

# ESE5320: System-on-a-Chip Architecture

Day 26: December 2, 2024  
Wrapup



Penn ESE5320 Fall 2024 -- DeHon

1

## Today

- Part 1:
  - What was course about
  - Final
- Part 2: Review w/ Simple Models
- Part 3
  - Other Courses
  - Questions/Discussion

Penn ESE5320 Fall 2024 -- DeHon

2

2

## Goal

- How to design/select/map to SoC to reduce ~~Energy~~ Area/Delay.

Penn ESE5320 Fall 2024 -- DeHon

3

3

Day 1

## Outcomes

- Design, optimize, and program a modern System-on-a-Chip.
- Analyze, identify bottlenecks, design-space
- Decompose into parallel components
- Characterize and develop real-time solutions
- Implement both hardware and software solutions
- Formulate hardware/software tradeoffs, and perform hardware/software codesign

Penn ESE5320 Fall 2024 -- DeHon

4

4

Day 1

## Outcomes

- Understand the system on a chip from gates to application software, including:
  - on-chip memories and communication networks, I/O interfacing, RTL design of accelerators, processors, firmware and OS/infrastructure software.
- Understand and estimate key design metrics and requirements including:
  - area, latency, throughput, ~~energy, power,~~ predictability, and ~~reliability.~~

Penn ESE5320 Fall 2024 -- DeHon

5

5

Day 2

## 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 ESE5320 Fall 2024 -- DeHon

6

6

## Abstract Approach

- Identify requirements, bottlenecks
- Decompose Parallel Opportunities
  - At extreme, how parallel could make it?
  - What forms of parallelism exist?
    - Thread-level, data parallel, instruction-level
- Design space of mapping
  - Choices of where to map, area-time tradeoffs
- Map, analyze, refine

## SoC Designer Hardware Building Blocks

- Computational blocks
  - Adders, multipliers, dividers, ALUs
- Data Storage
  - Registers
  - Memory blocks
- Communication blocks
  - Busses
  - Multiplexers, Crossbars
  - DMA Engines
- Processors
- I/O blocks
  - Ethernet, USB, PCI, DDR, Gigabit serial links

## Final

## Final. From Code

- Analysis
  - Bottleneck
  - Amdahl's Law Speedup
  - Computational requirements
  - Resource Bounds
  - Critical Path
  - Latency/throughput/II
- Model/estimate speedup
- Forms of Parallelism
- Dataflow, SIMD, **VLW**, hardware pipeline, threads, **reduce**
- Map/schedule task graph to (multiple) target substrates
- Memory assignment and movement
- Area-time points

## Final

- Like midterm
  - Content
  - Style
  - No book, notes
  - Calculators allowed (encouraged)
- Previous years final and solutions
- Friday, December 13, 3-5pm, DRLB 4
- 2 hours (standard final exam period)

## Review

### Part 2

## Sequential Computation

- Computation requires a collection of operations
  - Arithmetic
  - Logical
  - Data storage/retrieval

$$T = \sum T_{op_i}$$

Penn ESE5320 Fall 2024 -- DeHon

13

13

## Computations in C

- Express computations in a programming language, such as C
- Can execute computation in many different ways
  - Sequential, parallel, spatial hardware

```
while(true) {
    for (i=0;i<16;i++)
        y[i]=a[i]*b[i];
    z[0]=y[0];
    for (i=1;i<16;i++)
        z[i]=z[i-1]+y[i];
}
```

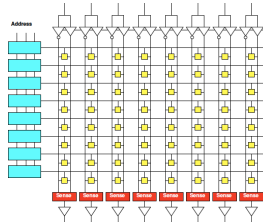
Penn ESE5320 Fall 2024 -- DeHon

14

14

## Memory Characteristics

- Small memories
  - Fast
  - Low energy
  - Not dense
    - (high area per bit)
- Large memories
  - Slow
  - High energy
  - Dense (less area per bit)



