

154B Discussion 5

February 11th, 2022

Goals

- Overview of Assignment 4.
- Some memory system stuff (AMAT, memory array calculation)

Assignment 4: Overview

- Implement forwarding logic and hazard detection logic for dual-issue pipelined CPU.
- Run full applications with/without loops unrolled on implemented CPU designs.
- Reason about the performance of different CPU designs on different workloads.

Assignment 4: Dual-issue Pipelined CPU

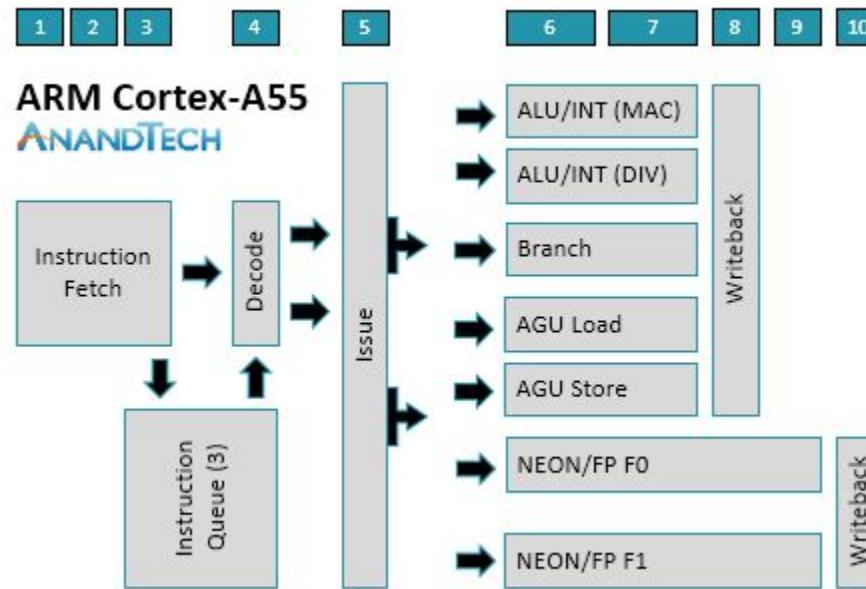
- Issuing: occurs in the Fetch stage, decides how many instructions will enter the pipeline in that cycle.
 - Different from fetching, which means loading instructions from memory.
 - Can only issue instructions fetched from memory.
- Multiple-issue: issuing multiple instructions in a cycle.
 - Dual-issue: issuing *upto* 2 instructions per cycle.

Assignment 4: Why Dual-issue?

- Further exploit the idea of instruction-level parallelism.
 - Maximizing the number of instructions could be executed at the same time.
- For DINOCPU, in terms of number of cycles need to complete a program, dual-issue design always needs fewer cycles than the original pipeline.

