

# ESE532: System-on-a-Chip Architecture

Day 2: September 4, 2019  
Analysis, Metrics, and Bottlenecks

Work Preclass  
Lecture start 10:35am



Penn ESE532 Spring 2019 -- DeHon

## Today: Analysis

- How do we quickly estimate what's possible?
  - Before (with less effort than) developing a complete solution
- How should we attack the problem?
  - Achieve the performance, energy goals?
- When we don't like the performance we're getting, how do we understand it?
- Where should we spend our time?

Penn ESE532 Spring 2019 -- DeHon

2

## Today: Analysis

- Throughput
- Latency
- Bottleneck
- Computation as a Graph, Sequence
- Critical Path
- Resource Bound
- 90/10 Rule (time permitting)

Penn ESE532 Spring 2019 -- DeHon

3

## Message for Day

- Identify the **Bottleneck**
  - May be in compute, I/O, memory, data movement
- Focus and reduce/remove bottleneck
  - More efficient use of resources
  - More resources
- Repeat

Penn ESE532 Spring 2019 -- DeHon

4

## Latency vs. Throughput

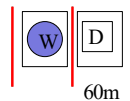
- **Latency:** Delay from inputs to output(s)
- **Throughput:** Rate at which can produce new set of outputs
  - (alternately, can introduce new set of inputs)

Penn ESE532 Spring 2019 -- DeHon

5

## Preclass Washer/Dryer Example

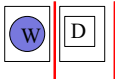
- 10 shirt capacity
- 1 Washer Takes 30 minutes
- 1 Dryer Takes 60 minutes
- How long to do one load of wash?
  - → Wash latency
- Cleaning Throughput?



Penn ESE532 Spring 2019 -- DeHon

6

## Pipeline Concurrency



- Break up the computation graph into stages
  - Allowing us to
    - reuse resources for new inputs (data),
    - while older data is still working its way through the graph
      - Before it has exited graph
  - Throughput > (1/Latency)
- Relate liquid in pipe
  - Doesn't wait for first drop of liquid to exit far end of pipe before accepting second drop

Penn ESE532 Spring 2019 -- DeHon

7

## Escalator



Image Source: [https://commons.wikimedia.org/wiki/File:Tanforan\\_Target\\_escalator\\_1.JPG](https://commons.wikimedia.org/wiki/File:Tanforan_Target_escalator_1.JPG)

Penn ESE532 Spring 2019 -- DeHon

8

## Escalator



- Moves 2 ft/second
- Assume for simplicity one person can step on escalator each second
- Escalator travels 30 feet (vertical and horizontal)
- Latency of escalator trip?
- Throughput of escalator: people/hour ?

Penn ESE532 Spring 2019 -- DeHon

9

## Bottleneck

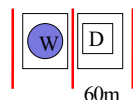
- What is the rate limiting item?
  - Resource, computation, ....

Penn ESE532 Spring 2019 -- DeHon

10

## Preclass Washer/Dryer Example

- 1 Washer Takes 30 minutes
  - Isolated throughput 20 shirts/hour
- 1 Dryer Takes 60 minutes
  - Isolated throughput 10 shirts/hour
- Where is bottleneck in our cleaning system?

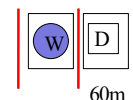


Penn ESE532 Spring 2019 -- DeHon

11

## Preclass Washer/Dryer Example

- 1 Washer \$500
  - Isolated throughput 20 shirts/hour
- 1 Dryer \$500
  - Isolated throughput 10 shirts/hour
- How do we increase throughput with \$500 investment

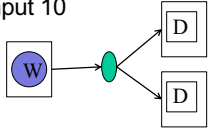


Penn ESE532 Spring 2019 -- DeHon

12

### Preclass Washer/Dryer Example

- 1 Washer \$500
  - Isolated throughput 20 shirts/hour
- 2 Dryers \$500
  - Isolated single dryer throughput 10 shirts/hour
- Latency?
- Throughput?