Penn ESE5320 Fall 2024 -- DeHon

15

15

## Memory and Compute

- Computation involves both arith/logical ops and memory ops

$$T = \sum T_{op_i}$$

$$T = \sum (op == mem) \times T_{mem} + \sum (op == alu) \times T_{alu}$$

$$T = N_{memop} \times T_{mem} + N_{alu} \times T_{alu}$$

Penn ESE5320 Fall 2024 -- DeHon

16

16

## Memory and Compute

- Computation involves both arith/logical ops and memory ops
- Either can dominate
  - Be bottleneck

$$T = \sum T_{op_i}$$

$$T = N_{memop} \times T_{mem} + N_{alu} \times T_{alu}$$

Penn ESE5320 Fall 2024 -- DeHon

17

17

## Memory and Compute

- Timing
  - $T_{mem}=10$
  - $T_{alu}=1$
- $N_{mpy} =$
- $N_{add} =$
- $N_{mem} =$
- Total Time, T?

```
while(true) {
    for (i=0;i<16;i++)
        y[i]=a[i]*b[i];
    z[0]=y[0];
    for (i=1;i<16;i++)
        z[i]=z[i-1]+y[i];
}
```

$$T = N_{memop} \times T_{mem} + N_{alu} \times T_{alu}$$

(not count loop, indexing costs in this sequence)

Penn ESE5320 Fall 2024 -- DeHon

18

18

## Memory and Compute

- Where bottleneck?

- $T_{mem}=10$
- $T_{alu}=1$
- $N_{mpy}=16$
- $N_{add}=15$
- $N_{mem}=95$

```
while(true) {
    for (i=0;i<16;i++)
        y[i]=a[i]*b[i];
    z[0]=y[0];
    for (i=1;i<16;i++)
        z[i]=z[i-1]+y[i];
}
```

$$T = N_{memop} \times T_{mem} + N_{alu} \times T_{alu}$$

## Data Reuse

- Reduce memory operations to large memory by storing in
  - Registers
  - Small memories

$$T = N_{memop} \times T_{mem} + N_{alu} \times T_{alu}$$

$$T = N_{smem} \times T_{smem} + N_{lmem} \times T_{lmem} + N_{alu} \times T_{alu}$$

$$T_{lmem} > T_{smem}$$

## Memory and Compute

- a,b in large mem
- y,z in small mem
  - $T_{lmem}=10$
  - $T_{smem}=1$
  - $T_{alu}=1$
- $N_{mpy}=16$
- $N_{add}=15$
- $N_{lmem}=?$
- $N_{smem}=?$

```
while(true) {
    for (i=0;i<16;i++)
        y[i]=a[i]*b[i];
    z[0]=y[0];
    for (i=1;i<16;i++)
        z[i]=z[i-1]+y[i];
}
```

$$T = N_{smem} \times T_{smem} + N_{lmem} \times T_{lmem} + N_{alu} \times T_{alu}$$

## Memory and Compute

- a,b in large mem
- y,z in small mem
  - $T_{lmem}=10$
  - $T_{smem}=1$
  - $T_{alu}=1$
- $N_{mpy}=16$
- $N_{add}=15$
- $N_{lmem}=32$
- $N_{smem}=63$

```
while(true) {
    for (i=0;i<16;i++)
        y[i]=a[i]*b[i];
    z[0]=y[0];
    for (i=1;i<16;i++)
        z[i]=z[i-1]+y[i];
}
```

- Performance impact?

$$T = N_{smem} \times T_{smem} + N_{lmem} \times T_{lmem} + N_{alu} \times T_{alu}$$

## Task Parallel

- Can run unrelated tasks concurrently

$$T = \sum T_{op_i}$$

$$T = \max \left( \sum T_{op_i} (op_i \in Task_1), \sum T_{op_i} (op_i \in Task_2) \right)$$

$$\text{Ideal: } T(p) = T(1)/p$$

## Data Parallel

- Can run same task on independent data in parallel

$$T = \sum T_{op_i}$$

$$\text{Ideal: } T(p) = T(1)/p$$

## Data Parallel

- Classify loops?
  - Sequential
  - Data parallel
  - Reduce
  - Parallel prefix

```
while(true) {
  for (i=0;i<16;i++)
    y[i]=a[i]*b[i];
  z[0]=y[0];
  for (i=1;i<16;i++)
    z[i]=z[i-1]+y[i];
}
```

Penn ESE5320 Fall 2024 -- DeHon

25

25

## Reduce

- (note slightly different second loop)
- Common pattern is a reduce operation
  - Combine a set of data into a single value
- Latency bound for second loop?

