

ESE5320: System-on-a-Chip Architecture

Day 7: September 23, 2024
Pipelining



Penn ESE5320 Fall 2024 -- DeHon

1

Previously

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

Penn ESE5320 Fall 2024 -- DeHon

2

2

Today

Pipelining details (for gates, primitive ops)

- Systematic Approach (Part 1)
- Justify Operator and Interconnect Pipelining (Part 2)
- Loop Bodies
- Cycles in the Dataflow Graph (Part 3)

Penn ESE5320 Fall 2024 -- DeHon

3

3

Message

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

Penn ESE5320 Fall 2024 -- DeHon

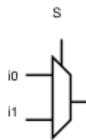
4

4

Multiplexer Gate

- MUX

- When $S=0$, output= i_0
- When $S=1$, output= i_1



S	i0	i1	Mux2(S,i0,i1)
0	0	0	0
0	0	1	0
0	1	0	1
0	1	1	1
1	0	0	0
1	0	1	1
1	1	0	0
1	1	1	1

Penn ESE5320 Fall 2024 -- DeHon

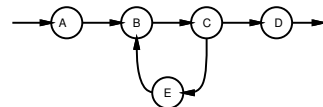
5

5

Cycle

Two uses of term in this lecture:

- Repetitive waveform
 - E.g. sine wave or square wave
- Graph cycle



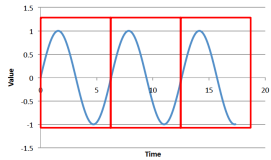
Penn ESE5320 Fall 2024 -- DeHon

6


6

Waveform Cycles

- How many cycles showing of sine wave?



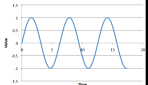
- How many cycles of square wave?




Penn ESE5320 Fall 2024 -- DeHon

7

Waveform Cycles



- How many cycles showing of sine wave?
- How many cycles of square wave?



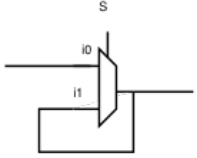
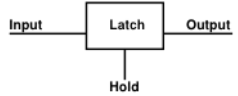
- Note: clock on which we pipeline is a square wave
 - Talk about what happens in a clock cycle
 - Talk about number of clock cycles

Penn ESE5320 Fall 2024 -- DeHon

8

Latch

- Element that can hold a previous value of an input

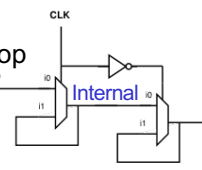
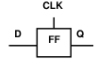



Penn ESE5320 Fall 2024 -- DeHon

9

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 ?
- What output Q until next 0 to 1 CLK transition?

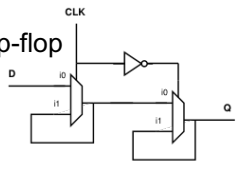
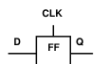
Penn ESE5320 Fall 2024 -- DeHon

10

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


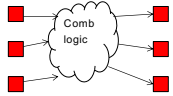



Penn ESE5320 Fall 2024 -- DeHon

11

Synchronous Circuit Discipline

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

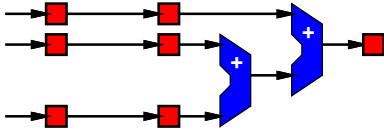



Penn ESE5320 Fall 2024 -- DeHon

12

Preclass 1

- Delay between registers as shown?



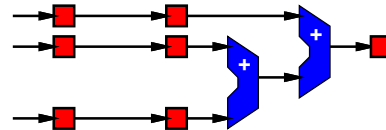
Penn ESE5320 Fall 2024 -- DeHon

13

13

Preclass 1

- Latency of pipeline?
 - Cycles
 - nanoseconds



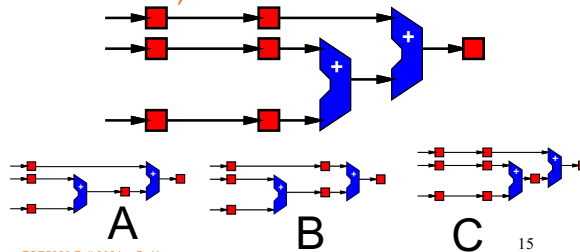
Penn ESE5320 Fall 2024 -- DeHon

14

14

Preclass 1

- Move registers so can clock at adder delay? (which is of choices at bottom is correct = same behavior)



Penn ESE5320 Fall 2024 -- DeHon

15

15

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 ~ 10 transistors/bit
 - Adder $\sim 40\text{--}50$ transistors/bit



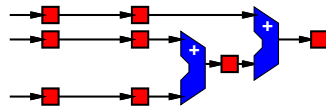
Penn ESE5320 Fall 2024 -- DeHon

16

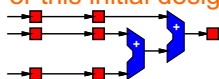
16

Preclass 2: What Happens?

