

ESE532: System-on-a-Chip Architecture

Day 7: September 28, 2020
Pipelining



Previously

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

Today

Pipelining details (for gates, primitive ops)

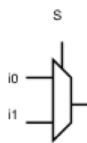
- Systematic Approach (Part 1)
- Justify Operator and Interconnect Pipelining (Part 2)
- Loop Bodies
- Cycles in the Dataflow Graph (Part 3)
- C-slow [probably separate record] (Part 4)

Message

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

Multiplexer Gate

- MUX
 - When $S=0$, output= $i0$
 - When $S=1$, output= $i1$

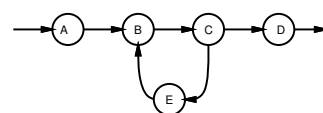


S	i0	i1	Mux2(S,i0,i1)
0	0	0	0
0	0	1	0
0	1	0	1
0	1	1	1
1	0	0	0
1	0	1	1
1	1	0	0
1	1	1	1

Cycle

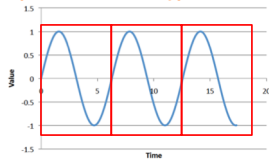
Two uses of term in this lecture:

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



Waveform Cycles

- How many cycles showing of sine wave?



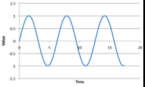
- How many cycles of square wave?



Penn ESE532 Fall 2020 -- 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

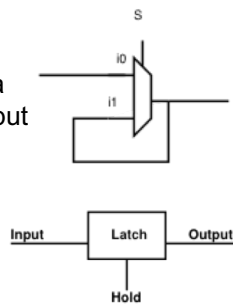
- Talk about what happens in a clock cycle
- Talk about number of clock cycles

Penn ESE532 Fall 2020 -- DeHon

8

Latch

- Element that can hold a previous value of an input

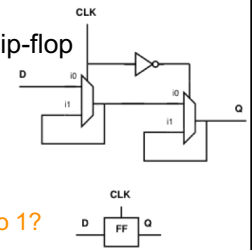


Penn ESE532 Fall 2020 -- DeHon

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 ?

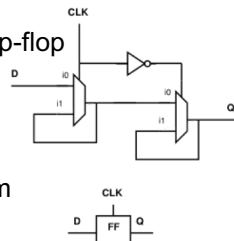


Penn ESE532 Fall 2020 -- DeHon

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

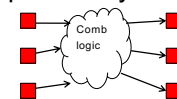


Penn ESE532 Fall 2020 -- DeHon

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

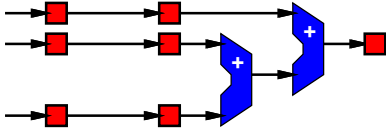


Penn ESE532 Fall 2020 -- DeHon

12

Preclass 1

- Delay between registers as shown?

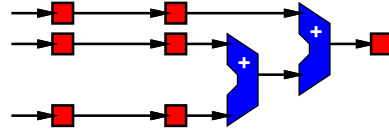


Penn ESE532 Fall 2020 -- DeHon

13

Preclass 1

- Move registers so can clock at adder delay? [show on Google Doc]



Penn ESE532 Fall 2020 -- DeHon

14

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

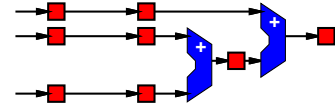


Penn ESE532 Fall 2020 -- DeHon

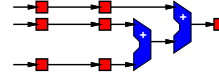
15

Preclass 2: What Happens?

- What would be wrong with this pipelining?