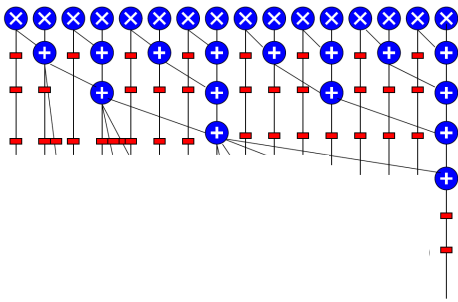
```
while(true) {
  for (i=0;i<16;i++)
    y[i]=a[i]*b[i];
  z=0;
  for (i=0;i<16;i++)
    z+=y[i];
}
```

Penn ESE5320 Fall 2024 -- DeHon

26

26

## Log-Depth Associative Reduce



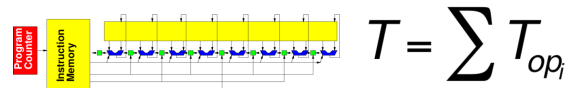
Penn ESE5320 Fall 2024 -- DeHon

27

27

## Vector/SIMD

- Can perform same operation on a set of data items



Ideal:  $T(VL) = T(1)/VL$   
(Vector Length)

Penn ESE5320 Fall 2024 -- DeHon

28

28

## Vector/SIMD

- Can perform same operation on a set of data items
  - Not everything vectorizable

$$T = \sum T_{op_i}$$

$$T = \left(\frac{1}{VL}\right) \sum T_{op_i}(\text{vectorize}(op_i)) + \sum T_{op_i}(\overline{\text{vectorize}(op_i)})$$

Penn ESE5320 Fall 2024 -- DeHon

29

29

## Vector/SIMD

- What's vectorizable?
- Speedup vector piece VL=4
  - (assume both memory and compute speedup)
- Overall speedup
- Amdahl's Law speedup for VL → infinity?

```
while(true) {
  for (i=0;i<16;i++)
    y[i]=a[i]*b[i];
  z[0]=y[0];
  for (i=1;i<16;i++)
    z[i]=z[i-1]+y[i];
}
```

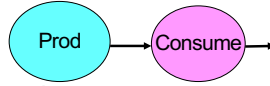
414 total cycles  
First loop: 352  
Second loop: 62

Penn ESE5320 Fall 2024 -- DeHon

30

30

## Dataflow Pipeline



- With no cycles in flowgraph
  - (no feedback)
  - (no loop carried dependencies)
- Producer and consumer can operate concurrently

$$T = \sum T_{op_i}$$

$$T = \max\left(\sum T_{op_i} (op_i \in Prod), \sum T_{op_i} (op_i \in Consume)\right)$$

Penn ESE5320 Fall 2024 -- DeHon

31

31

## Data Flow Pipeline

- **Producer/consumer here?**
- **Time each**
  - Assuming  $T_{mem}=1$
- **Granularity of dataflow pipeline**
- **Size of buffer needed to allow concurrent operation?**
- **Time on 2 processors dataflow**

```
while(true) {
    for (i=0;i<16;i++)
        y[i]=a[i]*b[i];
    z[0]=y[0];
    for (i=1;i<16;i++)
        z[i]=z[i-1]+y[i];
}
```

Penn ESE5320 Fall 2024 -- DeHon

32

32

## Data Flow Pipeline

- (note changed loops)
- **Granularity of dataflow pipeline**
- **Size of buffer needed to allow concurrent operation?**

```
for (j=0;j<MAX;j++) {
    for (i=0;i<16;i++)
        y[i]=a[j,i]*b[j,i];
    z[j,15]=y[j,15];
    for (i=14;i>=0;i--)
        z[j,i]=z[j,i+1]+y[i];
}
```

Penn ESE5320 Fall 2024 -- DeHon

33

33

## Spatial Pipeline

- Can build spatial pipeline of hardware operators to compute a dataflow graph

$$T = \sum T_{op_i}$$

- With no feedback:  $T_{loop}=1$
- With feedback limit II:  $T_{loop}=II_{cycle}$

Penn ESE5320 Fall 2024 -- DeHon

34

34

## Initiation Interval

- **What's II?**
  - a, b now  $T_{mem}=1$
  - Separate memory banks for a,b,z
  - y in local registers
  - z[i-1] communicated through local register

```
while(true) {
    for (i=0;i<16;i++)
        y[i]=a[i]*b[i];
    z[0]=y[0];
    for (i=1;i<16;i++)
        z[i]=z[i-1]+y[i];
}
```

Penn ESE5320 Fall 2024 -- DeHon

35

35

## Initiation Interval

- **What's II?**
  - a, b, c  $T_{mem}=1$
  - Separate memory banks for a,b,c,z
  - y in local registers
  - z[i-1] communicated through local register

```
while(true) {
    for (i=0;i<16;i++)
        y[i]=a[i]*b[i]+c[i];
    z[0]=y[0];
    for (i=1;i<16;i++)
        z[i]=z[i-1]+y[i];
}
```

