# ESE532:
# System-on-a-Chip Architecture

Day 4: September 11, 2019

Parallelism Overview

Pickup:
1 Preclass
1 Lego instructions
1 feecback
1 bag of legos

---

## Today

- Compute Models
  - How do we express and reason about parallel execution freedom
- Types of Parallelism
  - How can we slice up and think about parallelism?

2

---

## Message

- Many useful models for parallelism
  - Help conceptualize
- One-size does not fill all
  - Match to problem

3

---

## Parallel Compute Models

4

---

## 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)

5

---

## Parallelism can be explicit

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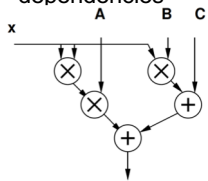- Multiply, add for quadratic equation

| cycle | mpy | add |
|---|---|---|
| 1 | B,x | |
| 2 | x,x | (Bx)+C |
| 3 | A,$x^2$ | |
| 4 | | A$x^2$+(Bx+C) |

6

---

1

## 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

7

## Implicit Parallelism

- $d=(x1-x2)*(x1-x2) + (y1-y2)*(y1-y2)$

- What parallelism exists here?

8

## 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?

9

## Term: Operation

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

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

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
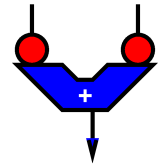
12

## Token Examples?

- How ethernet know when a packet shows up?
  - Versus when no packets are arriving?
- How serial link know character present?
- How signal miss in processor data cache and processor needs to wait for data?
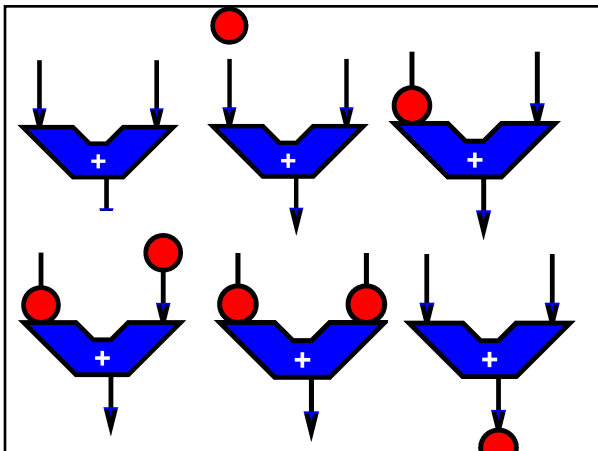
13

## Operation

- Takes in one or more inputs
- Computes on the inputs
- Produces results

- Logically **self-timed**
  - "Fires" only when input set present
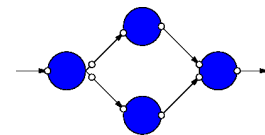  - Signals availability of output

14

## Dataflow Graph
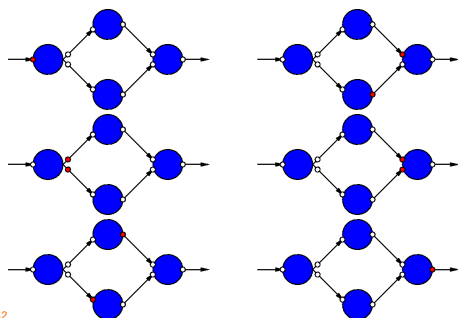
- Represents
  - computation sub-blocks
  - linkage
- Abstractly
  - controlled by data presence

16

## Dataflow Graph Example

17

## Sequential / FSM
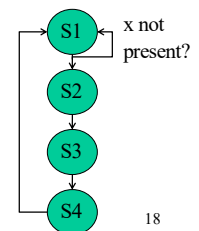
- FSM is degenerate dataflow graph where there is exactly one token

| cycle | mpy | add | next |
|-------|-----|-----|------|
| S1 | B,x | | x-->S2, else S1 |
| S2 | x,x | (Bx)+C | S3 |
| S3 | A,$x^2$ | | S4 |
| S4 | | A$x^2$+(Bx+C) | S1 |

x not present?