- What would be wrong with this pipelining?



- For this initial design:

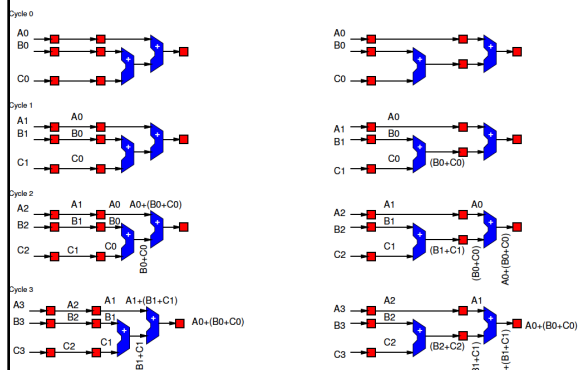


Penn ESE5320 Fall 2024 -- DeHon

17

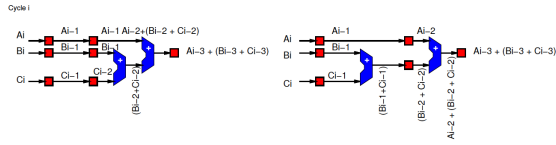
17

Behavior



18

Equations

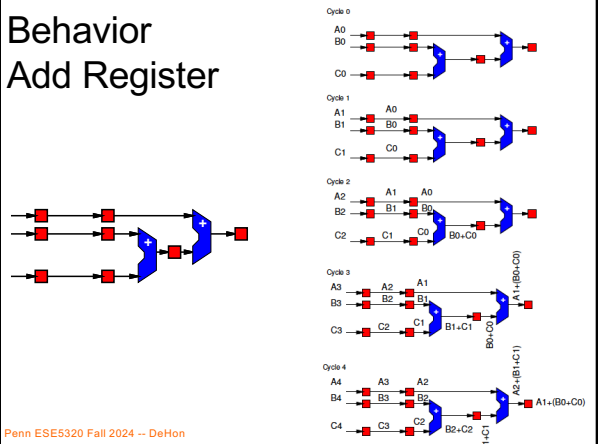


Penn ESE5320 Fall 2024 -- DeHon

19

19

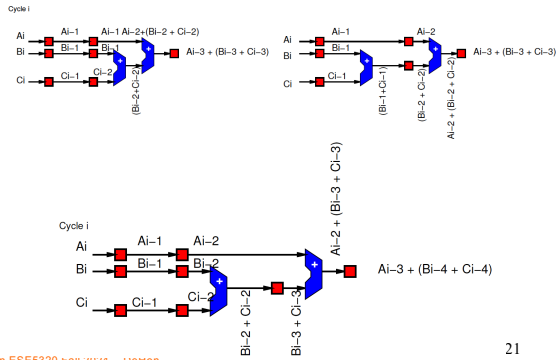
Behavior Add Register



Penn ESE5320 Fall 2024 -- DeHon

20

Equations



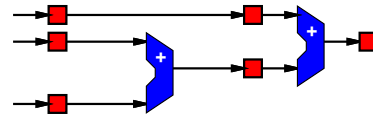
Penn ESE5320 Fall 2024 -- DeHon

21

21

Note Registers on Links

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



Penn ESE5320 Fall 2024 -- DeHon

22

22

Consistent Pipelining

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

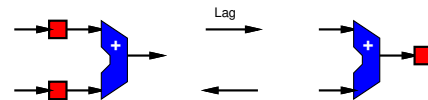
Penn ESE5320 Fall 2024 -- DeHon

23

23

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 ESE5320 Fall 2024 -- DeHon

24

24

Preclass 1

• Retime using Lead/Lag

Penn ESE5320 Fall 2024 -- DeHon

25

Preclass 1 (revisited)

Penn ESE5320 Fall 2024 -- DeHon

26

25

26

Rounup

• Result

- 2 clock cycle pipeline latency
- 1ns delay between registers
- 1GHz clock
- $2 * 1\text{ns} = 2\text{ns}$ total latency

Penn ESE5320 Fall 2024 -- DeHon

27

27

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 ESE5320 Fall 2024 -- DeHon

28

28

Add Registers at Input

Penn ESE5320 Fall 2024 -- DeHon

29

29

Add Registers at Input and Retime

Penn ESE5320 Fall 2024 -- DeHon

30

30

Add Registers at Input

Throughput and Latency?