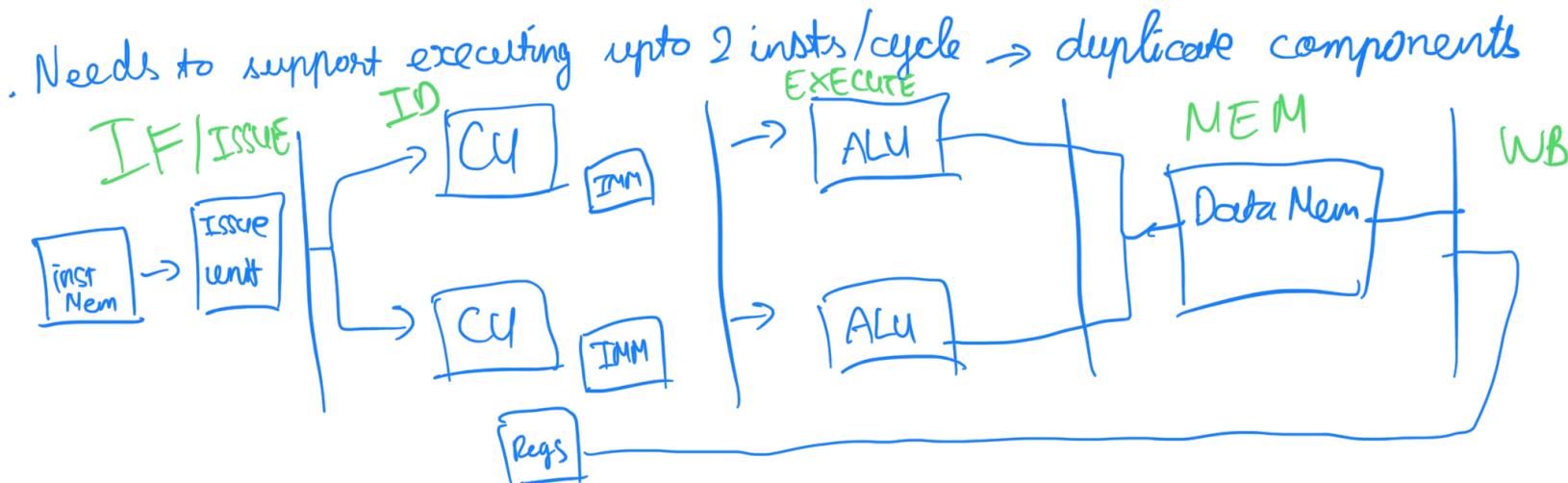
Dual-issue in Real World

- Used in ARM's low power core Cortex A55 (2017).



Assignment 4: How to implement dual-issue pipeline?

- Duplicate almost every component and all stage registers.
- Modifying forwarding unit and hazard detection unit.



Assignment 4: How fetched instructions are issued?

- There are two pipelines in dual-issue design: pipeA and pipeB.
 - If the issue unit decides to issue 2 instructions, the instruction at PC will be issued in pipeA, and the instruction at PC+4 will be issued in pipeB.
 - If there is only one instruction being issued, that instruction will enter pipeA, and the instruction in pipeB will be marked as invalid.
- Issue constraints:
 - The issue unit will only issue 2 instructions if,
 - The instruction at PC is not a load/store/jump/branch.
 - The instruction at PC will not write to a register that will be read by instruction at PC+4.
- The constraints lead to this property,
 - If the instruction at a stage in pipeA can update PC, the instruction at that stage in pipeB is invalid.

$\xrightarrow{\text{BRANCH/JUMP}}$
↳ that instruction should not be executed

Assignment 4: Loop Unrolling

- Compiler optimization option, trade-off between the binary size and potential instruction-level parallelism.
- Main idea: doing more work per iteration of a loop.

// A loop *Before being unrolled*

```
for (int i = 0; i < 16; i++)
{
    arr[i] = c[i] + 1;
}
```

// Unrolled loop *fewer iterations*

```
for (int i = 0; i < 16; i+=4)
{
    arr[i] = c[i] + 1;
    arr[i+1] = c[i+1] + 1;
    arr[i+2] = c[i+2] + 1;
    arr[i+3] = c[i+3] + 1;
}
```

more work per iteration

Assignment 4: Loop Unrolling

- Loops, in machine code.

```
// A loop  
  
label1:  
    addi a1, a1, 1      a 1 → current index  
  
    // do 1 thing  
  
    blt a1, 16, label1  ← this is the loop branch  
                         executed 16 times  
  
label2:  
...  
  
# Wasted cycles : 15 × 3 = 45 cycles
```

```
// Unrolled loop  
  
label1:  
    addi a1, a1, 4  
  
    // do 1 thing  
  
    // do 1 thing  
  
    // do 1 thing  
  
    blt a1, 16, label1  ← loop branch  
                         executed 4 times  
  
label2:  
...  
  
→ Not taken: 1 time  
→ Taken : 3 times  
  
→ #wasted cycles ; 3 × 3 = 9 cycles
```

