

A) A processor has the following cache hierarchy: L1 I-cache and L1 D-cache each 8 KB, 2-way associative, 16-byte block; unified L2, 512 KB, 4-way associative, 32-byte block; unified L3 1MB, 4-way associative, 64-byte block. The latencies (disregarding virtual memory TLB) expressed in clock cycles are: 2 in L1, 3 in L2, 8 in L3.

The processor executes the IA-32 ISA (but with 32-bit addressing mode).

a1) Assuming initially empty and invalidated cache lines throughout the hierarchy, show the content of each D-cache after the execution of the instruction LD R1,(0000AFAF_{hex}) (further assumption: the instruction is in the I-cache).

a2) determine hexadecimal values X (other than 0000AFAF_{hex}) such that the instruction ST (X),R2 has the following effect

- hit L1
- miss L1, hit L2
- miss L1, miss L2, hit L3
- miss throughout the cache hierarchy.

B) The processor runs at 2.33 GHz and is connected to RAM memory through a 64-bit bus interfaced to a DDR2-1066 2-way interleaved memory, 16-byte block length. The bus/DDR2 interface guarantees the maximum transfer rate allowed by the DDR2-1066 chips. Compute the cost of a miss measured in processor clocks, assuming that addressing the memory requires a single bus clock cycle, and that each row activation in the memory banks requires 3 bus clock cycles.

C) The processor has a superscalar, 2-way pipeline, that fetches, decodes issues and retires (commits) bundles containing each 2 instructions. The front-end in-order section (fetch and decode) consists of 4 stages. The issue logic takes 1 clock cycle, if the instructions in the bundle are independent, otherwise it takes 2 clock cycles. The architecture supports dynamic speculative execution, and control dependencies from branches are solved only at commit time, even if the outcome of the branch logic is available earlier. The execution model obeys the attached state transition diagram.

There are 2 functional units (FUs) Int1-INT2 for integer arithmetics (arithmetic and local instructions, branches and jumps, no multiplication), 2 FUS FAdd1-Fadd2 for floating point addition/subtraction, a FU FMolt1 for floating point multiplication, and a FU for division, FDiv1. There are 12 integer (R0-R11) and 12 floating point (F0-F11) registers.

Speculation is handled through a 8-entry ROB, a pool of 4 Reservation Stations (RS) Rs1-4 shared among all FUs, 2 load buffers Load1-2, 2 store buffers Store1-2 (see the attached execution model): an instruction bundle is first placed in the ROB (if two entries are available), then each instruction is dispatched to one of the shared RS (if available) and then executed in the proper FU. FUs are *pipelined* (not the Fdiv one) and have the latencies quoted in the following table:

Int - 1	Fadd - 3
Fmolt - 5	Fdiv - 6

Further assumption

- caches are described in point A) and are assumed empty and invalidated.

c1) show state transitions for the instructions of the first iteration of the following code fragment (assume a conventional L3 miss time of 6 clock cycles), highlighting conflicts, if any:

```
PC01 MOVI R1,0X0000AFAF
PC02 MOVI R2,0X0000000A
PC03 LD F1,16(R1)
PC04 LD F2,8(R1)
PC05 LD F3,0(R1)
PC06 MULTF F4,F3,F0      -- F0 HAS BEEN PRE-LOADED
PC07 MULTF F5,F2,F0
PC08 MULTF F6,F1,F0
PC09 ADDF F7,F4,F5
PC10 ADDF F7,F7,F6
PC11 SD 0(R1),F7
PC12 ADDI R1,R1,8
PC13 SUBI R2,R2,1
PC14 BNEZ R2,PC03
```

c2) show ROB, RS and buffer status at the issue of the BNEZ of the first iteration.

D) The code fragment is executed on a statically scheduled pipeline having the following structure (A1-A2 FP addition/subtraction, M1-M4 FP multiplication):

```
IF1A, ID1A, EXA,                ,MEM1A, MEM2A, WBA
    A1, A2
    M1, M2, M3, M4
```

- d1) Assuming BNEZ decides on the condition in EXA, schedule the branch delay slot;
- d2) Using b1), produce a schedule for the first iteration, assumin caches always hit with a latency of 1 clock cycle;
- d3) computed the CPI of the algorithm;
- d4) evaluate maximum unrolling (if possible), with 32 registers (both INT e FP);
- d5) unroll once and schedule the unrolled code, and compute the new CPI.