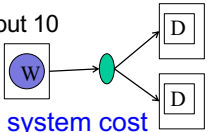


13

Penn ESE532 Spring 2019 -- DeHon

### Preclass Washer/Dryer Example

- 1 Washer \$500
  - Isolated throughput 20 shirts/hour
- 2 Dryers \$500
  - Isolated single dryer throughput 10 shirts/hour
- Able to double the throughput without doubling system cost

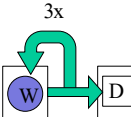


14

Penn ESE532 Spring 2019 -- DeHon

### Preclass Stain Example

- 1 Washer Takes 30 minutes
  - Isolated throughput 20 shirts/hour
- 1 Dryer Takes 60 minutes
  - Isolated throughput 10 shirts/hour
- Shirt need 3 wash cycles
- Latency?
- Throughput?
  - (assuming reuse single washer)



15

Penn ESE532 Spring 2019 -- DeHon

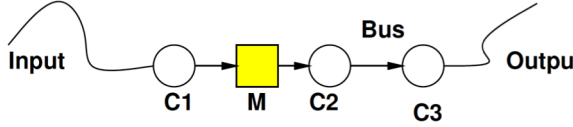
## Beyond Computation

16

Penn ESE532 Spring 2019 -- DeHon

### Bottleneck

- May be anywhere in path
  - I/O, compute, memory, data movement

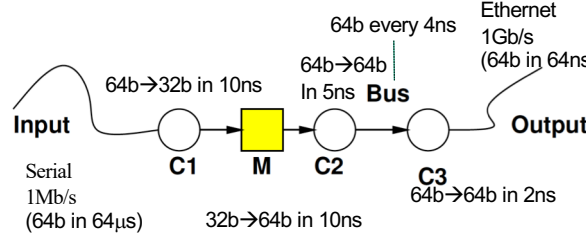


17

Penn ESE532 Spring 2019 -- DeHon

### Bottleneck

- Where bottleneck?

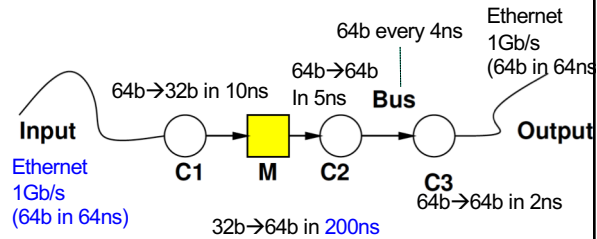


18

Penn ESE532 Spring 2019 -- DeHon

## Bottleneck

- Where bottleneck?

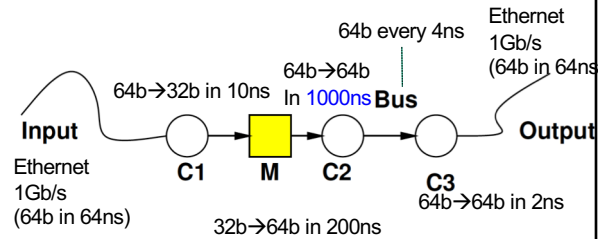


Penn ESE532 Spring 2019 -- DeHon

19

## Bottleneck

- Where bottleneck?



Penn ESE532 Spring 2019 -- DeHon

20

## Feasibility / Limits

- First things to understand
  - Obvious limits in system?
- Impossible?
- Which aspects will demand efficient mapping?
- Where might there be spare capacity?

Penn ESE532 Spring 2019 -- DeHon

21

## Generalizing

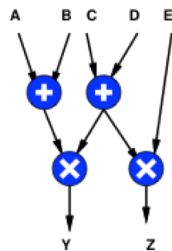
(to more general task graphs)

Penn ESE532 Spring 2019 -- DeHon

22

## Computation as Graph

- Shown “simple” graphs (pipelines) so far
- $Y = (A+B) * (C+D)$
- $Z = (C+D) * E$



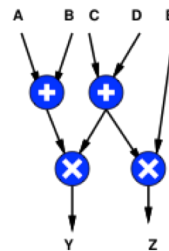
Note: HW2 ask you to draw a dataflow graph. Here's an example...more to come.

Penn ESE532 Spring 2019 -- DeHon

23

## Computation as Graph

- Nodes have multiple input/output edges
- Edges may fanout
  - Results go to multiple successors

