

# ESE532: System-on-a-Chip Architecture

Day 7: September 23, 2019  
Pipelining



Penn ESE532 Fall 2019 -- DeHon

## Previously

- Pipelining in the large
  - Not just for gate-level circuits
- Throughput and Latency
- Form of parallelism

2

Penn ESE532 Fall 2019 -- DeHon

## Today

Pipelining details (for gates, primitive ops)

- Systematic Approach
- Justify Operator and Interconnect Pipelining
- Loop Bodies
- Cycles in the Dataflow Graph
- C-slow [bonus if have time]

Penn ESE532 Fall 2019 -- DeHon

3

## Message

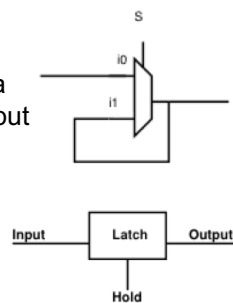
- Pipelining is an efficient way to reuse hardware to perform the **same** set of operations at high throughput

Penn ESE532 Fall 2019 -- DeHon

4

## Latch

- Element that can hold a previous value of an input

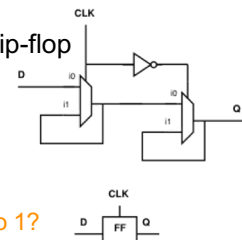


ESE150 Spring 2019

5

## Register

- Use a pair to create a flip-flop
  - Also call register
- What happens when
  - CLK is low (0) ?
  - CLK is high (1) ?
  - CLK transitions from 0 to 1?

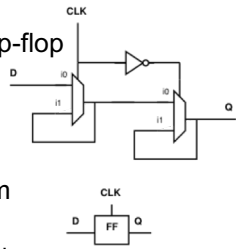


ESE150 Spring 2019

6

## Register

- Use a pair to create a flip-flop
  - Also call register
- Sample D input on 0→1 transition of clock (CLK)
- Never an open path from D→Q
  - One of the mux latches always in hold state

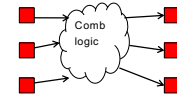


ESE150 Spring 2019

7

## Synchronous Circuit Discipline

- Registers that sample inputs at clock edge and hold value throughout clock period
- Compute from registers-to-registers
- Cycle time large enough for longest logic path between registers
- Min cycle = Max path delay between registers

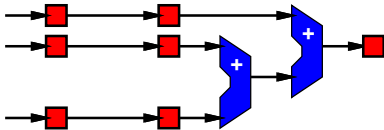


Penn ESE532 Fall 2019 -- DeHon

8

## Preclass 1

- Delay between registers as shown?

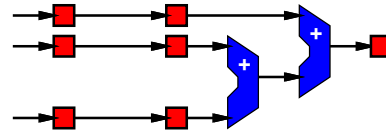


Penn ESE532 Fall 2019 -- DeHon

9

## Preclass 1

- Move registers so can clock at adder delay?



Penn ESE532 Fall 2019 -- DeHon

10

## Pipeline Reuse

- Lower delay between clocks
  - Higher clock rate
  - Higher potential throughput
  - Faster we reuse our logic
  - More capacity get out of design
    - Assuming registers cheap in area and time overhead
      - $T_{\text{setup}}, T_{\text{clk} \rightarrow \text{q}} \sim 20\text{ps}$ ,  $T_{\text{add}} \sim 500\text{ps}$
      - Registers  $\sim 10$  transistors/bit
      - Adder  $\sim 40\text{--}50$  transistors/bit

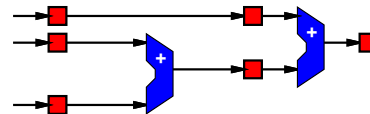


Penn ESE532 Fall 2019 -- DeHon

11

## Note Registers on Links

- Some links end up with multiple registers.
- Why?

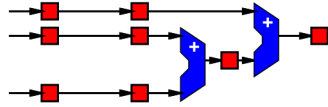


Penn ESE532 Fall 2019 -- DeHon

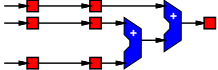
12

## Preclass 2: What Happens?

- What would be wrong with this pipelining?