Dynamic speculative execution
Decoupled ROB RS execution model

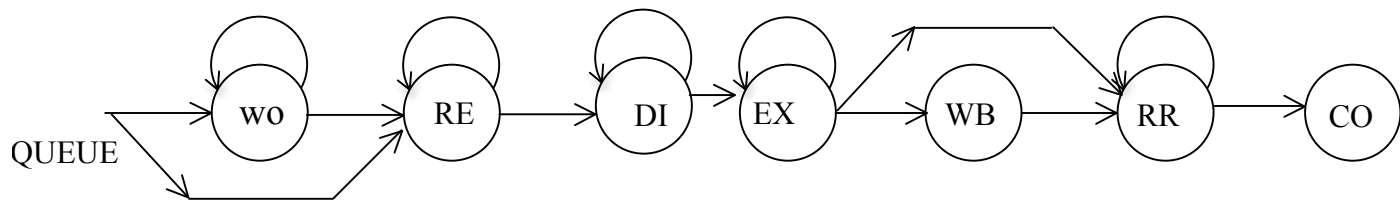
[illegible]

	Reservation station and load/store buffers							
	Busy	Op	Vj	Vk	ROB _j	ROB _k	ROB pos	Address
Rs1								
Rs2								
Rs3								
Rs4								
Load1								
Load2								
Store1								
Store2								

ROB_j ROB_k: sources not yet available
ROB pos: ROB entry number where instruction is located

	Result Register status													
Integer	R0	R1	R2	R3	R4	R5	R6	R7	R8	R9	R10	R11		
ROB pos														
state														
Float.	F0	F1	F2	F3	F4	F5	F6	F7	F8	F9	F10	F11		
ROB pos														
state														

Reorder Buffer (ROB)					
ROB Entry#	Busy	Op	Status	Destination	Value
1					
2					
3					
4					
5					
6					
7					
8					



Decoupled execution model for bundled (paired) instructions

The state diagram depicts the model for a dynamically scheduled, speculative execution microarchitecture equipped with a Reorder Buffer (ROB) and a set of Reservation Stations (RS). The RSs are allocated during the ISSUE phase, denoted as RAT (Register Alias Allocation Table) in INTEL microarchitectures, as follows: a bundle (2 instructions) if fetched from the QUEUE of decoded instructions and ISSUED if there is a free couple of consecutive entries in the ROB (head and tail of the ROB queue do not match); each *single* instruction is moved into a RS (if available) when all of its operands are available. Access memory instructions are allocated in the ROB and then moved to a load/store buffer (if available) when operands (address and data, if proper) are available .

States are labelled as follows:

WO:	Waiting for Operands (at least one of the operands is not available)
RE:	Ready for Execution (all operands are available)
DI:	Dispatched (posted to a free RS or load/store buffer)
EX:	Execution (moved to a load/store buffer or to a matching and free UF)
WB:	Write Back (result is ready and is returned to the Rob by using in exclusive mode the Common Data Bus CDB)
RR:	Ready to Retire (result available or STORE has completed)
CO:	Commit (result is copied to the final ISA register)

State transitions happen at the following events:

<i>from</i> QUEUE <i>to</i> WO:	ROB entry available, operand missing
<i>from</i> QUEUE <i>to</i> RE:	ROB entry available, all operands available
<i>loop at</i> WO:	waiting for operand(s)
<i>from</i> WO <i>to</i> RE:	all operands available
<i>loop at</i> RE:	waiting for a free RS or load/store buffer
<i>from</i> RE <i>to</i> DI:	RS or load/store buffer available
<i>loop on</i> DI:	waiting for a free UF
<i>from</i> DI <i>to</i> EX:	UF available
<i>loop at</i> EX:	multi-cycle execution in a UF, or waiting for CDB
<i>from</i> EX <i>to</i> WB:	result written to the ROB with exclusive use of CDB
<i>from</i> EX <i>to</i> RR:	STORE completed, branch evaluated
<i>loop at</i> RR:	instruction completed, not at the head of the ROB, or bundled with a not RR instruction
<i>from</i> RR <i>to</i> CO:	bundle of RR instructions at the head of the ROB, no exception raised

Resources

Register-to-Register instructions hold resources as follows:

- ROB: from state WO (or RE) up to CO, inclusive;
- RS: state DI
- UF: EX and WB

Load/Store instructions hold resources as follows:

- ROB: from state WO (or RE) up to CO, inclusive;
- Load buffer: from state WO (or RE) up to WB
- Store buffer: from state (or RE) up to EX (do not use WB)

Forwarding: a write on the CDB (WB) makes the operand available to the consumer in the same clock cycle. If the consumer is doing a state transition from QUEUE to WO or RE, that operand is made available; if the consumer is in WO, it goes to RE in the same clock cycle of WB for the producer.

Branches: they compute Next-PC and the branch condition in EX and optionally forward Next-PC to the “in-order” section of the pipeline (Fetch states) in the next clock cycle. They do not enter WB and go to RR instead.