Penn ESE5320 Fall 2024 -- DeHon

31

Add Register and Retime

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

Penn ESE5320 Fall 2024 -- DeHon

32

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 ESE5320 Fall 2024 -- DeHon

33

Justify Pipelining

(or composing pipelined operators)
Part 2

Penn ESE5320 Fall 2024 -- DeHon

34

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 ESE5320 Fall 2024 -- DeHon

35

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 ESE5320 Fall 2024 -- DeHon

36

Interconnect Delay

- Chips >> 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 ESE5320 Fall 2024 -- DeHon 37

37

Interconnect Example

Penn ESE5320 Fall 2024 -- DeHon 38

38

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 ESE5320 Fall 2024 -- DeHon 39

39

Model

- 3-stage Multiplier
- Interconnect Delay

Penn ESE5320 Fall 2024 -- DeHon 40

40

Pipeline Loop

(and use for justify pipeline example)

Penn ESE5320 Fall 2024 -- DeHon 41

41

Preclass 4

- Logical (unpipelined) dataflow graph for loop body

Penn ESE5320 Fall 2024 -- DeHon 42

42

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 (or adder)

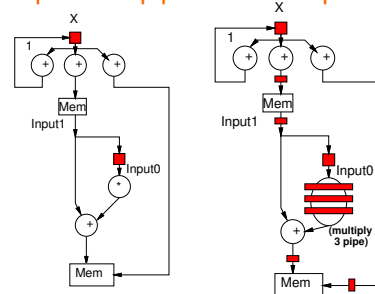
Penn ESE5320 Fall 2024 -- DeHon

43

43

Illustrate Need

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



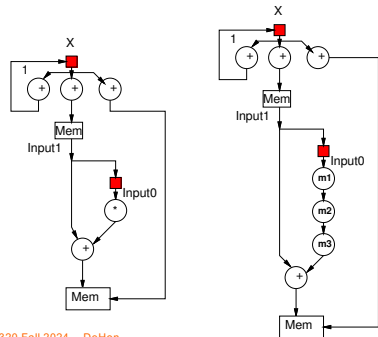
Penn ESE5320 Fall 2024 -- DeHon

44

44

Model Graph

- Revised graph for modeling



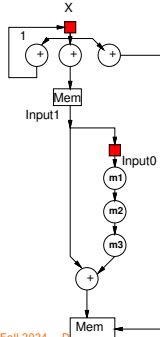
Penn ESE5320 Fall 2024 -- DeHon

45

45

Pipeline Graph

- Result after first retime (top register)



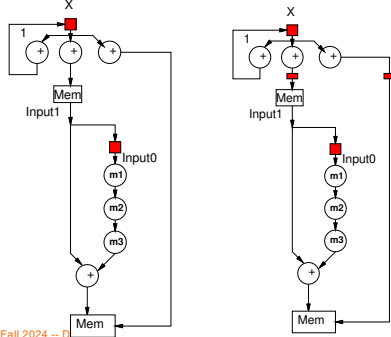
Penn ESE5320 Fall 2024 -- DeHon

46

46

Pipeline Graph

- Result after first retime



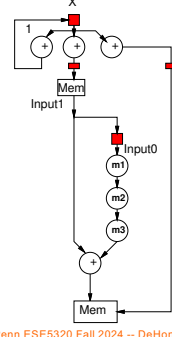
Penn ESE5320 Fall 2024 -- DeHon

47

47

Pipeline Graph

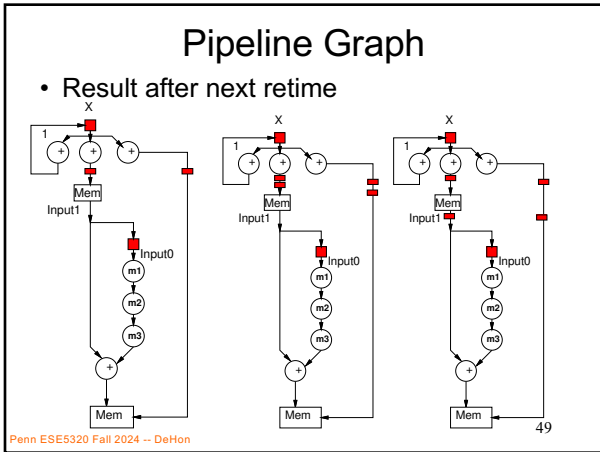
- Result after next retime (top register)?