Assignment 4: Loop Unrolling

- In the assignment, the binaries with loops unrolled only have the loops of known size unrolled.
- What happens when the size is unknown?

```
// A loop  
  
while (cond)  
{  
    // do something  
}
```

```
// Unrolled unknown size loop, machine code  
  
label1:  
    // do 1 thing  
    blt a1, a2, label1  
    // do 1 thing  
    blt a1, a2, label1  
    // do 1 thing  
    blt a1, a2, label1  
    // do 1 thing  
    blt a1, a2, label1  
  
label2:  
    ...
```

*have to check
stopping condition*

*multiple times
per iteration*

AMAT: Average Memory Access Time

- 1-level cache system

$$\text{AMAT} = \text{L1_latency} + (1 - \text{L1_cache_hit_rate}) * \text{memory_latency}$$

- 2-level cache system

$$\text{AMAT} = \text{L1_latency} + (1 - \text{L1_hit_rate}) \times [\text{L2_latency} + (1 - \text{L2_hit_rate}) \times \text{memory_latency}]$$

- 3-level cache system

$$\begin{aligned}\text{AMAT} = & \text{L1_latency} + (1 - \text{L1_hit_rate}) \times (\text{L2_latency} + (1 - \text{L2_hit_rate}) \times (\text{L3_latency} \\ & + (1 - \text{L3_hit_rate}) \times \text{memory_latency}))\end{aligned}$$

Generalized: n-level cache:

$$\text{AMAT} = \text{L1_latency} + (1 - \text{L1_hit_rate}) \times \text{AMAT}_{L2}$$

$$\text{AMAT}_{L2} = \text{L2_latency} + (1 - \text{L2_hit_rate}) \times \text{AMAT}_{L3}$$

$$\vdots$$
$$\text{AMAT}_{Ln} = \text{Ln_latency} + (1 - \text{Ln_hit_rate}) \times \text{memory_latency}$$

AMAT: Example

- 2-Level Cache System

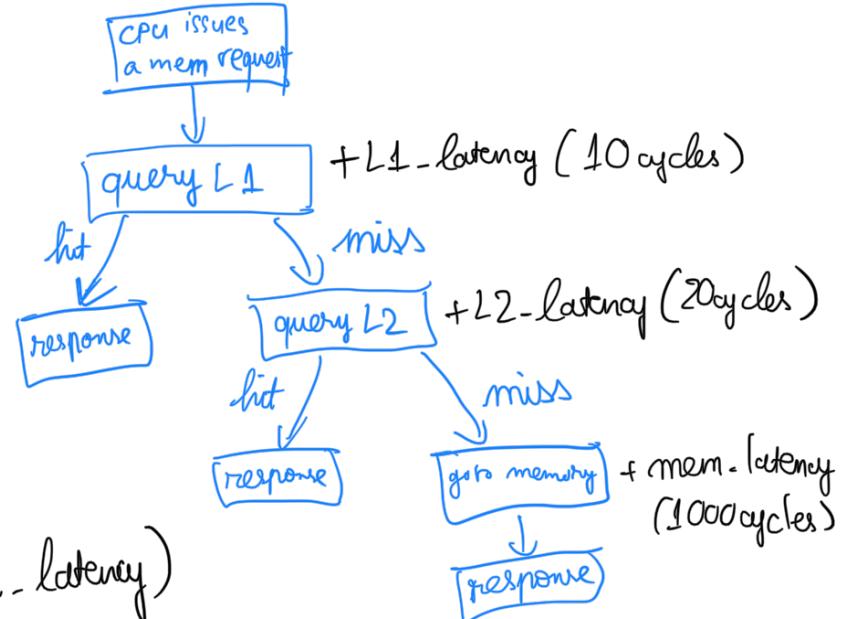
L1 hit rate: 30%; L1 latency: 10 cycles

L2 hit rate: 60%; L2 latency: 20 cycles

Memory latency: 1000 cycles

From graph:

$$\begin{aligned} \text{AMAT} &= \text{L1-latency} \\ &+ (1 - \text{L1-hit-rate}) \times (\text{L2-latency} + (1 - \text{L2-hit-rate}) \times \text{mem-latency}) \\ &= 10 + (1 - 0.3) \times (20 + (1 - 0.6) \times 1000) \\ &= 304 \text{ cycles} \end{aligned}$$



general idea:

→ accessing cache or memory costs time regardless whether it results in a hit or a miss.