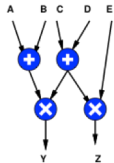


Penn ESE532 Spring 2019 -- DeHon

24

## Computation as Sequence

- Shown “simple” graphs (pipelines) so far
- $Y=(A+B)*(C+D)$
- $Z=(C+D)*E$



$$T1=A+B$$

$$T2=C+D$$

$$Y=T1*T2$$

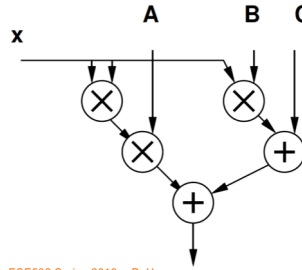
$$Z=T1*E$$

Penn ESE532 Spring 2019 -- DeHon

25

## Computation as Graph

$$Y=Ax^2+Bx+C$$



$$T1=x*x$$

$$T2=A*T1$$

$$T3=B*x$$

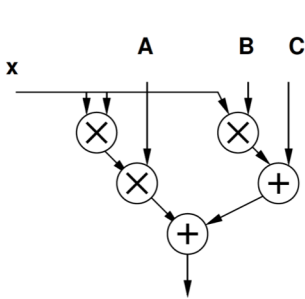
$$T4=T2+T3$$

$$Y=C+T4$$

Penn ESE532 Spring 2019 -- DeHon

26

## Computation as Graph



- Latency multiply = 1
- Latency add = 1/3
- Latency from B to output?
- Latency from x to output?
  - Through  $Ax^2$  ?
  - Through  $Bx$  ?

Penn ESE532 Spring 2019 -- DeHon

27

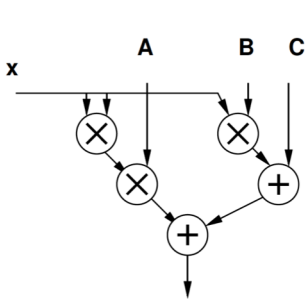
## Delay in Graphs

- **Observe:** There are multiple paths from inputs to outputs
- Need to complete all of them to produce outputs
- Limited by longest path
- **Critical path:** longest path in the graph

Penn ESE532 Spring 2019 -- DeHon

28

## Computation as Graph



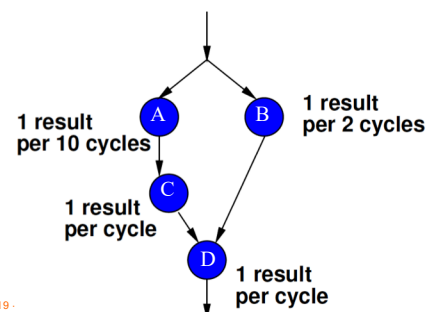
- Latency multiply = 1
- Latency add = 1/3
- Critical Path?

Penn ESE532 Spring 2019 -- DeHon

29

## Bottleneck

- Where is the bottleneck?



Penn ESE532 Spring 2019 .

## Time and Space

Penn ESE532 Spring 2019 -- DeHon

31

## Space

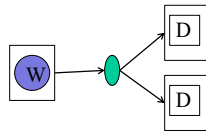
- “Space” is an abstract term for physical resources
  - On VLSI chip: Area – mm<sup>2</sup> of silicon
  - On our FPGA: # of LUTs used
  - More abstractly: # of Adders, multipliers
  - Laundry example
    - \$\$ to spend on laundry equipment
    - Physical space (sq. ft) in laundry room

Penn ESE532 Spring 2019 -- DeHon

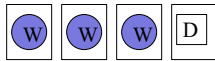
32

## Space-Time

- In general, we can spend resources to reduce time
  - Increase throughput



Three wash stain removal case

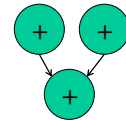


Penn ESE532 Spring 2019 -- DeHon

33

## Space Time

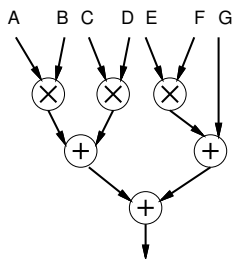
- Computation
  - $A = x0 + x1$
  - $B = x2 + x3$
  - $C = A + B$
- Adder takes one cycle
  - Latency on one adder?
  - Latency on 3 adders?



Penn ESE532 Spring 2019 -- DeHon

34

## Computation as Graph



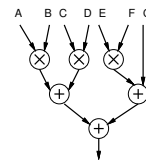
- Latency multiply = 1
- Space multiply = 3
- Latency add = 1
- Space add = 1
- Latency and Space
  - 3 mul, 2 add

Penn ESE532 Spring 2019 -- DeHon

35

## Schedule 3 mul, 2 add

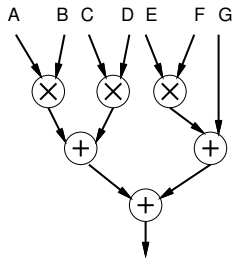
Cycle	Mul	Mul	Mul	Add	Add
0	A*B	C*D	E*F		
1				A*B+C*D	E*F+G
2				(A*B+C*D)	+(E*F+G)