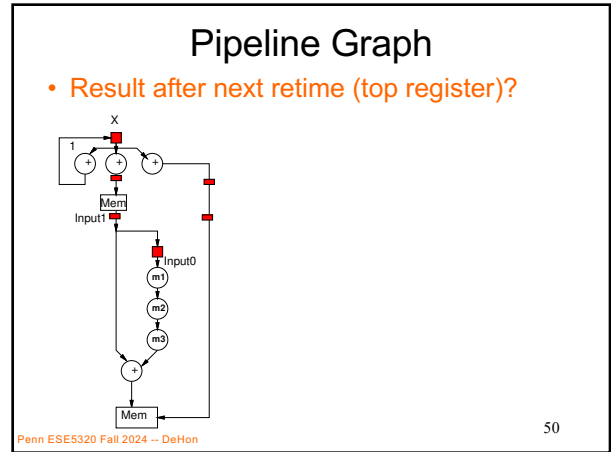
Penn ESE5320 Fall 2024 -- DeHon

48

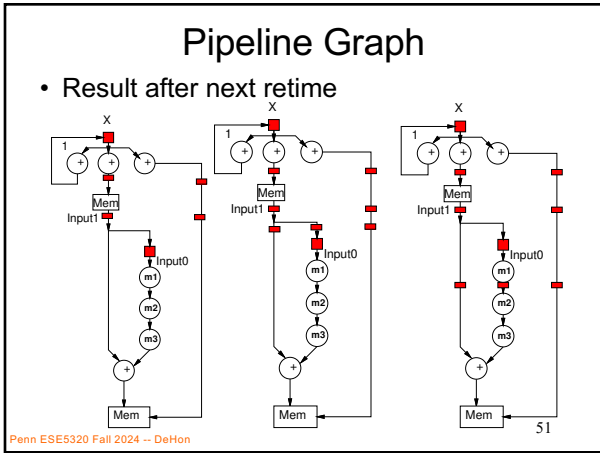
48



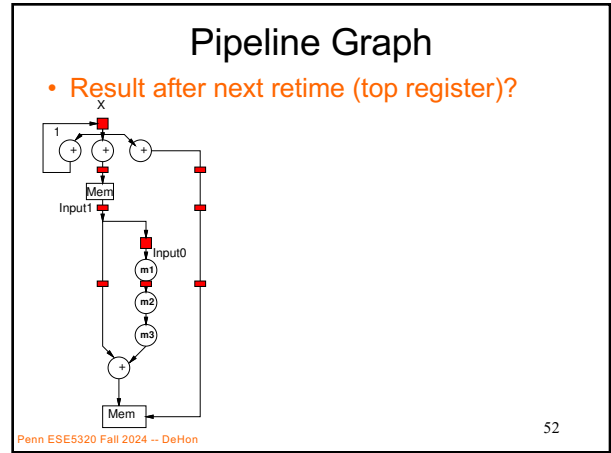
49



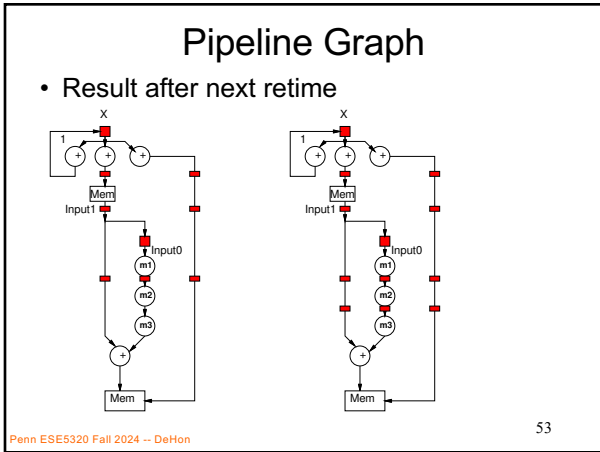
50



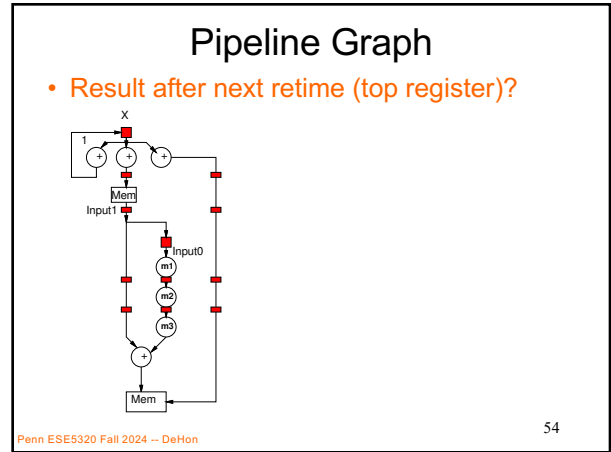
51



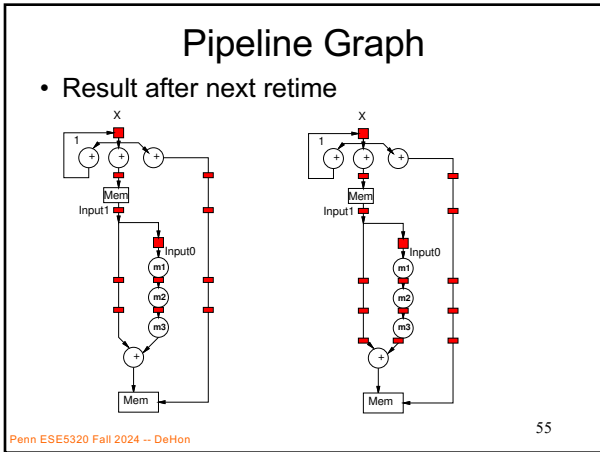
52



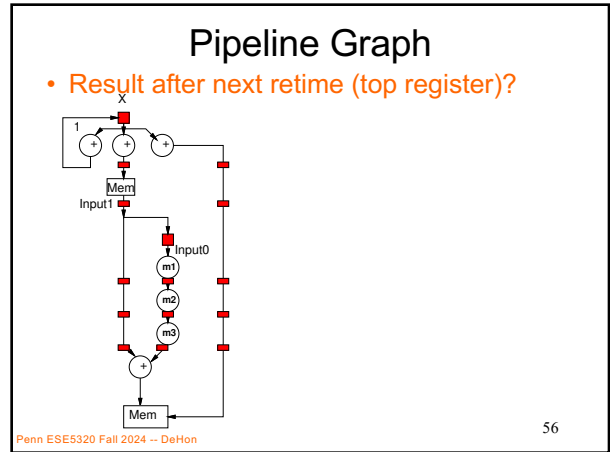
53



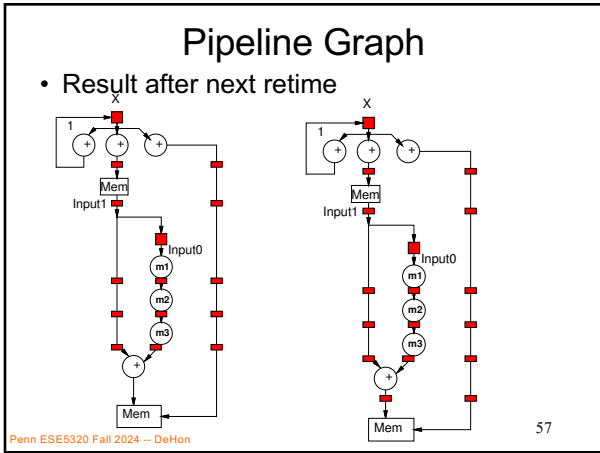
54



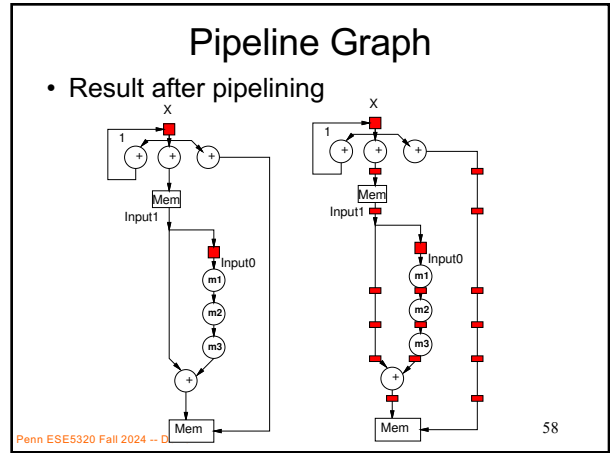
55



56



57



58

- ### Pipelining Lesson
- Can always pipeline an **acyclic graph** (no graph cycles) to fixed frequency target (fixed clock cycle period)
 - fixed pipelining of primitive operators
 - Pipeline interconnect delays
 - Need to keep track of registers to balance paths
