

# ESE5320: System-on-a-Chip Architecture

Day 4: September 14, 2022  
Parallelism Overview

Masks  
Board holders pickup boards  
Preclass



Penn ESE5320 Fall 2022 -- DeHon

1

## Today

- Compute Models (Part 1)
  - How do we *express* and reason about parallel execution freedom
- Types of Parallelism (Part 2)
  - How can we slice up and think about parallelism?
  - How *exploit* parallelism

Penn ESE5320 Fall 2022 -- DeHon

2

2

## Message

- Many useful models for parallelism
  - Help conceptualize
- One-size does not fill all
  - Match to problem
  - Will want to exploit all of them

Penn ESE5320 Fall 2022 -- DeHon

3

3

## Parallel Compute Models

Control Flow, Dataflow  
Combining  
Explicit, Implicit Parallelism

Penn ESE5320 Fall 2022 -- DeHon

4

4

## Term: Operation

- **Operation** – logic computation to be performed

Penn ESE5320 Fall 2022 -- DeHon

5

5

## Sequential Control Flow

### Control flow

- Program is a sequence of operations
- Operation reads inputs and writes outputs into common store (memory)
- One operation runs at a time
  - defines successor

Model of correctness is sequential execution

### Examples

- C (Java, ...)
- Finite-State Machine (FSM)
- Finite Automata (FA)
- assembly code (ISA)

Penn ESE5320 Fall 2022 -- DeHon

6

6

## Parallelism can be explicit

- State which operations occur on a cycle
- Multiply, add for quadratic equation

cycle	mpy	add
1	B,x	
2	x,x	(Bx)+C
3	A,x <sup>2</sup>	
4		Ax <sup>2</sup> +(Bx+C)

Penn ESE5320 Fall 2022 -- DeHon

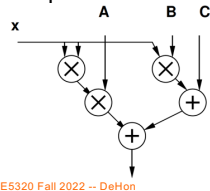
7

7

## Parallelism can be implicit

- Sequential expression
- Infer data dependencies

T1=x\*x  
T2=A\*T1  
T3=B\*x  
T4=T2+T3  
Y=C+T4



- Or
- $$Y=A*x*x+B*x+C$$

Penn ESE5320 Fall 2022 -- DeHon

8

8

## Implicit Parallelism

- $d=(x1-x2)*(x1-x2) + (y1-y2)*(y1-y2)$
- What parallelism exists here?

Penn ESE5320 Fall 2022 -- DeHon

9

9

## Parallelism can be implicit

- Sequential expression
- Infer data dependencies

for (i=0;i<100;i++)  
y[i]=A\*x[i]\*x[i]+B\*x[i]+C

Why can these operations be performed in parallel?

Penn ESE5320 Fall 2022 -- DeHon

10

10

## Dataflow / Control Flow

### Dataflow

- Program is a graph of operations
- Operation consumes **tokens** and produces tokens
- All operations run concurrently

### Control flow (e.g. C)

- Program is a sequence of operations
- Operation reads inputs and writes outputs into common store
- One operation runs at a time
  - defines successor

Penn ESE5320 Fall 2022 -- DeHon

11

11

## Token

- Data value with presence indication
  - May be conceptual
    - Only exist in high-level model
    - Not kept around at runtime
  - Or may be physically represented
    - One bit represents presence/absence of data

Penn ESE5320 Fall 2022 -- DeHon

12

12

## FIFO

Write  
DataIn

Empty  
DataOut  
Read

**FIFO**

- Hardware Block
- Outputs data in order received
  - First-In, First-Out
- Tell it when you are providing data
  - Write
  - May choose not to insert on a cycle
    - Need to signal
- Tell it when you are consuming data
  - Read
- Tells you when it's **empty** and has no data to provide
- Tells you when it's **full** and can hold nothing else

What are data presence indicators here?

Penn ESE5320 Fall 2022 -- DeHon 13

## Token Examples?

- How serial link know character present?
- How signal miss in processor data cache and processor needs to wait for data?

