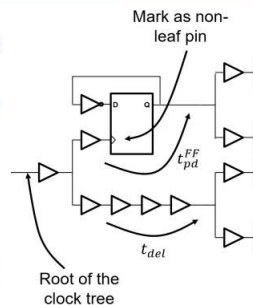- **Clock box**: top-level entity that generates all clock signals



---

## Clock Dividers

- Derive a slow clock from a fast clock by dividing it by an integer number (counter)
- Remember: all clocks (including divided clocks) need to be free of hzards and glitches
  - » Must assume that any logic can prodice glitches
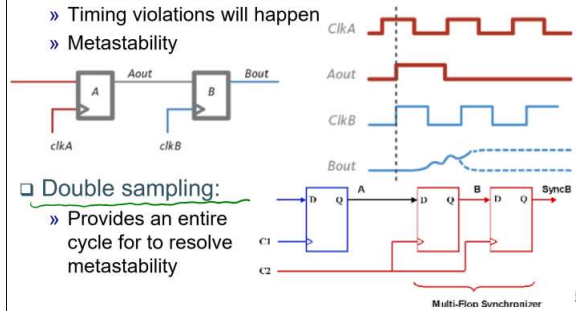  - » Generated (divided) clocks MUST come directly from a FF output

## Clock Dividers in the Clock Tree

□ Clock and derived clock must be phase aligned at the clock roots

□ Clock divider has a delay (FF propagation delay)
» Need to consider clock divider during clock-tree generation

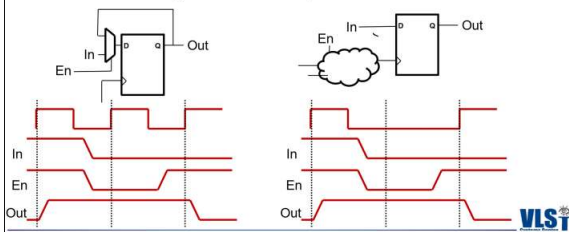□ Modern clock tree synthesis tools can trace through clock dividers

Mark as non-leaf pin

$t_{pd}^{FF}$

$t_{del}$

Root of the clock tree

## Crossing Clock Domains

□ Arbitrary phase relationship
» Timing violations will happen
» Metastability

ClkA
Aout
ClkB
Bout

□ Double sampling:
» Provides an entire cycle for to resolve metastability
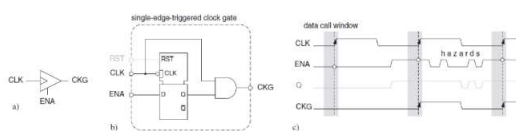
Multi-Flop Synchronizer

## Clock Gating

□ Flip-Flops with ENABLE: prevent FlipFlop(s) from capturing new data
» Functional: preserve content over multiple cycles
» Power savings: avoid activity in a block that is inactive

In
En
Out

## Clock Gating

□ Advantages:
» Saves a multiplexer per FlipFlop with Enable (Area, Delay, and Power advantage)
» Avoids activity on the clock net leading to the FlipFlop
» Avoids activity on the FlipFlops clock pin reducing internal power consumption

□ Disadvantages:
» Need for additional logic to suppress the clock while FlipFlop is disabled (area and power penalty)
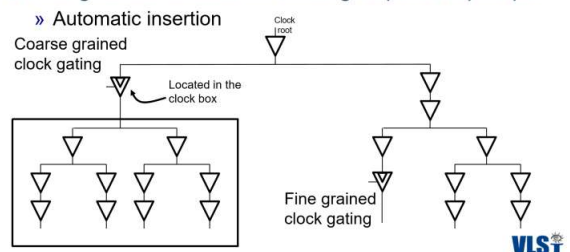» Need to ensure that the clock signal is free of glitches

□ A safe strategy to realize clock gating

single-edge-triggered clock gate

data call window

» Latch on the Enable signal shields glitches during the sensitive period of the clock
» Implemented with individual cells: some timing constraints need to be observed
» Often realized as single dedicated library cell

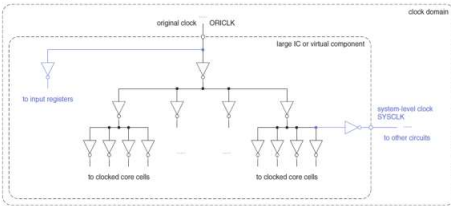## Clock Gating

□ Coarse grained: disable entire blocks of the design
» Inserted manually

□ Fine grained: enable for small groups of FlipFlops
» Automatic insertion

Clock root

Coarse grained clock gating

Located in the clock box

Fine grained clock gating

11

## Improving IO Timing

❑ Chip provides a clock output that is aligned with the clock at the leaves of the FlipFlops
» Output of the clock output pad is declared as clock leaf
» Delayed clock is used as a reference for rest of system



## Improving IO Timing

❑ Delay locked loop ( DLL )
» Generates a phase shifted clock such that the reference input is phase aligned with the input clock
» Clock reference is taken from the leaves of the clock tree

❑ Internal clock at the leaves is aligned with clock input