

ESE5320: System-on-a-Chip Architecture

Day 25: November 25, 2024
Real Time



Penn ESE5320 Fall 2024 -- DeHon

1

Today

Real Time

- Part 1: Demands
- Part 2: Challenges
 - Algorithms
 - Architecture
- Part 3: Disciplines to achieve

Penn ESE5320 Fall 2024 -- DeHon

2

2

Message

- Real-Time applications demand different discipline from best-effort tasks
- Look more like synchronous circuits
- Can sequentialize, like processor
 - But must avoid/rethink typical general-purpose processor common-case optimizations

Penn ESE5320 Fall 2024 -- DeHon

3

3

Real Time

- “Real” – refers to physical time
 - Connection to Real or Physical World
- Contrast with “virtual” or “variable” time
- Handles events with absolute guarantees on timing

Penn ESE5320 Fall 2024 -- DeHon

4

4

Time Constants

- Many mechanical sense/response times are ~ 5-10ms
- Human reaction times ~20ms

Penn ESE5320 Fall 2024 -- DeHon

5

5

Breaking

- Car traveling 60 km/h
- Visibility of 30 m
- How long do you have to stop as soon as something comes into view?
 - Simple total
 - See/recognize, decide, apply breaks, breaks respond

Penn ESE5320 Fall 2024 -- DeHon

6

6

Real-Time Tasks

- What timing guarantees might you like for the following tasks?
 - Self-driving car detects an object in its path
 - Delay from object appearing to detection
 - Pacemaker stimulates your heart
 - Turn steering wheel on a drive-by-wire car
 - Delay to recognized and car turns
 - Video playback (frame to frame delay)

Penn ESE5320 Fall 2024 -- DeHon

7

7

Real-Time Guarantees

- Attention/processing within fixed interval
 - Sample new value every XX ms
 - Produce new frame every 30 ms
 - Both: schedule to act and complete action
- Bounded response time
 - Respond to keypress within 20 ms
 - Detect object within 100 ms
 - Return search results within 200 ms

Penn ESE5320 Fall 2024 -- DeHon

8

8

Computer Response

- What do these things indicate?
 - When will the computer complete the task?



https://en.wikipedia.org/wiki/File:Windows_8_%2B_10_wait_cursor.gif

<https://en.wikipedia.org/wiki/File:WaitCursor-300p.gif>

Penn ESE5320 Fall 2024 -- DeHon

9

9

Real-Time Response

- What if your car gave you a spinning wait wheel for 5 seconds when you
 - Turned the wheel?
 - Stepped on the brakes?



Penn ESE5320 Fall 2024 -- DeHon

10

10

Synchronous Circuit Model

- A simple synchronous circuit is a good “model” for real-time task
 - Run at fixed clock rate
 - Take input every “cycle” (application cycle)
 - Produce output every “cycle” (application cycle)
 - Complete computation between input and output
 - Designed to run at fixed-frequency
 - Critical path meets frequency requirement

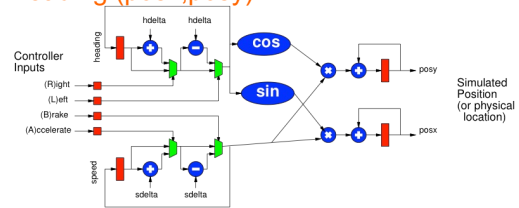
Penn ESE5320 Fall 2024 -- DeHon

11

11

Preclass 2

- (Circuit) Cycle time could operate?
- Assume clocked at 100Hz (application cycle)
- Worst-case delay from (L)eft press to change in heading (posx, posy)?



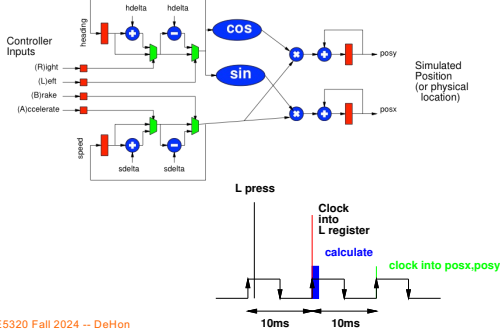
Penn ESE5320 Fall 2024 -- DeHon

12

12

Preclass 2

- Assume clocked at 100Hz (application cycle)



Penn ESE5320 Fall 2024 -- DeHon

13

13

Historically

- Real-Time concerns grew up in EE
 - Because an analog circuit was the only way could meet frequency demands
 - ...later a dedicated digital circuit...
- Applications
 - Signal processing, video, control, ...