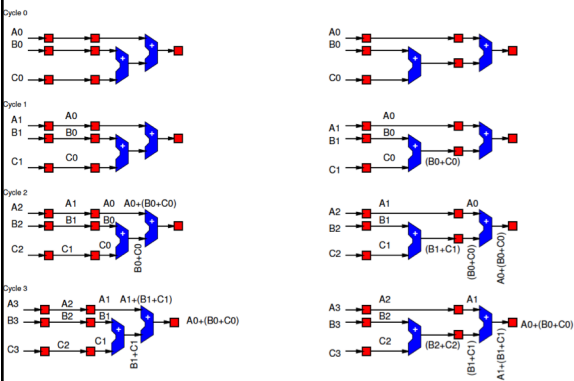
- For this initial design:



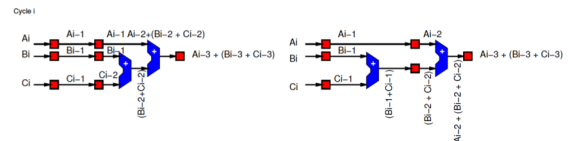
Penn ESE532 Fall 2020 -- DeHon

16

Behavior



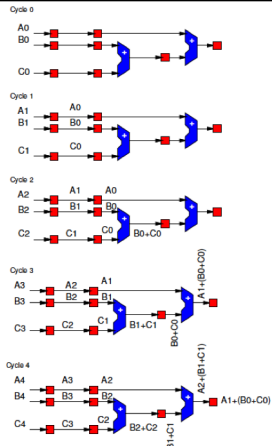
Equations



Penn ESE532 Fall 2020 -- DeHon

18

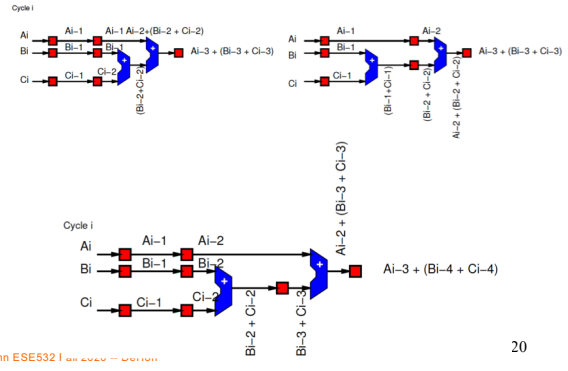
Behavior Add Register



Penn ESE532 Fall 2020 -- DeHon

19

Equations

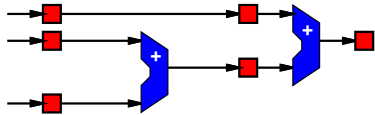


Penn ESE532 Fall 2020 -- DeHon

20

Note Registers on Links

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



Penn ESE532 Fall 2020 -- DeHon

21

Consistent Pipelining

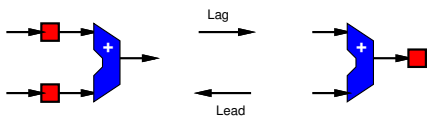
- Makes sure a consistent input set arrives at each gate/operator
 - Don't get mixing between input sets

Penn ESE532 Fall 2020 -- DeHon

22

Legal Register Moves

- Retiming Lag/Lead

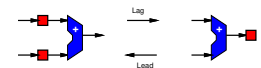


- Lag: remove register every input
add register every output
- Lead: remove register every output
add register every input

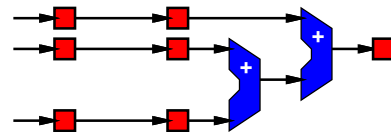
Penn ESE532 Fall 2020 -- DeHon

23

Preclass 1



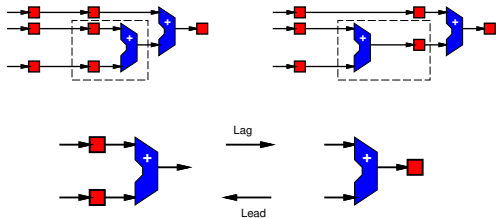
- Retime using Lead/Lag



Penn ESE532 Fall 2020 -- DeHon

24

Preclass 1 (revisited)