Penn ESE5320 Fall 2022 -- DeHon 14

13

14

## Operation

- Takes in one or more inputs
- Computes on the inputs
- Produces results
- Logically **self-timed**
  - “Fires” only when input set present
  - Signals availability of output

Penn ESE5320 Fall 2022 -- DeHon 15

15

Penn ESE5320 Fall 2022 -- DeHon 16

16

## Dataflow Graph

- Represents
  - computation sub-blocks
  - linkage
- Abstractly
  - controlled by data presence

Penn ESE5320 Fall 2022 -- DeHon 17

17

## Dataflow Graph Example

Penn ESE532 18

18

## Dataflow / Control Flow

### Dataflow

- Program is a graph of operations
- Operation consumes **tokens** and produces tokens
- All operations run concurrently

### Control flow (e.g. C)

- Program is a sequence of operations
- Operation reads inputs and writes outputs into common store
- One operation runs at a time
  - defines successor

Penn ESE5320 Fall 2022 -- DeHon

19

19

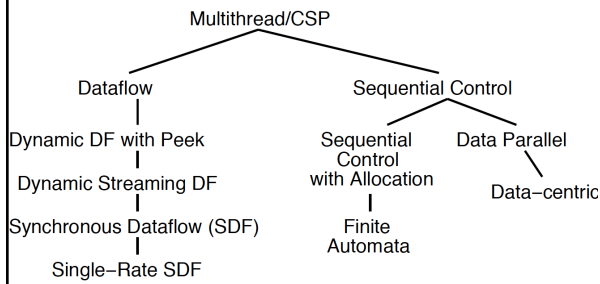
## Communicating Threads

- Computation is a collection of sequential/control-flow “threads”
- Threads may communicate
  - Through dataflow I/O
  - (Through shared variables)
- View as hybrid or generalization
  - Of control flow and dataflow
- CSP – Communicating Sequential Processes → canonical model example<sub>20</sub>

Penn ESE5320 Fall 2022 -- DeHon

20

## Compute Models



Penn ESE5320 Fall 2022 -- DeHon

22

22

## All Used

- All of these things get used in modern CPUs and SoCs
  - Sequential control flow
  - Operation parallelism
  - Data presence and data-driven flow
  - Multiple threads
  - Data Parallel

Penn ESE5320 Fall 2022 -- DeHon

23

23

## Value of Multiple Models

- When you have a big enough hammer, everything looks like a nail.
- Many stuck on single model
  - Try to make all problems look like their nail
- Value to diversity / heterogeneity
  - One size does not fit all



Penn ESE5320 Fall 2022 -- DeHon

24

24

## Types of Parallelism

Part 2

Penn ESE5320 Fall 2022 -- DeHon

25

25

## Types of Parallelism

- **Data Level** – Perform same computation on different data items
- **Thread or Task Level** – Perform separable (perhaps heterogeneous) tasks independently
- **Instruction Level** – Within a single sequential thread, perform multiple operations on each cycle.

Penn ESE5320 Fall 2022 -- DeHon

26

26

## Pipeline Parallelism

- Pipeline – organize computation as a spatial sequence of concurrent operations
  - Can introduce new inputs before finishing
  - Instruction- or thread-level
  - Use for data-level parallelism
  - Can be directed graph

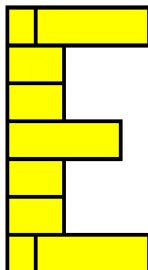
Penn ESE5320 Fall 2022 -- DeHon

27

27

## Sequential

- Single person build E
- Latency?
- Throughput?



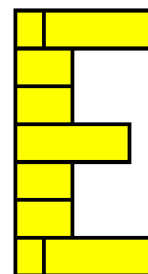
Penn ESE5320 Fall 2022 -- DeHon

28

28

## Data Parallel

- Everyone in class build own E
- Latency?
- Throughput?
- Ideal speedup?
- Resource Bound?
  - 100 Es, 12 people
- When useful?



Penn ESE5320 Fall 2022 -- DeHon

29

29

## Data-Level Parallelism

- **Data Level** – Perform same computation on different data items
- Resource Bound:  $T_{dp} = T_{seq}/P$
- (with enough independent problems, match our resource bound computation)

Penn ESE5320 Fall 2022 -- DeHon

30

30

## Thread Parallel

- Each person build distinct letter or number (e.g. E, S, 5, 3, 2, 0)
- Latency? (assume each has  $\leq 9$  bricks)
- Throughput?
  - Build 6 distinct letters
  - Using whole class ( $\geq 6$  people)
  - (distinct letters/time-unit)
- Speedup over sequential build of 6 letters?

