

# ESE532: System-on-a-Chip Architecture

Day 2: September 5, 2018  
Analysis, Metrics, and Bottlenecks

Work Preclass  
Lecture start 10:35pm



Penn ESE532 Spring 2018 -- 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 2018 -- DeHon

2

## Today: Analysis

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

Penn ESE532 Spring 2018 -- 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 2018 -- DeHon

4

## Latency vs. Throughput

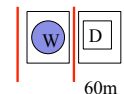
- **Latency:** Delay from inputs to output(s)
- **Throughput:** Rate at which can produce new set of outputs
  - (alternately, can introduce new set of inputs)

Penn ESE532 Spring 2018 -- 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 2018 -- 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 2018 -- 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 2018 -- 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 2018 -- DeHon 9

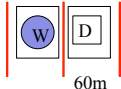
## Bottleneck

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

Penn ESE532 Spring 2018 -- 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?

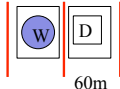


60m

Penn ESE532 Spring 2018 -- 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

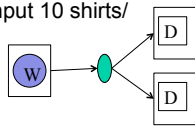


60m

Penn ESE532 Spring 2018 -- 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?

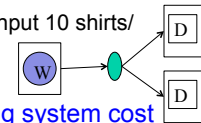


Penn ESE532 Spring 2018 -- DeHon

13

## 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

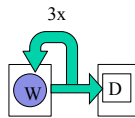


Penn ESE532 Spring 2018 -- DeHon

14

## 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 share)?



Penn ESE532 Spring 2018 -- DeHon

15

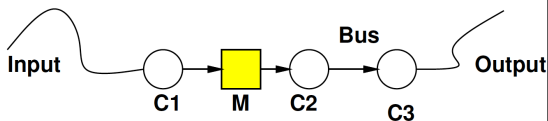
## Beyond Computation

Penn ESE532 Spring 2018 -- DeHon

16

## Bottleneck

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

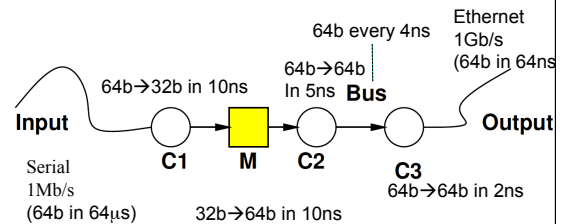


Penn ESE532 Spring 2018 -- DeHon

17

## Bottleneck

- Where bottleneck?



Penn ESE532 Spring 2018 -- DeHon

18

### Bottleneck

- Where bottleneck?

Input: Ethernet 1Gb/s (64b in 64ns)

C1: 64b → 32b in 10ns

M: 32b → 64b in 200ns

C2: 64b → 64b In 5ns

Bus: 64b every 4ns

C3: 64b → 64b in 2ns

Output: Ethernet 1Gb/s (64b in 64ns)

Penn ESE532 Spring 2018 -- DeHon 19

### Bottleneck

- Where bottleneck?

Input: Ethernet 1Gb/s (64b in 64ns)

C1: 64b → 32b in 10ns

M: 32b → 64b in 200ns

C2: 64b → 64b In 1000ns

Bus: 64b every 4ns

C3: 64b → 64b in 2ns

Output: Ethernet 1Gb/s (64b in 64ns)

Penn ESE532 Spring 2018 -- 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 2018 -- DeHon 21

### Generalizing

Penn ESE532 Spring 2018 -- DeHon 22

### Computation as Graph

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

Penn ESE532 Spring 2018 -- DeHon 23

### Computation as Graph

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

Penn ESE532 Spring 2018 -- DeHon 24

## Computation as Sequence

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

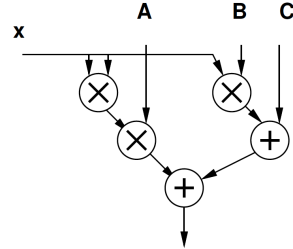
$$\begin{aligned} T1 &= A+B \\ T2 &= C+D \\ Y &= T1*T2 \\ Z &= T1*E \end{aligned}$$

Penn ESE532 Spring 2018 -- DeHon

25

## Computation as Graph

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

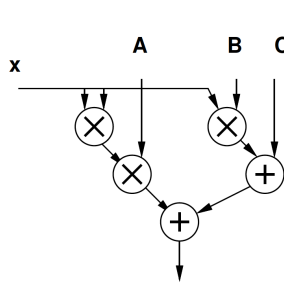


$$\begin{aligned} T1 &= x*x \\ T2 &= A*T1 \\ T3 &= B*x \\ T4 &= T2+T3 \\ Y &= C+T4 \end{aligned}$$