 - So see consistent delays to operators
- Penn ESE5320 Fall 2024 -- DeHon

59

Graph Cycles

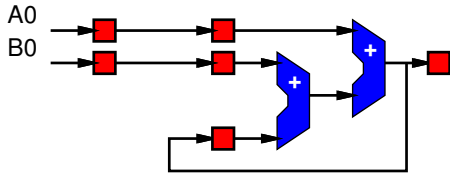
Watch: Clock cycle
Cycle time
Cycle in Graph
Part 3

Penn ESE5320 Fall 2024 -- DeHon

60

Preclass 3

- Can we retime to reduce clock cycle time?



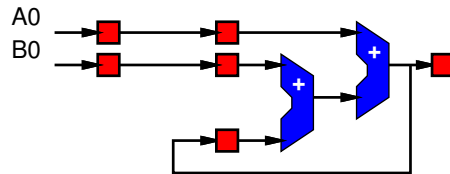
Penn ESE5320 Fall 2024 -- DeHon

61

61

Retiming Limits?

- What prevents us from retiming?

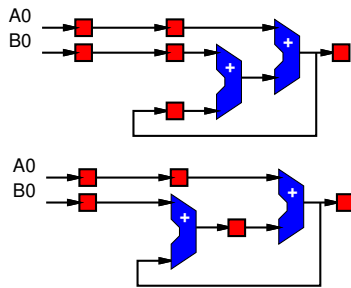


Penn ESE5320 Fall 2024 -- DeHon

62

62

Retiming Limits?



Penn ESE5320 Fall 2024 -- DeHon

63

63

(Graph) Cycle Observation

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

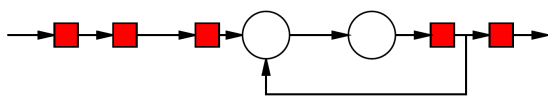
Penn ESE5320 Fall 2024 -- DeHon

64

64

Simple Graph Cycle