Penn ESE5320 Fall 2022 -- DeHon

31

31

## Thread-Level Parallelism

- **Thread or Task Level** – Perform separable (perhaps heterogeneous) tasks independently
- Resource Bound:  $T_{tp} = T_{seq}/P$
- $T_{tp} = \max(T_{t1}, T_{t2}, T_{t3}, \dots)$ 
  - Less speedup than ideal if not balanced
- Can produce a diversity of calculations
  - Useful if have limited need for the **same** calculation

Penn ESE5320 Fall 2022 -- DeHon

32

32

## Instruction-Level Parallelism

- Build single letter in lock step
- Group of 3
  - [2 volunteers; steps up front]
- Resource Bound for 3 people building 9-brick letter?
- Announce steps from slide
  - Stay in step with slides

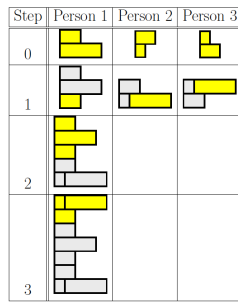
Penn ESE5320 Fall 2022 -- DeHon

33

33

## Group Communication

- Groups of 3
- Note who was person 1 task
- 2, 3 will need to pass completed substructures

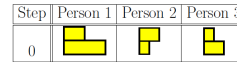