→ if an access is a miss, needs to query the next level

Memory Array

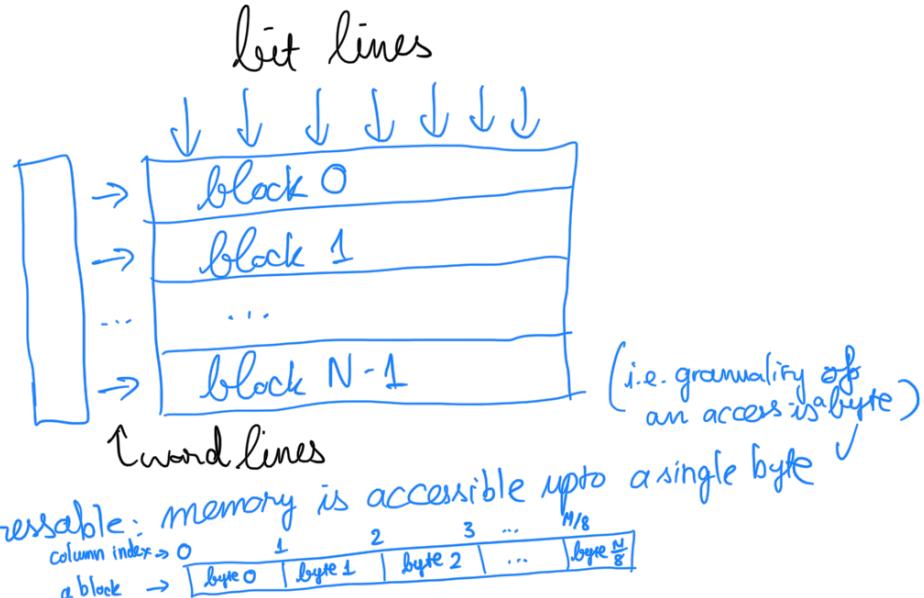
- Row: Word lines
- Column: Bit lines
- Address: (row bits) (column bits)

→ # of bits for selecting a row:
 $\log_2(\text{how many different rows}) = \log_2(N)$

→ # of bits for selecting columns assuming byte addressable:
 $\log_2(\text{how many possibilities for selecting columns})$
 $\log_2(\text{how many possibilities for selecting columns}) = \log_2(M) - \log_2(8)$
 $= \log_2\left(\frac{M}{8}\right) = \log_2(M) - 3$
 $= \log_2(M) - 3$

Mem Array:

- N different rows
- each row has 1 block of M bits
- bit lines = selecting columns
- word lines: selecting 1 row
- to get data: select a row first, then selected columns



Memory Array: Example

- Memory Array
 - Byte-addressable
 - 4096-bit word per block
 - Capacity: 262144 bits

LED question

- . refresh rate of 90 Hz \rightarrow the screen needed to be updated 90 times per second (or every 1/90 sec)
- . Watt = Joule / second

Bandwidth: $\frac{1920 \times 1080 \times 24}{8} \div 1024^2 \times 90$

of bits of framebuffer | bits | bytes | bytes \rightarrow MiB

needs to update the framebuffer

MiB needed per update

90 times / second

nJ per bit

Power (Joule / second) = $\frac{1920 \times 1080 \times 24}{8} \times 90 \times 2 \times 10^{-9}$

of bits per update | # updates / second | nJ \rightarrow J