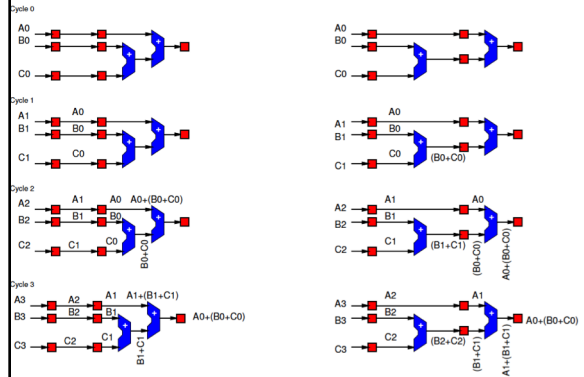
- For this initial design:



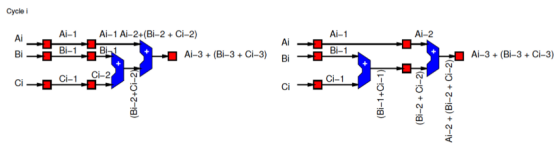
Penn ESE532 Fall 2019 -- DeHon

13

## Behavior



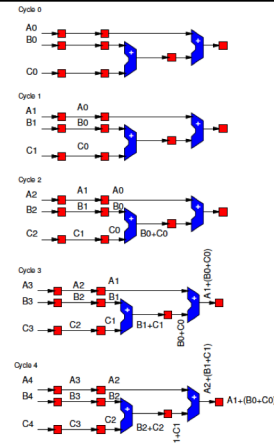
## Equations



Penn ESE532 Fall 2019 -- DeHon

15

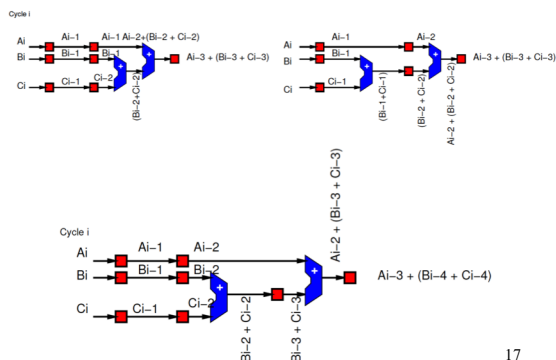
## Behavior Add Register



Penn ESE532 Fall 2019 -- DeHon

16

## Equations



Penn ESE532 Fall 2019 -- DeHon

17

## Consistent Pipelining

- Makes sure a consistent input set arrives at each gate/operator
  - Don't get mixing between input sets

Penn ESE532 Fall 2019 -- DeHon

18

## Legal Register Moves

- Retiming Lag/Lead

- Lag: remove register every input  
add register every output
- Lead: remove register every output  
add register every input

Penn ESE532 Fall 2019 -- DeHon 19

## Preclass 1

- Retime using Lead/Lag

Penn ESE532 Fall 2019 -- DeHon 20

## Preclass 1 (revisited)

Penn ESE532 Fall 2019 -- DeHon 21

## Add Registers and Move

- If we're willing to add pipeline delay
  - Add any number of pipeline registers at input
  - Move registers into circuit to reduce cycle time
    - Reduce max delay between registers

Penn ESE532 Fall 2019 -- DeHon 22

## Add Registers at Input

Penn ESE532 Fall 2019 -- DeHon 23

## Add Register and Retime

- Add chain of registers on every input
- Retime registers into circuit
  - Minimizing delay between registers

Penn ESE532 Fall 2019 -- DeHon 24

## Add Registers and Retime

- Lets us think about behavior
  - What the pipelining is doing to cycles of delay
- Separate from details of how redistribute registers
- Behavioral equivalence between the registers-at-front and properly retimed version of circuit

Penn ESE532 Fall 2019 -- DeHon

25

## Justify Pipelining

(or composing pipelined operators)

Penn ESE532 Fall 2019 -- DeHon

26

## Handling Pipelined Operators

- Given a pipelined operator
  - (or a pipelined interconnect)
- Discipline of picking a frequency target and designing everything for that
  - May be necessary to pipeline operator since its delay is too high
- Due to hierarchy
  - Pipelined this operator and now want to use it as a building block

Penn ESE532 Fall 2019 -- DeHon

27

## Examples

- Run at 500MHz
- Floating-point unit that takes 9ns
  - Can pipeline into 5, 2ns stages
- Multiplier that takes 6ns
- Memory can access in 2ns
  - Only if registers on address/inputs and output
  - i.e. exist in own clock stage

Penn ESE532 Fall 2019 -- DeHon

28

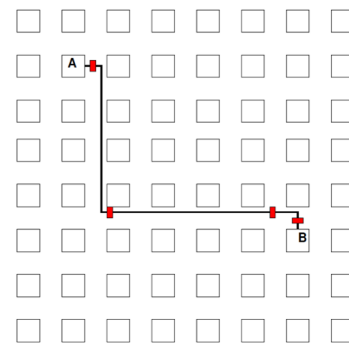
## Interconnect Delay

- Chips  $\gg$  Clock Cycles
- May have chip 100s of Operators wide
- May only be able to reach across 10 operators in a 2ns cycle
- Must pipeline long interconnect links