S1
S2
S3
S4

18

3

## Sequential / FSM

- FSM is degenerate dataflow graph where there is exactly one token

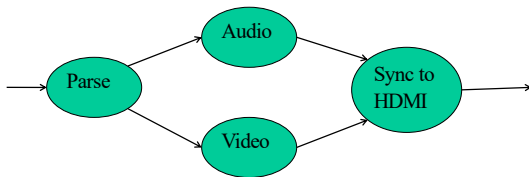| cycle | mpy | add | next |
|---|---|---|---|
| S1 | B,x | | x-->S2, else S1 |
| S2 | x,x | (Bx)+C | S3 |
| S3 | A,$x^2$ | | S4 |
| S4 | | $Ax^2+(Bx+C)$ | S1 |

S1 → S2 → S3 → S4

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
- CSP – Communicating Sequential Processes → canonical model example

20

## Video Decode

Parse → Audio, Video → Sync to HDMI →
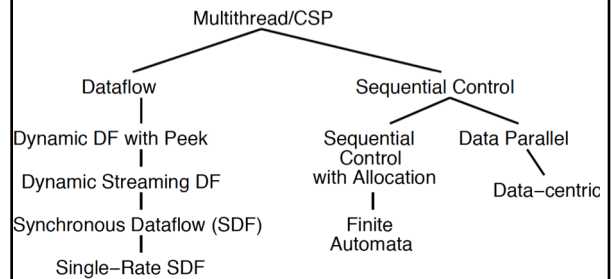
- Why might need to synchronize to send to HDMI?

21

## Compute Models

Multithread/CSP

Dataflow

Dynamic DF with Peek

Dynamic Streaming DF

Synchronous Dataflow (SDF)

Single–Rate SDF

Sequential Control

Sequential Control with Allocation

Finite Automata

Data Parallel

Data–centric

22

## 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

23

## Types of Parallelism

24

4

## 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.

25

## 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
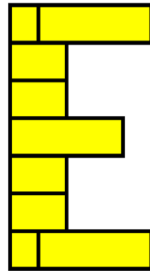  - Can be directed graph

26

**Build 1**

## Sequential

- Single person build E
- Latency?
- Throughput?
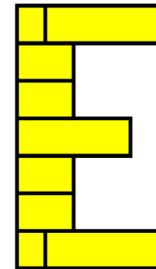
27

## Data Parallel

- Everyone in class build own E
- Latency?
- Throughput?

- Ideal speedup?
- Resource Bound?
  - 100 Es, 12 people
- When useful?

28

## Data-Level Parallelism

- **Data Level** – Perform same computation on different data items

- Ideal: $T_{dp} = T_{seq}/P$
- (with enough independent problems, match our resource bound computation)

29

**Build 2**

## Thread Parallel

- Each person build indicated letter
- Latency?
- Throughput?
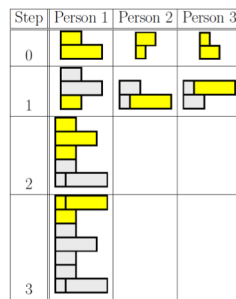- Speedup over sequential build of 6 letters?

30

## Thread-Level Parallelism

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

31

---

## Instruction-Level Parallelism

- Build single letter in lock step
- Groups of 3
- Resource Bound for 3 people building 9-brick letter?
- Announce steps from slide
  - Stay in step with slides

32

---

## Group Communication

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

| Step | Person 1 | Person 2 | Person 3 |
|------|----------|----------|----------|
| 0 | | | |
| 1 | | | |
| 2 | | | |
| 3 | | | |

33

---

## Step 0

34

---

## Step 1

35

---

## Step 2

36

---

## Step 3

37

## Instruction-Level Parallelism (ILP)

- Latency?
- Throughput?

- Can reduce **latency** for single letter
- Ideal: $T_{latency} = T_{seqlatency}/P$
  - …but **critical path bound** applies, dependencies may limit

38

---

**Build 4**

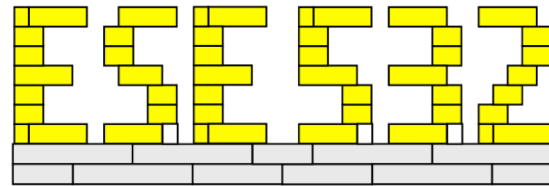## Bonus (time permit): Instruction-Level Pipeline

- Each person adds one brick to build
- Resources? (people in pipeline?)
- Run pipeline once alone
- Latency? (brick-adds to build letter)
- Then run pipeline with 5 inputs
- Throughput? (letters/brick-add-time)

39

## Thread Graph

- How would we build with task level parallelism?
  - Tasks?
  - Dependencies?

40

---

## 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.

41

## 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

42

---

## Big Ideas

- Many parallel compute models
  - Sequential, Dataflow, CSP
- Find natural parallelism in problem
- Mix-and-match

## Admin

- Reading Day 5 on web
- HW2 due Friday
- HW3 out

- Return Legos ☺
- Recitation in here at noon
  - Will take questions after class in hall