Penn ESE532 Fall 2020 -- DeHon

25

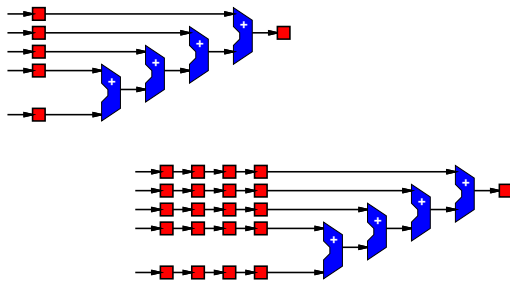
Add Registers and Move

- 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 2020 -- DeHon

26

Add Registers at Input



Penn ESE532 Fall 2020 -- DeHon

27

Add Register and Retime

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

Penn ESE532 Fall 2020 -- DeHon

28

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 2020 -- DeHon

29

Justify Pipelining

(or composing pipelined operators)

Part 2

Penn ESE532 Fall 2020 -- DeHon

30

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 2020 -- DeHon

31

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 2020 -- DeHon

32

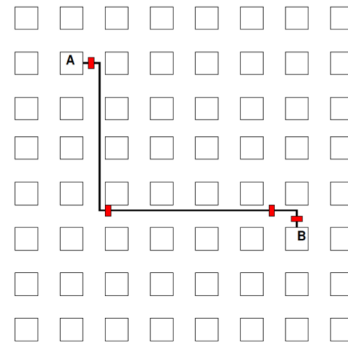
Interconnect Delay

- 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 ESE532 Fall 2020 -- DeHon

33

Interconnect Example



Penn ESE532 Fall 2020 -- DeHon

34

Methodology: Pipelined Operator Graph

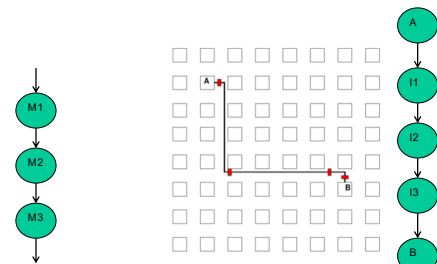
- Start with logical, unpipelined graph
- Treat each pipelined operator as a set of unit-delay operators of mandatory depth
- Treat each interconnect pipeline stage as a unit-delay buffer
- Add registers at input
- Retime into graph

Penn ESE532 Fall 2020 -- DeHon

35

Model

- 3-stage Multiplier
- Interconnect Delay



Penn ESE532 Fall 2020 -- DeHon

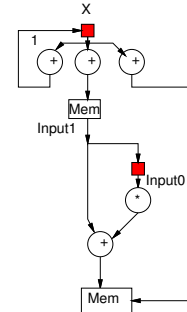
36

Pipeline Loop

(and use for justify pipeline example)