Penn ESE532 Fall 2019 -- DeHon

29

## Interconnect Example



Penn ESE532 Fall 2019 -- DeHon

30

## Methodology: Pipelined Operator Graph

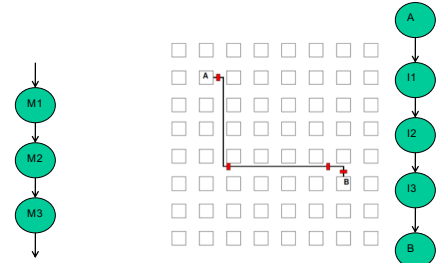
- Start with logical, unpipelined graph
- Treat each pipelined operator as a set of unit-delay operators of mandatory depth
- Treat each interconnect pipeline stage as a unit-delay buffer
- Add registers at input
- Retime into graph

Penn ESE532 Fall 2019 -- DeHon

31

## Model

- 3-stage Multiplier
- Interconnect Delay



Penn ESE532 Fall 2019 -- DeHon

32

## Pipeline Loop

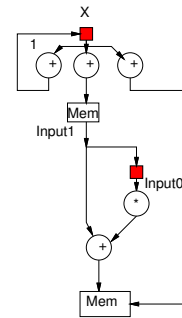
(and use for justify pipeline example)

Penn ESE532 Fall 2019 -- DeHon

33

## Preclass 4

- Logical (unpipelined) dataflow graph for loop body



Penn ESE532 Fall 2019 -- DeHon

34

## Example Operators

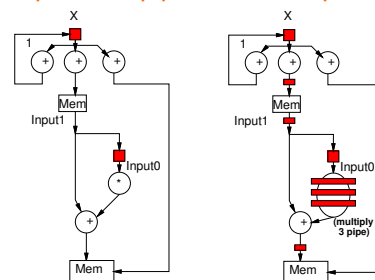
- Operator and Interconnect delays
  - Multiplier 3 cycles
  - Reading from Input array
    - Memory op is cycle after computing address
    - Takes one cycle delay bring data back to multiplier

Penn ESE532 Fall 2019 -- DeHon

35

## Illustrate Need

- What happens if just use graph as is (with operators pipelined as required)?



Penn ESE532 Fall 2019 -- DeHon

36

### Model Graph

- Revised graph for modeling

Penn ESE532 Fall 2019 -- DeHon

37

### Pipeline Graph

- Result after pipelining?

Penn ESE532 Fall 2019 -- DeHon

38

### Pipelining Lesson

- Can always pipeline an **acyclic** graph to fixed frequency target
  - fixed pipelining of primitive operators
  - Pipeline interconnect delays
- Need to keep track of registers to balance paths
  - So see consistent delays to operators

Penn ESE532 Fall 2019 -- DeHon

39

### Graph Cycles

Watch: Clock cycle  
Cycle time  
Cycle in Graph

Penn ESE532 Fall 2019 -- DeHon

40

### Preclass 3

- Can we retime to reduce cycle time

Penn ESE532 Fall 2019 -- DeHon

41

### Retiming Limits?

- What prevents us from retiming?

Penn ESE532 Fall 2019 -- DeHon

42

## (Graph) Cycle Observation

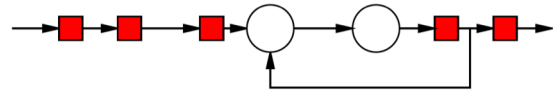
- Retiming does not allow us to change the *number of registers inside a cycle*.
- Limit to cycle time
  - Max delay in cycle / Registers in cycle
- Pipelining doesn't help inside cycle
  - Cannot push registers into cycle

Penn ESE532 Fall 2019 -- DeHon

43

## Simple Cycle

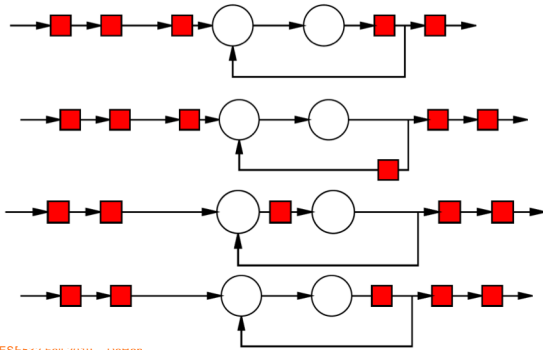
- Delay of cycle?
- Registers in cycle?
- What happens to cycle if try to apply lead/lag?