Penn ESE532 Spring 2019 -- DeHon

36

## Computation as Graph



- Latency multiply = 1
- Space multiply = 3
- Latency add = 1
- Space add = 1
- Latency and Space – 2 mul, 1 add

Penn ESE532 Spring 2019 -- DeHon

37

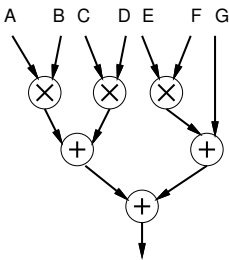
## Schedule 2 mul, 1 add

Cycle	Mul	Mul	Add
0	A*B	C*D	
1	E*F		(A*B+C*D)
2			E*F+G
3			(A*B+C*D)+(E*F+G)

Penn ESE532 Spring 2019 -- DeHon

38

## Computation as Graph



- Latency multiply = 1
- Space multiply = 3
- Latency add = 1
- Space add = 1
- Latency and Space – 1 mul, 1 add

Penn ESE532 Spring 2019 -- DeHon

39

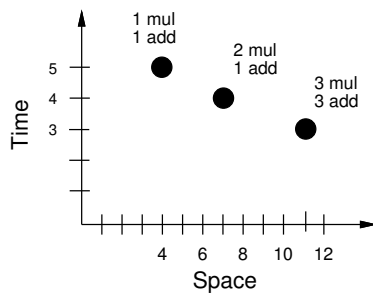
## Schedule 1 mul, 1 add

Cycle	Mul	Add
0	A*B	
1	C*D	
2	E*F	A*B+C*D
3		E*F+G
4		(A*B+C*D)+(E*F+G)

Penn ESE532 Spring 2019 -- DeHon

40

## Space-Time Graph



(corrected after lecture)

Penn ESE532 Spring 2019 -- DeHon

41

## Two Bounds

(still in Time and Space)

Penn ESE532 Spring 2019 -- DeHon

42

## Problem

- Coming up with an exact time count can be hard (human/computer time consuming)
  - Technically a hard problem
    - NP-Complete: no known non-exponential solution
- Requires reasoning about structure of graph
- Would be nice to have a quick (easy) answer on what is feasible
  - ...and what is not feasible → impossible.

Penn ESE532 Spring 2019 -- DeHon

43

## Bounds

- Establish the feasible range
  - Must be larger than LB
  - Must be smaller than UB
  - Solution will be between LB and UB
  - $LB \leq ActualTime \leq UB$
- Bounds in sports
  - Ball landing in-bounds or out-of bounds

Penn ESE532 Spring 2019 -- DeHon

44

## Bounds

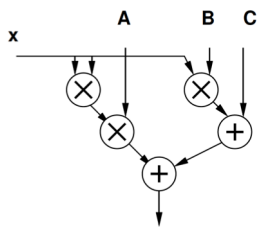
- Quick **lower** bounds (LB) can estimate
  - $LB \leq ActualTime$
- Two:
  - CP: Critical Path
    - Sometimes call it “Latency Bound”
  - RB: Resource Capacity Bound
    - Sometimes call it “Throughput Bound” or “Compute Bound”

Penn ESE532 Spring 2019 -- DeHon

45

## Critical Path Lower Bound

- Critical path assuming infinite resources
- Certainly cannot finish any faster than that
- $CP \leq ActualTime$
- Ignores resource limits



Penn ESE532 Spring 2019 -- DeHon

46

## Resource Capacity Lower Bound

- Sum up all capacity required per resource:  $TotalOps = \sum Ops$ 
  - E.g. number of multiplications, additions, memory lookups
- Divide by total resource (for type)
  - E.g., number of multipliers, adders, memory ports
  - $RB = \lceil TotalOps / Operators \rceil \leq ActualTime$
- Lower bound on compute
  - (best can do is pack all use densely)
  - Ignores data dependency constraints

Penn ESE532 Spring 2019 -- DeHon

47

## Multiple Resource Types

- $RB = \text{Max}(\lceil TotalOps_1 / Operators_{s_1} \rceil, \lceil TotalOps_2 / Operators_{s_2} \rceil, \dots) \leq ActualTime$
- Combine Critical Path Lower Bound
 
$$\text{Max}(CP, \lceil TotalOps_1 / Operators_{s_1} \rceil, \lceil TotalOps_2 / Operators_{s_2} \rceil, \dots) \leq ActualTime$$