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

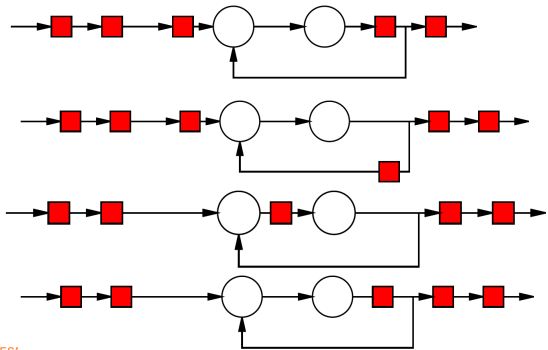


Penn ESE5320 Fall 2024 -- DeHon

65

65

Retiming

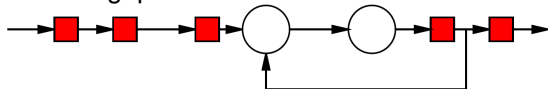


Penn ESE5320 Fall 2024 -- DeHon

66

Initiation Interval (II)

- Cyclic dependencies in a dataflow graph 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 ESE5320 Fall 2024 -- DeHon

67

67

Loop

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

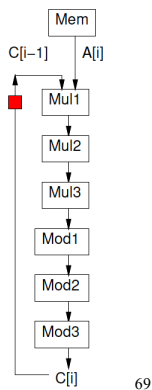
Penn ESE5320 Fall 2024 -- DeHon

68

68

Loop

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



Penn ESE5320 Fall 2024 -- DeHon

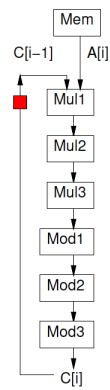
69

69

Loop

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

- Initiation Interval?