Penn ESE532 Fall 2019 -- DeHon

44

## Retiming



Penn ESE532 Fall 2019 -- DeHon

## Loop

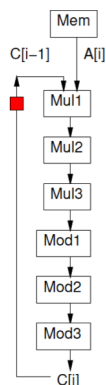
- Consider
  - [multiply and mod each take 3 cycles]
- For  $(i=0; i<N; i++)$   
 $C[i] = (C[i-1] * A[i]) \% N;$

Penn ESE532 Fall 2019 -- DeHon

46

## Loop

- For  $(i=0; i<N; i++)$   
 $C[i] = (C[i-1] * A[i]) \% N;$

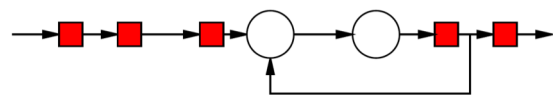


Penn ESE532 Fall 2019 -- DeHon

47

## Initiation Interval (II)

- Cyclic dependencies can limit throughput
- Due to data dependent cycles in graph,
  - May not be able to initiate a new computation on every clock cycle
- II – clock cycles (delay) before can initiate
- Throughput =  $1/II$



Penn ESE532 Fall 2019 -- DeHon

48



## Loop

- For (i=0;i<N;i++)  
 $C[i]=(C[i-1]*A[i])\%N;$
- Initiation Interval?

49

Penn ESE532 Fall 2019 -- DeHon

## Initial Interval

- Delay Around graph cycle?  
 – Assume multiply 3, add 1
- Registers in graph cycle?
- Retiming clock cycle bound = II ?
- Achievable?

50

Penn ESE532 Fall 2019 -- DeHon

## II and Latency

- Actually is a cycle  
 – II?  
 – Latency?

51

Penn ESE532 Fall 2019 -- DeHon

## II and Latency

- II? (assume willing to pipeline inputs)
- Latency?

52

Penn ESE532 Fall 2019 -- DeHon

## C-Slow

(if time permits)

53

Penn ESE532 Fall 2019 -- DeHon

## Problem

- Pipelining cannot push registers into cycle
- Graph cycles can prevent running at full pipeline target (maximum frequency)
- If not reusing operators at full pipeline target are underutilizing resources
- Can we use the resources for something?

54

Penn ESE532 Fall 2019 -- DeHon

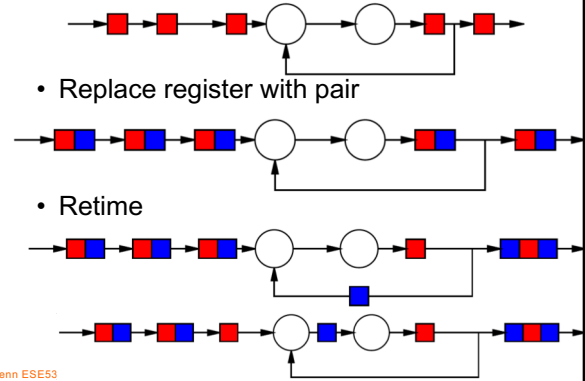
## C-Slow

- **Observation:** if we have data-level parallelism, can use to solve independent problems on same hardware
- **Transformation:** make C copies of each register
- **Guarantee:** C computations operate independently
  - Do not interact with each other

Penn ESE532 Fall 2019 -- DeHon

55

## 2-Slow Simple Cycle



Penn ESE53

## 2-Slow Simple Cycle

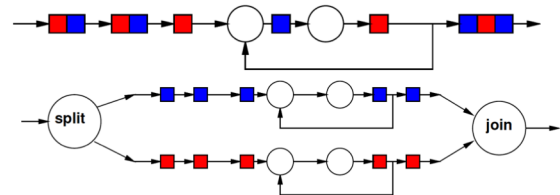
- 
- Replace register with pair
  - Retime
  - Observe independence of red/blue computations

Penn ESE532 Fall 2019 -- DeHon

57

## Equivalence

- The 2-slow operator is equivalent to two data parallel operators running at half the speed



Penn ESE532 Fall 2019 -- DeHon

58

## Automation

- No mainstream tool today will perform C-slow transformation for you automatically
- Synthesis tools will retime registers

Penn ESE532 Fall 2019 -- DeHon

59

## Lesson

- Cyclic dependencies limit throughput on single task or data stream
  - Cycle-length / registers-in-cycle
- Can use on C independent (data parallel) tasks

Penn ESE532 Fall 2019 -- DeHon

60

## Big Ideas

- Pipeline computations to reuse hardware and maximize computational capacity
- Can compose pipelined operators and accommodate fixed-frequency target
  - Be careful with data retiming
- Cycles limit pipelining on single stream
- C-slow to share hardware among multiple, data-parallel streams

Penn ESE532 Fall 2019 -- DeHon

61

## Admin

- Reading for Day 8 on web
- HW4 due Friday
  - HW4 lighter, HW5 heavier
  - Suggest planning to get early start on HW5

Penn ESE532 Fall 2019 -- DeHon

62