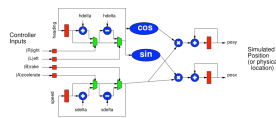
Penn ESE5320 Fall 2024 -- DeHon

14

14

Technological Change

- Area units for spatial design shown (preclass 2c)
- Fraction of processor capacity required (Preclass 2d)
- Why might prefer using a processor to using the spatial circuit?
 - Hint: What does preclass 2c,d suggest?



Penn ESE5320 Fall 2024 -- DeHon

15

15

Performance Scaling

- As circuit speeds increased
 - Can meet real-time performance demands with heavy sequentialization
- Circuit and processor clocks
 - from MHz to GHz
- Many real-time task rates unchanged
 - 44KHz audio, 33 frames/second video
- Even 100MHz processor
 - Can implement audio in a small fraction of its computational throughput capacity

Penn ESE5320 Fall 2024 -- DeHon

16

16

HW/SW Co-Design

- Computer Engineers – know can implement anything as hardware or software
- Want freedom to move between hardware and software to meet requirements
 - Performance, costs, energy

Penn ESE5320 Fall 2024 -- DeHon

17

17

Real-Time Challenge

- Meet real-time demands / guarantees
 - Economically using programmable architectures
- Sequentialize and share resources with deterministic, guaranteed timing
- Spatial (all hardware, HLS synthesized) implementations are good at meeting real-time guarantees, but may be bigger than necessary

Penn ESE5320 Fall 2024 -- DeHon

18

18

Part 2

CHALLENGES

Penn ESE5320 Fall 2024 -- DeHon

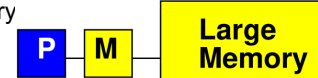
19

19

Day 3

Processor Data Caches

- Traditional Processor Data Caches are a heuristic instance of this
 - Add a small memory local to the processor
 - It is fast, low latency
 - Store anything fetched from large/remote memory in local memory
 - Hoping for reuse in near future
 - On every fetch, check local memory before go to large memory



Penn ESE5320 Fall 2024 -- DeHon

20

Day 3

Processor Data Caches

- Demands more than a small memory
 - Need to sparsely store address/data mappings from large memory
 - Makes more area/delay/energy expensive than just a simple memory of capacity
- Don't need explicit data movement
- Cannot control when data moved/saved
 - Bad for determinism
- Limited ability to control what stays in small memory simultaneously

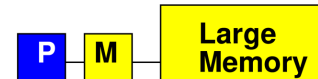
Penn ESE5320 Fall 2024 -- DeHon

21

21

Processor Data Caches

- Traditional Processor Data Caches are a heuristic instance of this
 - Store anything fetched from large/remote memory in local memory
 - Hoping for reuse in near future
 - On every fetch, check local memory before go to large memory
 - Stall processor while waiting for data



Penn ESE5320 Fall 2024 -- DeHon

22

Preclass 3:

Processor Cache Timing

- Assume
 - cache miss (go to large memory) takes 10 cycles
 - Cache hit (small memory) takes 1
 - Start with empty cache
- Due to memory delay, **how long to execute:**

$b = a[0] + a[1];$	$b = a[i] + a[j];$
$c = a[1] + a[2];$	$c = a[k] + a[l];$
$d = a[2] + a[0];$	$d = a[m] + a[n];$

Penn ESE5320 Fall 2024 -- DeHon

23

23

Scratchpad

- Recall, scratchpad memory
 - Small
 - Explicitly managed (not dynamic like cache)
- If move (DMA) data to scratchpad memory, would be deterministic

$b = a[0] + a[1];$	$b = a[i] + a[j];$
$c = a[1] + a[2];$	$c = a[k] + a[l];$
$d = a[2] + a[0];$	$d = a[m] + a[n];$

Penn ESE5320 Fall 2024 -- DeHon

24

24

Observe

- Instructions on “General Purpose” processors take variable number of cycles

Penn ESE5320 Fall 2024 -- DeHon

25

25

Preclass 4

- How many cycles?