Preclass 4

- Logical (unpipelined) dataflow graph for loop body

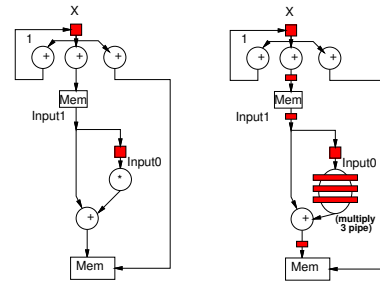


Example Operators

- 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

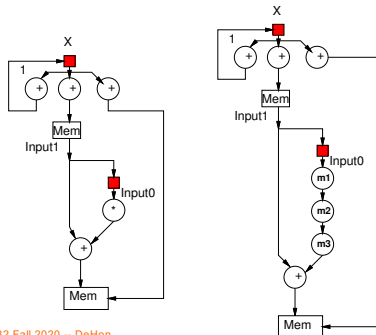
Illustrate Need

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



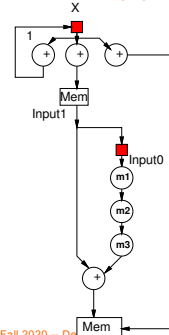
Model Graph

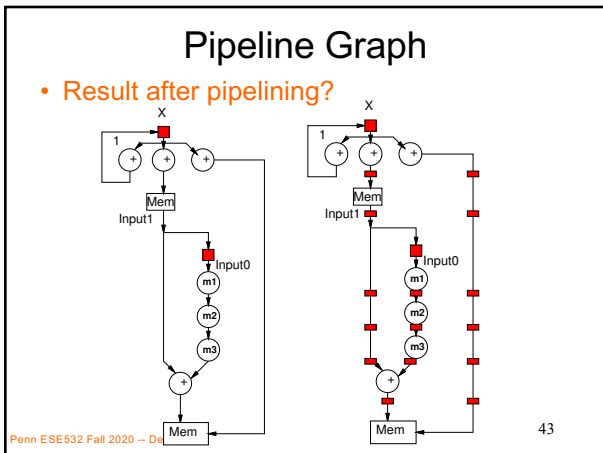
- Revised graph for modeling



Pipeline Graph

- Result after pipelining?



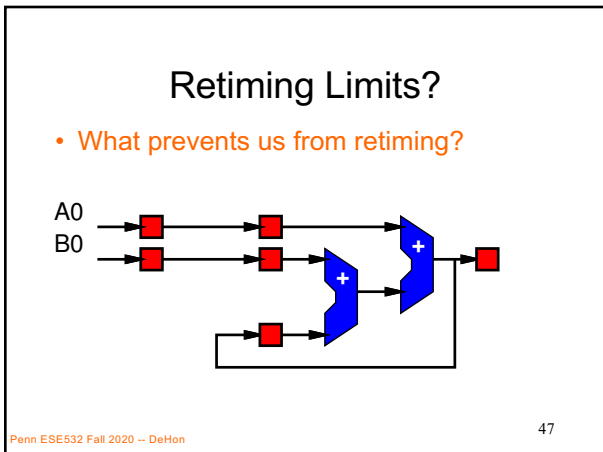
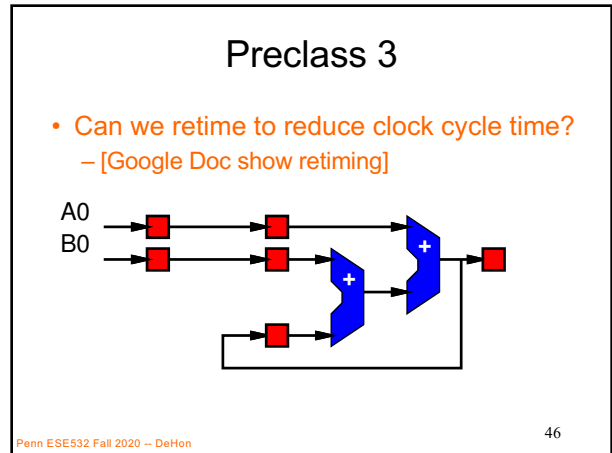


- ### Pipelining Lesson
- Can always pipeline an **acyclic** graph (no graph cycles) to fixed frequency target
 - fixed pipelining of primitive operators
 - Pipeline interconnect delays
 - Need to keep track of registers to balance paths
 - So see consistent delays to operators
- Penn ESE532 Fall 2020 -- DeHon 44

Graph Cycles

Watch: Clock cycle
Cycle time
Cycle in Graph
Part 3

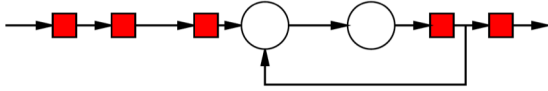
Penn ESE532 Fall 2020 -- DeHon 45



- ### (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 ESE532 Fall 2020 -- DeHon 48

Simple Graph Cycle

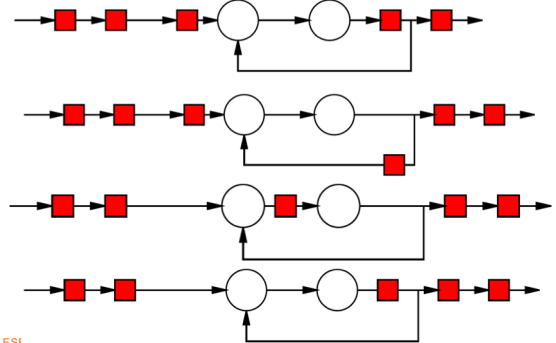
- Delay of graph cycle?
- Registers in graph cycle?
- What happens to graph cycle if try to apply lead/lag?