Penn ESE532 Spring 2018 -- DeHon

26

## Computation as Graph



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

Penn ESE532 Spring 2018 -- DeHon

27

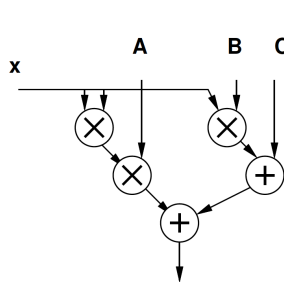
## Delay in Graphs

- 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 2018 -- DeHon

28

## Computation as Graph



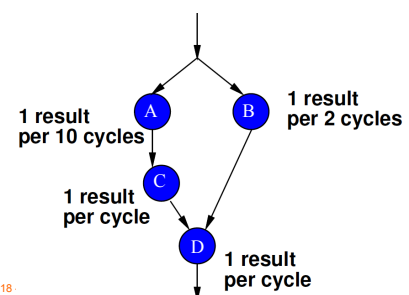
- Latency multiply = 3
- Latency add = 1
- **Critical Path?**

Penn ESE532 Spring 2018 -- DeHon

29

## Bottleneck

- **Where is the bottleneck?**



Penn ESE532 Spring 2018

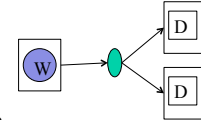
## Time and Space

Penn ESE532 Spring 2018 -- DeHon

31

## Space-Time

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



Three wash stain removal case



Penn ESE532 Spring 2018 -- DeHon

32

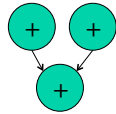
## Space Time

- Computation

- $A=x0+x1$
- $B=A+x2$
- $C=B+x3$

- Adder takes one cycle

- Throughput on one adder?
- Throughput on 3 adders?



Penn ESE532 Spring 2018 -- DeHon

33

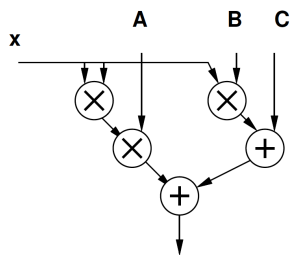
## Dependencies and S-T

- Dependencies may limit throughput acceleration
  - Give benefit less than  $1/\text{space}$

Penn ESE532 Spring 2018 -- DeHon

34

## Computation as Graph

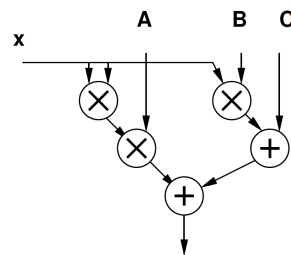


- Latency multiply = 1
- Thput mult = 1
- Space multiply = 3
- Latency add = 1/3
- Space add = 1
- Thput and Space
  - 3 mul, 2 add

Penn ESE532 Spring 2018 -- DeHon

35

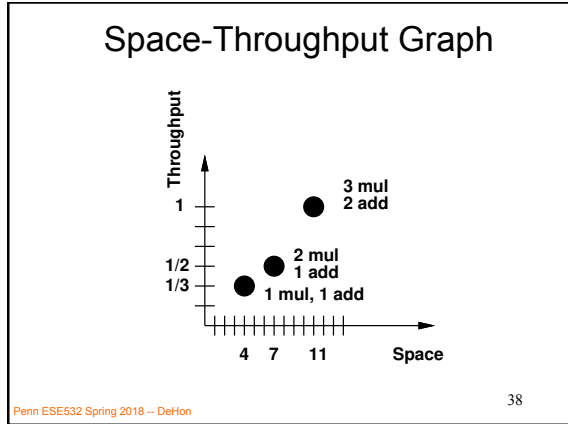
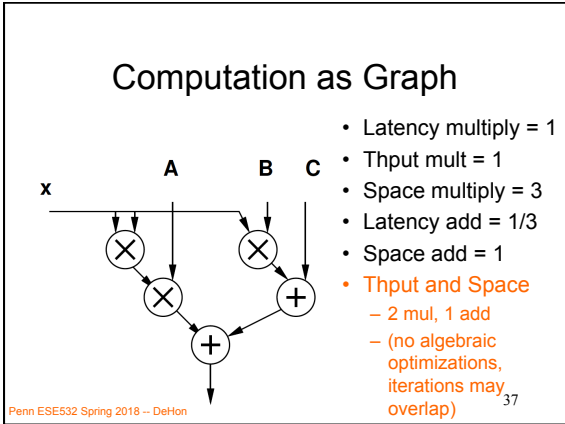
## Computation as Graph



- Latency multiply = 1
- Thput mult = 1
- Space multiply = 3
- Latency add = 1/3
- Space add = 1
- Thput and Space
  - 1 mul, 1 add
  - Where is bottleneck?