Penn ESE5320 Fall 2024 -- DeHon

36

36

## Initiation Interval

- First loop ( $T_{mem}=1$ )
    - Latency of loop body?
    - II?
- ```

while(true) {
    g=0;
    for (i=0;i<16;i++) {
        t=a[i]*b[i]+c[g];
        y[i]=t;
        g=t%16; }
    z[0]=y[0];
    for (i=1;i<16;i++)
        z[i]=z[i-1]+y[i];
}
    
```

Penn ESE5320 Fall 2024 -- DeHon

37

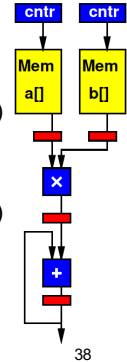
37

## Spatial Pipeline

- Back to original code
- Pipeline taking one  $a[i], b[i]$  per cycle
  - $T_{mem}=1$
- Cycles to complete?

```

while(true) {
    for (i=0;i<16;i++)
        y[i]=a[i]*b[i];
    z[0]=y[0];
    for (i=1;i<16;i++)
        z[i]=z[i-1]+y[i];
}
    
```



Penn ESE5320 Fall 2024 -- DeHon

38

38

## Critical Path

- Critical Path assuming associativity holds?
- ```

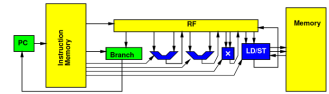
while(true) {
    for (i=0;i<16;i++)
        y[i]=a[i]*b[i];
    z[0]=y[0];
    for (i=1;i<16;i++)
        z[i]=z[i-1]+y[i];
}
    
```

Penn ESE5320 Fall 2024 -- DeHon

39

39

## VLIW



- Can control datapath to perform multiple, heterogeneous operations per cycle
  - Tune parallelism
- Op types: A, B, C
  - Ops  $N_A, N_B, N_C$
  - Number of Hardware Units  $H_A, H_B, H_C$
- $T_{RB} = \max(N_A/H_A, N_B/H_B, N_C/H_C, \dots) \leq T_{VLIW}$
- $T_{CP} \leq T_{VLIW}$

$$T = \sum T_{op_i}$$

Penn ESE5320 Fall 2024 -- DeHon

40

40

## VLIW

- VLIW
    - 1 load/store (a,b,z)
      - Assume single cycle
    - 1 mpy
    - 1 add
    - (ignoring increment for simplicity/consistency)
      - $y/z[i-1]$  in registers
  - RB
  - CP
- ```

while(true) {
    for (i=0;i<16;i++)
        y[i]=a[i]*b[i];
    z[0]=y[0];
    for (i=1;i<16;i++)
        z[i]=z[i-1]+y[i];
}
    
```

Penn ESE5320 Fall 2024 -- DeHon

41

41

## VLIW

- VLIW
    - 1 load/store (a,b,z)
      - Assume single cycle
    - 1 mpy
    - 1 add
    - 1 incr
- ```

while(true) {
    for (i=0;i<16;i++)
        y[i]=a[i]*b[i];
    z[0]=y[0];
    for (i=1;i<16;i++)
        z[i]=z[i-1]+y[i];
}
    
```

Cycle	mpy	add	Ld/st
0	a*b 0		a1
1		+y[i] 0	b1
2			Z0

Penn ESE5320 Fall 2024 -- DeHon

42

42

## VLIW

- VLIW
  - 4 load/store (a,b,z)
    - Assume single cycle
  - 2 mpy
  - 2 add
- RB
- CP

```
while(true) {
    for (i=0;i<16;i++)
        y[i]=a[i]*b[i];
    z[0]=y[0];
    for (i=1;i<16;i++)
        z[i]=z[i-1]+y[i];
}
```

Penn ESE5320 Fall 2024 -- DeHon

43

43

## Memory Bottleneck

- Memory can end up being bottleneck

$$T = \sum T_{op_i} = T_{comp} + T_{mem}$$

$$\text{Ideal: } T(p) = T_{comp}(1)/p + T_{mem}$$

Penn ESE5320 Fall 2024 -- DeHon

44

44

## Memory Bottleneck

- Pipeline to perform one multiply per cycle
- Memory supplies one data item from large memory (a,b) every 10 cycles
- Performance?
  - (number of cycles)

```
while(true) {
    for (i=0;i<16;i++)
        y[i]=a[i]*b[i];
    z[0]=y[0];
    for (i=1;i<16;i++)
        z[i]=z[i-1]+y[i];
}
```