Penn ESE532 Fall 2020 -- DeHon

49

Retiming



Penn ESE532 Fall 2020 -- DeHon

Loop

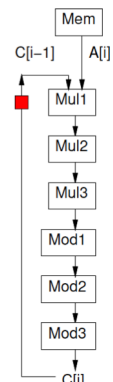
- Consider
 - [multiply and mod each take 3 cycles]
- For (i=0;i<N;i++)
 - $C[i] = (C[i-1] * A[i]) \% N;$

Penn ESE532 Fall 2020 -- DeHon

51

Loop

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

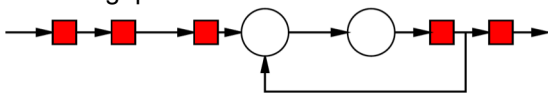


Penn ESE532 Fall 2020 -- DeHon

52

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 = $1/II$



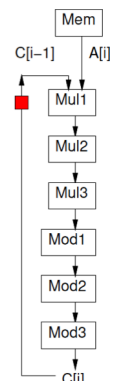
Penn ESE532 Fall 2020 -- DeHon

53

Loop

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

- Initiation Interval?

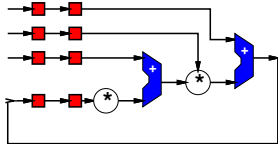


Penn ESE532 Fall 2020 -- DeHon

54

Initial Interval

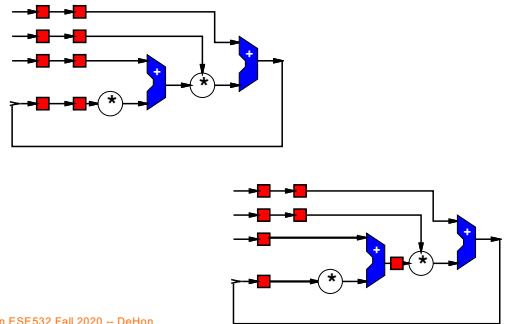
- Delay Around graph cycle?
 - Assume multiply 3, add 1
- Registers in graph cycle?
- Retiming clock cycle bound = II ?
- Achievable?



Penn ESE532 Fall 2020 -- DeHon

55

Retimed

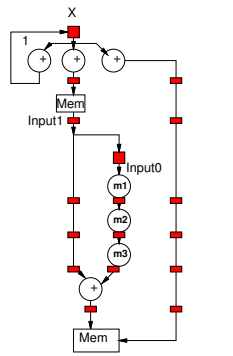


Penn ESE532 Fall 2020 -- DeHon

56

II and Latency

- Actually is a cycle
 - II?
 - Latency?

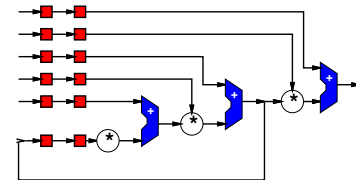


Penn ESE532 Fall 2020 -- DeHon

57

II and Latency

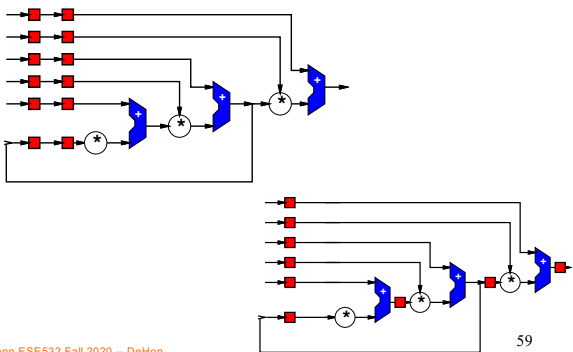
- II? (assume willing to pipeline inputs)
- Latency?