- sin, cos 100 cycles each
- Assignments take 1 cycle

```
old_sh=sh; old_ch=ch;
if (!left || !right)
    {sh=old_sh;ch=old_ch;}
else
    {sh=sin(heading);
    ch=cos(heading);}
```

Penn ESE5320 Fall 2024 -- DeHon

26

26

Preclass 5

- How many cycles?

```
sum=0;
for (i=0;i<32;i++) {
    sum+=(0-(b%2)) & a;
    b=b>>1;
    a=a<<1;
}
```

Penn ESE5320 Fall 2024 -- DeHon

27

27

Preclass 5

- How many cycles?

```
sum=0;
for (;b!=0;b=b>>1) {
    if (b%2==1)
        sum+=a;
    a=a<<1;
}
```

Penn ESE5320 Fall 2024 -- DeHon

28

28

Observe

- Data-dependent branching, looping
 - Means variable time for operations

Penn ESE5320 Fall 2024 -- DeHon

29

29

Two Challenges

1. Architecture – Hardware have variable (data-dependent) delay
 - Esp. for General-Purpose processors
 - Instructions take different number of cycles
2. Algorithm – computational specification have variable (data-dependent) operations
 - Different number of instructions

$$Time = \sum_i Cycles(i)$$

Penn ESE5320 Fall 2024 -- DeHon

30

30

Algorithm

- What programming constructs are data-dependent (variable delay)?

Penn ESE5320 Fall 2024 -- DeHon

31

31

Programming Constructs

- Conditionals: if/then/else
- Loops without compile-time determined bounds
 - While with termination expressions
 - For with data-dependent bounds
- Data-dependent recursion
- Interrupts
 - I/O events, time-slice
- Note: 1st three were issue for HLS
 - For same reason – how did we address?

Penn ESE5320 Fall 2024 -- DeHon

32

32

Hardware Architecture

- Some typical (4710,5710) processor “optimizations” can cause variable delay
 - Caches
 - Branch prediction
 - Common-case optimizations
 - Pipeline stalls
 - Speculative issue

Penn ESE5320 Fall 2024 -- DeHon

34

34

Part 3

DISCIPLINES TO ACHIEVE REAL-TIME

Penn ESE5320 Fall 2024 -- DeHon

35

35

Already Addressed

- Hardware pipelines are deterministic
 - HLS limitations for hardware already drove us to fixed timing
- Explicit scheduling of VLIW issue
 - Can be fixed timing

Penn ESE5320 Fall 2024 -- DeHon

36

36

Open Issues

- Processor Architecture
- Resource Sharing

Penn ESE5320 Fall 2024 -- DeHon

37

37

What can we do to make architecture more deterministic?

- Explicitly managed memory
- Eliminate Branching (too severe?)
- Unpipelined processors
- Fixed-delay pipelines
 - Offline-scheduled resource sharing
 - Multi-threaded
- Deadlines

Penn ESE5320 Fall 2024 -- DeHon

38

38

Explicitly Managed Memory

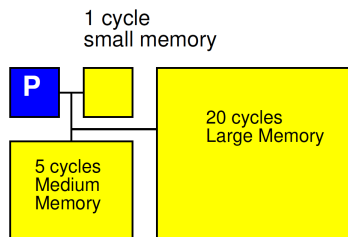
- Make memory hierarchy visible
 - Use Scratchpad memories instead of caches
- Explicitly move data between memories
 - E.g. movement into local memory
- Already do for Register File in Processor
 - Load/store between memory and RF slot
 - ...but don't do for memory hierarchy

Penn ESE5320 Fall 2024 -- DeHon

39

39

Explicitly Managed Memory



Penn ESE5320 Fall 2024 -- DeHon

40

40

Offline Schedule Resource Sharing

- Don't arbitrate
- Decide up-front when each shared resource can be used by each thread or processor
 - Simple fixed schedule
 - Detailed Schedule
- What
 - Memory bank, bus, I/O, network link, ...