Penn ESE5320 Fall 2024 -- DeHon

45

45

## Memory Bottleneck

- Memory can end up being bottleneck

$$T = \sum T_{op_i} = T_{comp} + T_{mem}$$

$$\text{Ideal: } T(p) = T_{comp}(1)/p + T_{mem}/\text{MemBw}$$

Penn ESE5320 Fall 2024 -- DeHon

46

46

## Memory Bottleneck

- Invest in higher bandwidth
  - Wider memory
  - More memory banks
- Exploit data reuse
  - Smaller memories may have more bandwidth
  - Smaller memories are additional banks
    - More bandwidth

Penn ESE5320 Fall 2024 -- DeHon

47

47

## Communication

- Once parallel, communication may be bottleneck

$$T = \sum T_{op_i} = T_{comp} + T_{mem} + T_{comm}$$

- Ideal (like VLIW):
 
$$T \leq \max(T_{comp}/p, T_{mem}/\text{membw}, T_{comm}/\text{netbw})$$

Penn ESE5320 Fall 2024 -- DeHon

48

48



## Multi-Objective Optimization

- Many forms of parallelism
- Given fixed area (resources, energy)
  - Maximize performance
  - Select best architecture
- Given fixed performance goal (Real Time)
  - Find architecture that achieves
  - With minimum area (energy, cost)

## Other Courses Wrapup, Questions

### Part 3

## Distinction

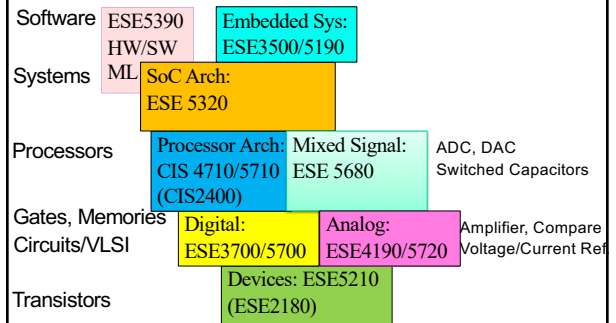
### CIS2400, 5710(4710)

- Best Effort Computing
  - Run as fast as you can
- Binary compatible
- ISA separation
- Shared memory parallelism
- Caching – automatic memory management
- Superscalar
- Pipelined processor

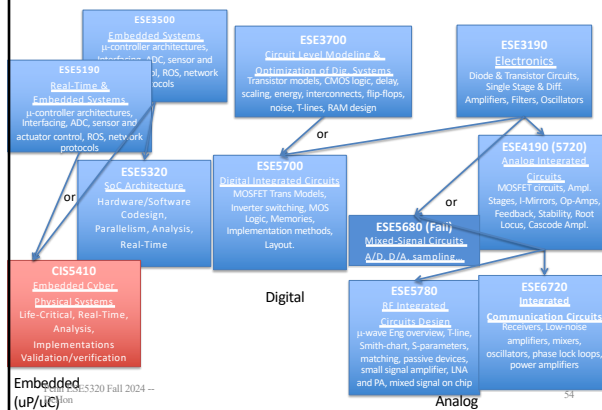
### ESE5320

- Hardware-Software codesign
  - Willing to recompile, maybe rewrite code
  - Define/refine hardware
- Real-Time
  - Guarantee meet deadline
- Non shared-memory models
- Explicit memory management
- VLIW

## Abstraction Stack



## Circuits and Computer Engineering



## Other Courses

- Security: CIS3310, CIS5510
  - ESE5370 (this Spring): Hardware Security
- Networking: ESE4070, CIS5530
- GPGPU (graphics focus): CIS5650
- Machine Learning: ESE5390
- Power Electronics: ESE5800

## Message

- Any interesting, challenging computation will require both hardware and software
- SoC powerful implementation platform
  - Target pre-existing
  - Design customized for problem
  - Exploit heterogeneous parallelism
- Understand and systematically remove bottlenecks
  - Compute, memory, communicate, I/O

Penn ESE5320 Fall 2024 -- DeHon

56

56

## Admin

- Feedback
- Project
  - Demo Day Monday, Dec. 9<sup>th</sup>
    - With board return
    - Signup for slot
  - Reports due Dec. 9<sup>th</sup>
- No DQ today
- Review: watch Ed Discuss
- Final: Friday, Dec. 13, 3pm, DRLB A4

Penn ESE5320 Fall 2024 -- DeHon

57

57

## Question/Discussion

Penn ESE5320 Fall 2024 -- DeHon

58

58