Penn ESE532 Fall 2020 -- DeHon

58

II and Latency



Penn ESE532 Fall 2020 -- DeHon

59

Vector Pipelines

- Data Parallel Vector Operations are interesting even when
 - Vector Lanes < Vector Length
- Within Vector operation, data parallel so no cyclic dependencies
 - So get an II=1 issuing Vector Lane operations
 - May have data dependences between Vector operations

Penn ESE532 Fall 2020 -- DeHon

60

Vector Pipelines

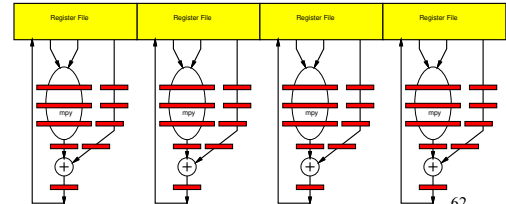
- Data Parallel Vector Operations are interesting even when Vector Lanes < Vector Length
- Within Vector operation, data parallel so no cyclic dependencies
 - So get an $II=1$ issuing Vector Lane operations
 - May have data dependences between Vector operations

Penn ESE532 Fall 2020 -- DeHon

61

Vector Pipeline Example

```
for (int i=0; i<32; i++)
    c[i] += a[i] * b[i]
```

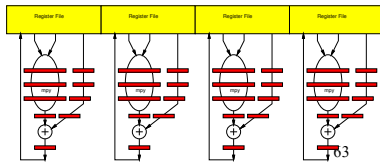


Penn ESE532 Fall 2020 -- DeHon

62

Dependence between Vector Operations

```
for (int i=0; i<32; i++)
    c[i] += a[i] * b[i]
for (int i=0; i<32; i++)
    c[i] += d[i] * e[i]
```



Penn ESE532 Fall 2020 -- DeHon

63

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
- C-slow to share hardware among multiple, data-parallel streams

Penn ESE532 Fall 2020 -- DeHon

64

Admin

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

Penn ESE532 Fall 2020 -- DeHon

65

C-Slow

(Probably record separately)
Part 4

Penn ESE532 Fall 2020 -- DeHon

66

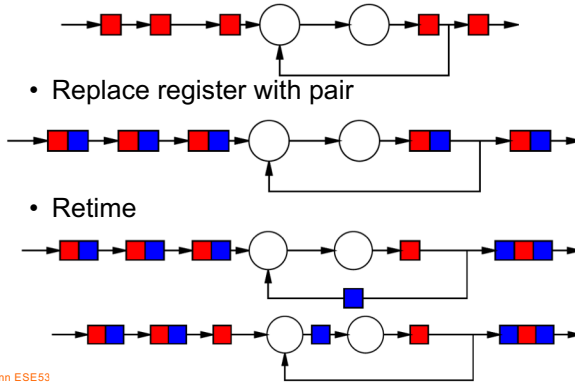
Problem

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

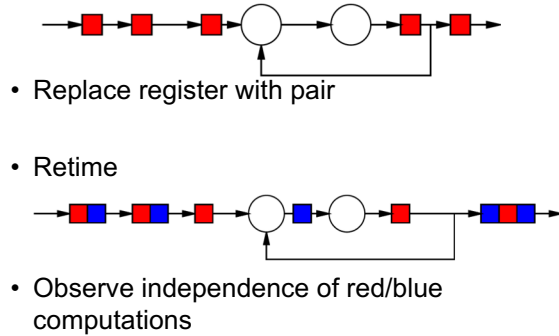
C-Slow

- **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

2-Slow Simple Cycle

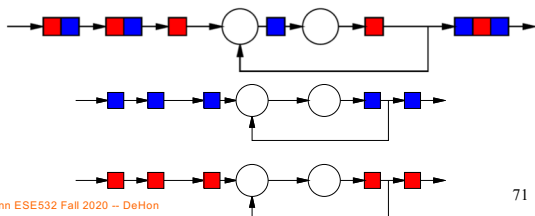


2-Slow Simple Cycle



Equivalence

- The 2-slow operator is equivalent to two data parallel operators running at half the speed
 - E.g. processing separate audio channels



Automation

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

Lesson

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