Penn ESE532 Spring 2018 -- DeHon

36



### Two Bounds

(still in Time and Space)

Penn ESE532 Spring 2018 -- DeHon 39

### Bounds

- Quick lower bounds can estimate
- Two:
  - CP: Critical Path
    - Sometimes call it "Latency Bound"
  - RB: Resource Bound
    - Sometimes call it "Throughput Bound" or "Compute Bound"

Penn ESE532 Spring 2018 -- DeHon 40

### Critical Path Lower Bound

- Critical path assuming infinite resources
- Certainly cannot finish any faster than that

Penn ESE532 Spring 2018 -- DeHon 41

### Resource Capacity Lower Bound

- Sum up all capacity required per resource
  - E.g. number of multiplications, additions, memory lookups
- Divide by total resource (for type)
  - E.g., number of multipliers, adders, memory ports
- Lower bound on compute
  - (best can do is pack all use densely)
  - **Ignores data dependency constraints**

Penn ESE532 Spring 2018 -- DeHon 42

### Example

Critical Path

Penn ESE532 Spring 2018 -- DeHon

43

### Example

Total capacity (yellow circle evaluations) needed?  
Resource Bound (2 resources)?

Penn ESE532 Spring 2018 -- DeHon

44

### Example

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

Resource Bound (2 resources)?

Penn ESE532 Spring 2018 -- DeHon

45

### Example

Resource Bound (4 resources)?

Penn ESE532 Spring 2018 -- DeHon

46

### Example

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

Legal Schedule?

Penn ESE532 Spring 2018 -- DeHon

47

## Resource Capacity Lower Bound

- Sum up all capacity required per resource
  - E.g. number of multiplications, additions, memory lookups
- Divide by total resource (for type)
  - E.g., number of multipliers, adders, memory ports
- Lower bound on compute
  - (best can do is pack all use densely)
  - Ignores data dependency constraints

Penn ESE532 Spring 2018 -- DeHon

48



### Example

Critical Path    3

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

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

Penn ESE532 Spring 2018 -- DeHon

49

### Critical Path

- Latency multiply = 3
- Thput mult = 1/3
- Space multiply = 3
- Latency add = 1
- Space add = 1
- **Critical Path?**

Penn ESE532 Spring 2018 -- DeHon

50

### Resource Bound

- Latency multiply = 3
- Thput mult = 1/3
- Space multiply = 3
- Latency add = 1
- Space add = 1
- **Resource Bound**
  - 1 mul, 1 add
  - 2 mul, 1 add
  - 3 mul, 2 add

Penn ESE532 Spring 2018 -- DeHon

51

### 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 2018 -- DeHon

52

### Amdahl's Law

- If you only speedup  $Y(\%)$  of the code, the most you can accelerate your application is  $1/(1-Y)$
- $T_{\text{before}} = 1*Y + 1*(1-Y)$
- Speedup by factor of  $S$
- $T_{\text{after}} = (1/S)*Y + 1*(1-Y)$
- Limit  $S \rightarrow \text{infinity}$   $T_{\text{before}}/T_{\text{after}} = 1/(1-Y)$

Penn ESE532 Spring 2018 -- DeHon

53

### Amdahl's Law

- $T_{\text{before}} = 1*Y + 1*(1-Y)$
- Speedup by factor of  $S$
- $T_{\text{after}} = (1/S)*Y + 1*(1-Y)$
- $Y=70\%$ 
  - Possible speedup ( $S \rightarrow \text{infinity}$ ) ?
  - Speedup if  $S=10$ ?

Penn ESE532 Spring 2018 -- DeHon

54

## Amdahl's Law

- If you only speedup  $Y(\%)$  of the code, the most you can accelerate your application is  $1/(1-Y)$
- Implications
  - Amdahl: good to have a fast sequential processor
  - Keep optimizing
    - $T_{\text{after}} = (1/S) * Y + 1 * (1-Y)$
    - For large  $S$ , bottleneck now in the  $1-Y$

Penn ESE532 Spring 2018 -- DeHon

55

## 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

Penn ESE532 Spring 2018 -- DeHon

56

## Admin

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

Penn ESE532 Spring 2018 -- DeHon

57