Penn ESE5320 Fall 2022 -- DeHon

34

34

## Step 0

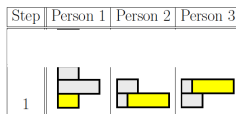


Penn ESE5320 Fall 2022 -- DeHon

35

35

## Step 1

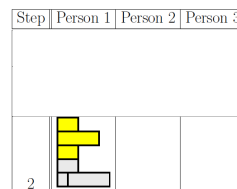


Penn ESE5320 Fall 2022 -- DeHon

36

36

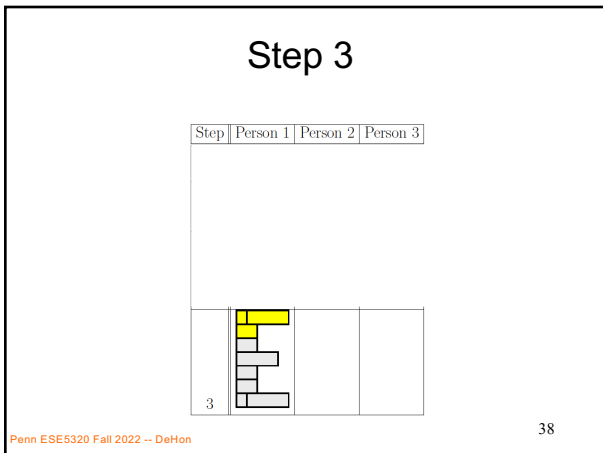
## Step 2



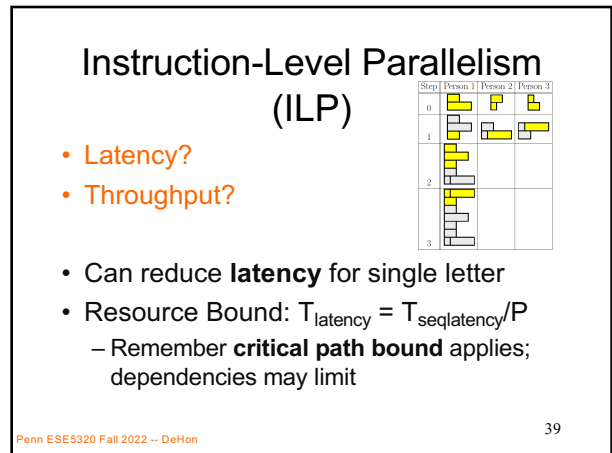
Penn ESE5320 Fall 2022 -- DeHon

37

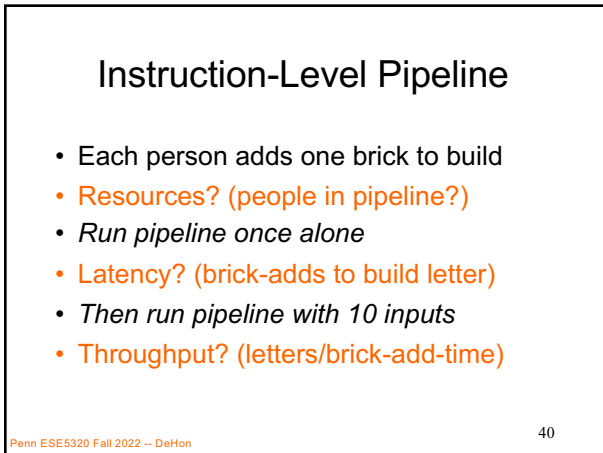
37



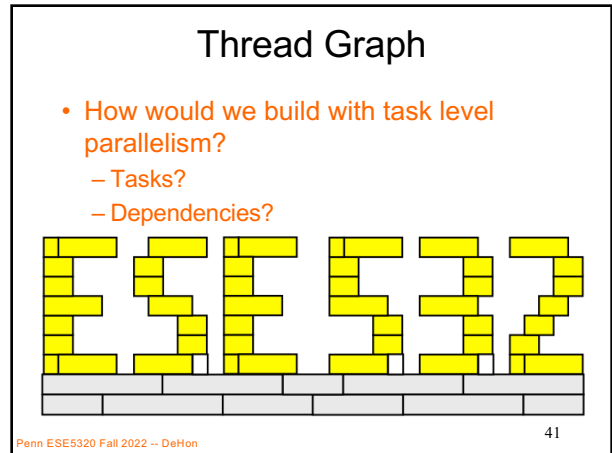
38



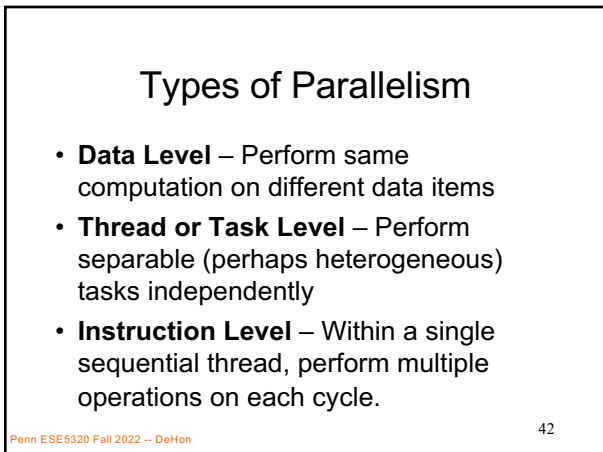
39



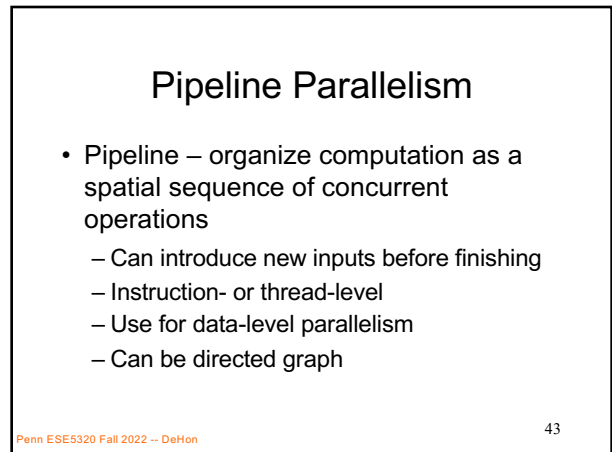
40



41



42



43

## Big Ideas

- Many parallel compute models
  - Sequential, Dataflow, CSP
- Find natural parallelism in problem
- Mix-and-match
- Likely to need all of them at some point

## Admin

- Board Holder: Board pickup if didn't get before lecture
- Reading Day 5 on web
- HW2 due Friday
- HW3 out
  - Including partner assignments on canvas
  - Board Holder reach out to partner ASAP