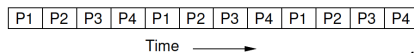
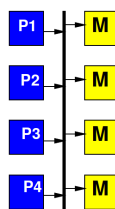
Penn ESE5320 Fall 2024 -- DeHon

41

41

Time-Multiplexed Bus

- Fixed by hardware master
- 4 masters share a bus
 - Each master gets to make a request on the bus every 4th cycle
 - If doesn't use it, goes idle



Penn ESE5320 Fall 2024 -- DeHon

42

42

Time-Multiplexed Bus

- Regular schedule
- Fixed bus slot schedule of length $N >$ masters
 - (probably a multiple)
- Assign owner for each slot
 - Can assign more slots to one
- E.g. $N=8$, for 4 masters
 - Schedule (1 2 1 3 1 2 1 4)

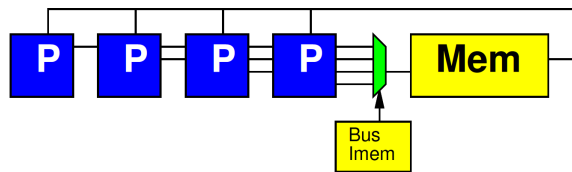
Penn ESE5320 Fall 2024 -- DeHon

43

43

Fully Scheduled

- At extreme, fully schedule which tasks gets resource on each cycle



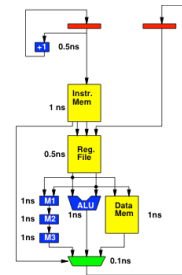
Penn ESE5320 Fall 2024 -- DeHon

44

44

Simple Deterministic Processor with Multiplier

- No branching
- Unpipelined
- Every operation completes in fixed time



- Cycle time as shown?
- Retimed cycle time?
- What's inefficient about this design?

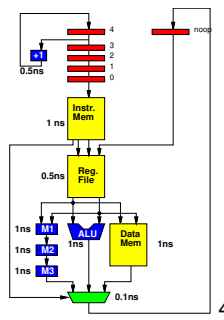
Penn ESE5320 Fall 2024 -- DeHon

45

45

Simple Deterministic Processor with some Pipelining

- No branching
- Every operation completes in fixed time
- Retimed cycle time?
 - Hint what are cycles?
- How pipelines added change behavior?
 - Hint: what is sequence of addresses into Instr. Mem?



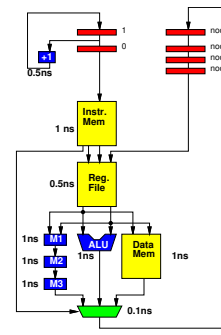
Penn ESE5320 Fall 2024 -- DeHon

46

46

Simple Deterministic Pipelined Processor

- No branching
- Every operation completes in fixed time
- How pipelines added change behavior?
 - Hint R1 value



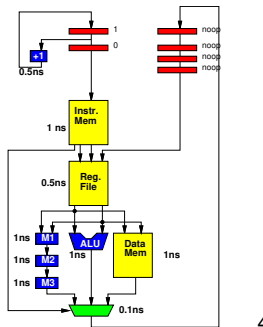
Penn ESE5320 Fall 2024 -- DeHon

47

47

Simple Deterministic Pipelined Processor

- No branching
- Every operation completes in fixed time
- Retimed cycle time?



Penn ESE5320 Fall 2024 -- DeHon

48

48

Deterministic Pipelines

- Pipeline data updates
 - Not available on next cycle, but depend on delay through logic
- Delay branches
 - Not happen immediately, a fixed number of clock cycles later
- Schedule like VLIW

Penn ESE5320 Fall 2024 -- DeHon

49

49

Deterministic Pipelines

- Not how ARM, Intel (4710, 5710) processor are pipelined
- Those include operations that make timing variable
 - dynamic data hazards, branch speculation
- Here, data becomes available after a predictable time
- Branches take effect at a fixed time
 - Likely delayed
- Schedule to delays to get correct data

Penn ESE5320 Fall 2024 -- DeHon

68

68

Deadline Instruction

- Deal with algorithmic (branching) variability
- Set a hardware counter for thread
- Decrement counter on each cycle
- Demand counter reach 0 before thread allowed to continue at deadline instruction
 - Stall if get there early
 - Similar to flip-flop on a logic path
 - Wait for clock edge to change or sample value
- Model: fixed execution time