Penn ESE532 Spring 2019 -- DeHon

48



## For Single Resource Type

- (and no communication time...)
- Can use to get upper bound:
- $ActualTime \leq CP + RB$
- Together:
- $Max(CP, RB) \leq ActualTime \leq CP + RB$

## Washer-Dryer Bounds

- Task: wash & dry 30 shirts
- Washer: 10 shirts/30 min.
- Dryer: 10 shirts/60 min.
- 2 Washers, 2 Dryers
- W-->D CP=90 minutes

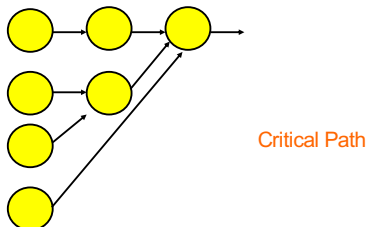
## Washer-Dryer Bounds

- Task: wash & dry 30 shirts CP=90 min
- Washer: 10 shirts/30 min.
- Dryer: 10 shirts/60 min.
- 2 Washers, 2 Dryers
- Washer Bound:
  - $WB = \lceil 30 \text{ shirts} / 2 \text{ washers} \rceil \times 10 \text{ shirts/washer} = 60 \text{ min}$
- Dryer
  - $DB = \lceil (30 \text{ shirts} / 2 \text{ dryers}) \rceil \times 10 \text{ shirts/dryer} = 120 \text{ minutes}$

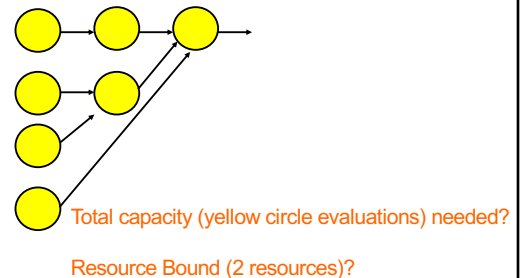
## Washer-Dryer Bounds

- Task: wash & dry 30 shirts CP=90 min
- Washer: 10 shirts/30 min.
- Dryer: 10 shirts/60 min.
- 2 Washers, 2 Dryers RB=120 min
- $Max(90, 60, 120) = 120 \leq TaskTime$
- TaskTime: 150
  - 0 start 2 washes
  - 30 start 2 dryers; start 1 wash
  - 90 finish 2 dryers; start last dryer load
  - 150 finish last dryer load

## Example



## Example



### Example

Cycle	Resource 1	Resource 2
0	A	B
1	C	D
2	E	F
3	G	

Resource Bound (2 resources)?

Penn ESE532 Spring 2019 -- DeHon

55

### Example

Resource Bound (4 resources)?

Penn ESE532 Spring 2019 -- DeHon

56

### Example

Cycle	R1	R2	R3	R4
0	A	B	C	D
1	E	F	G	

Legal Schedule?

Penn ESE532 Spring 2019 -- DeHon

57

### Resource Capacity Lower Bound

- Sum up all capacity required per resource:  $TotalOps = \sum Ops$ 
  - E.g. number of multiplications, additions, memory lookups
- Divide by total resource (for type)
  - E.g., number of multipliers, adders, memory ports
  - $RB = \lceil TotalOps / Operators \rceil \leq ActualTime$
- Lower bound on compute
  - (best can do is pack all use densely)
  - **Ignores data dependency constraints**

Penn ESE532 Spring 2019 -- DeHon

58

### Example

Critical Path     3

Resource Bound (2 resources)      $7/2=4$

Resource Bound (4 resources)      $7/4=2$

Penn ESE532 Spring 2019 -- DeHon

59

### 90/10 Rule (of Thumb)

- Observation that code is not used uniformly
- 90% of the time is spent in 10% of the code
- Knuth: 50% of the time in 2% of the code
- Implications
  - There will typically be a bottleneck
  - We don't need to optimize everything
  - We don't need to uniformly replicate space to achieve speedup
  - Not everything needs to be accelerated

Penn ESE532 Spring 2019 -- DeHon

60

## Big Ideas

- Identify the Bottleneck
  - May be in compute, I/O, memory ,data movement
- Focus and reduce/remove bottleneck
  - More efficient use of resources
  - More resources

## Admin

- Reading for Day 3 on web
- HW1 due Friday
- HW2 out
  - Partner assignment (see canvas)
- Remember feedback