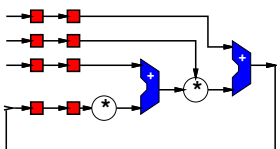
Penn ESE5320 Fall 2024 -- DeHon

70

70

Initiation Interval

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

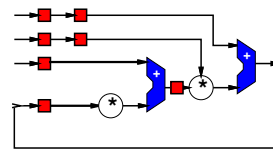
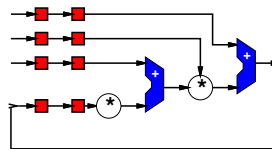


Penn ESE5320 Fall 2024 -- DeHon

71

71

Retimed



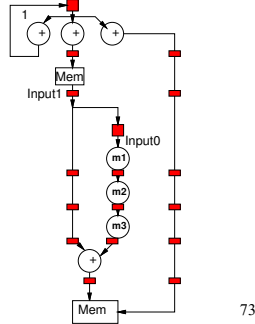
Penn ESE5320 Fall 2024 -- DeHon

72

72

II and Latency

- Actually is a graph cycle x
- II?
- Latency
- x to mem write?



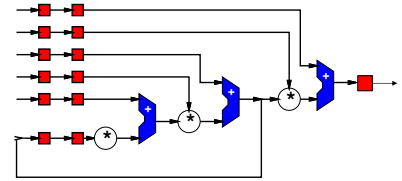
Penn ESE5320 Fall 2024 -- DeHon

73

73

II and Latency

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

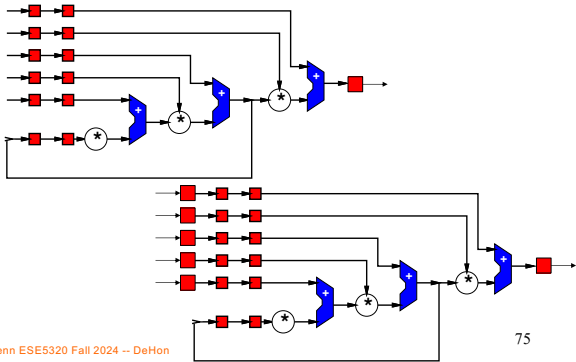


Penn ESE5320 Fall 2024 -- DeHon

74

74

II and Latency

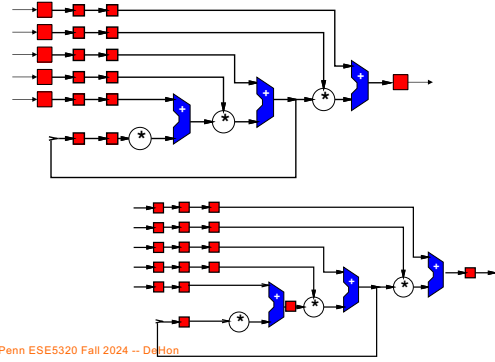


Penn ESE5320 Fall 2024 -- DeHon

75

75

II and Latency

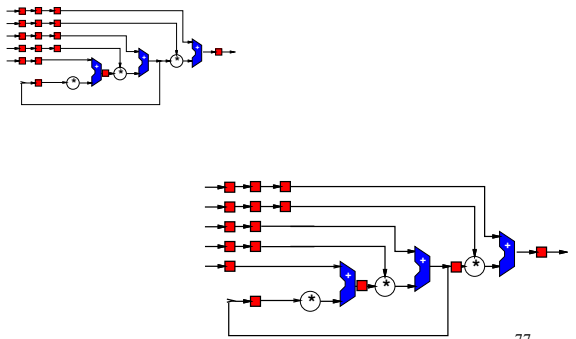


Penn ESE5320 Fall 2024 -- DeHon

76

76

II and Latency



Penn ESE5320 Fall 2024 -- DeHon

77

77

Lesson

- Graph Cyclic dependencies limit throughput on single task or data stream
 - Cycle Bound = Cycle-length / registers-in-cycle
 - (analog to latency bound)

Penn ESE5320 Fall 2024 -- DeHon

78

78

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
- Graph cycles limit pipelining on single stream – II (Initiation Interval)

Admin

- Remember Feedback form
 - Including HW3
- Reading for Day 8 on web
- HW4 due Friday