Penn ESE5320 Fall 2024 -- DeHon

69

69

WCET

- WCET – Worst-Case Execution Time
- Analysis when working with algorithms and architectures with data-dependent delay
 - Need to meet real time
 - Calculate the worst-case runtime of a task
 - Like calculating the critical path (but harder)
 - Worst-case delay of instructions
 - Worst-case path through code
 - Worst-case # loop iterations
 - Rationale for setting Deadlines
 - (like a cycle time)

Penn ESE5320 Fall 2024 -- DeHon

70

70

Different Goals

- | Real-Time | General Purpose/Best Effort |
|--|---|
| <ul style="list-style-type: none">• Willing to recompile to new hardware• Want time on hardware predictable• Willing to schedule for delays in particular hardware | <ul style="list-style-type: none">• ISA fixed• Want to run same assembly on different implementations• Tolerate different delays for different hardware• Run faster on newer, larger implementations |

Penn ESE5320 Fall 2024 -- DeHon

71

71

SoC Opportunity

- Can choose which resources are shared
- Can dedicate resources to tasks
- Isolate real-time tasks/portions of tasks from best-effort
 - Separate hardware/processors
 - Separate memories, network

Penn ESE5320 Fall 2024 -- DeHon

72

72

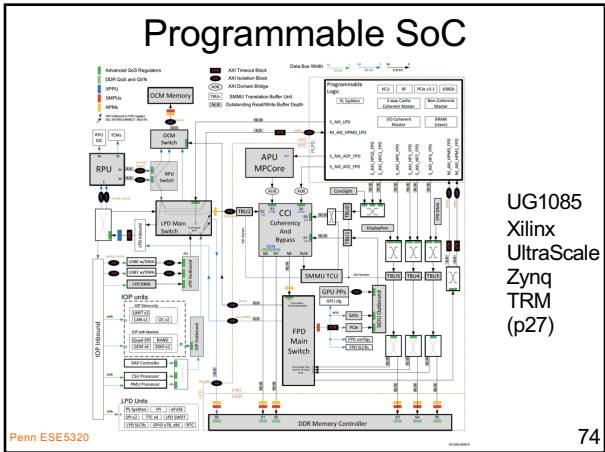
UltraScale+ Zynq

- Has 2 “Real-Time Processor”
 - ARM Cortex-R5
 - 32b (vs. 64b for A53 APU processor)
 - ARMv7-R (vs. ARMv8)
 - Single ALU, dual issue
 - Branch prediction
- Explicitly managed scratchpads
 - Tightly-Coupled Memories
 - On-Chip Memory (OCM)

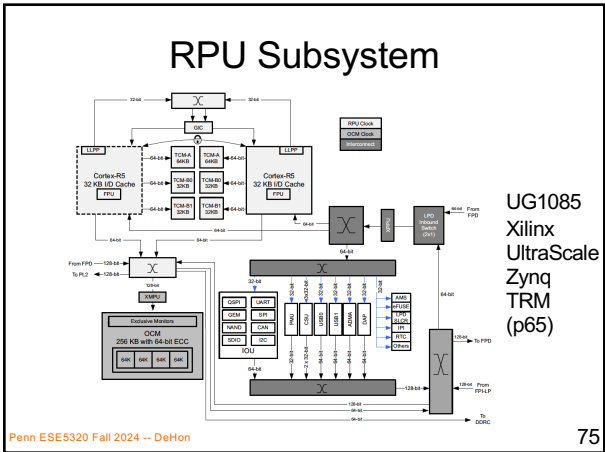
Penn ESE5320 Fall 2024 -- DeHon

73

73



UG1085
Xilinx
UltraScale
Zynq
TRM
(p27)



UG1085
Xilinx
UltraScale
Zynq
TRM
(p65)

- ### Big Ideas:
- Real-Time applications demand different discipline from best-effort tasks
 - Look more like synchronous circuits and hardware discipline
 - Avoid or use care with variable delay programming constructs
 - Can sequentialize, like processor
 - But must avoid/rethink typical processor common-case optimizations
 - Offline calculate static schedule for computation and sharing
 - Instead of dynamic arbitration, static

- ### Admin
- Feedback
 - Penn on Thursday-Friday schedule for tomorrow and Wednesday
 - → no lecture on Wednesday
 - Next lecture Monday
 - Wrapup lecture