

Lecture 9

Superscalar-processing

Stallings: Ch.14

- Instruction dependences
- Register renaming
- Pentium / PowerPC

Superscalar processors

- Goal
 - Concurrent execution of scalar instructions
 - Several independent pipelines
 - Not just more stages in one pipeline
 - Own functional units in each pipeline

| Reference | Speedup |
|-----------|---------|
| [TJAD70] | 1.8 |
| [KUCK72] | 8 |
| [WEIS84] | 1.58 |
| [ACOS86] | 2.7 |
| [SOHD90] | 1.8 |
| [SMIT89] | 2.3 |
| [DOUP89b] | 2.2 |
| [LEE91] | 7 |

Instruction-level parallelism

[Sta09 Fig 14.1, 1b14.1]

Computer Organization II, Spring 2010, Tiina Niklander 16.2.2010 2

Superscalar

Key: Fetch, Decode, Execute, Write

- Simple 4-stage pipeline: Only one instruction execute at one time
- Superpipelined: Each stage split into 2 "half-stages"
- Superscalar: Two instructions executed at the same time.

[Sta09 Fig 14.2]

Computer Organization II, Spring 2010, Tiina Niklander 16.2.2010 3

Superscalar processor

- Efficient memory usage
 - Fetch several instructions at once, prefetching (*ennaltanouto*)
 - Data fetch and store (read and write)
 - Concurrency
- Several instructions of the same process executed concurrently on different pipelines
 - Select executable instruction from the prefetched one following a policy (in-order issue/out-of-order issue)
- Finish more than one instruction during each cycle
 - Instructions may complete in different order than started (out-of-order completion)
- When can an instruction finish before the preceding ones?

Computer Organization II, Spring 2010, Tiina Niklander 16.2.2010 4

Dependencies (riippuvuus)

`add r1,r2
move r3,r1`

- True Data/Flow Dependency (*datariippuvuus*)
 - Read after Write (RAW)
 - The latter instruction needs data from former instruction
- Procedural/Control Dependency (*kontrolliriippuvuus*)
 - Instruction after the jump executed only, when jump does not happen
- Resource Conflict (*Resurssiiriippuvuus*)
 - One or more pipeline stage needs the same resource
 - Memory buffer, ALU, access to register file, ...

`JNZ R2, 100
ADD R1, =1`

Computer Organization II, Spring 2010, Tiina Niklander 16.2.2010 5

Effect of dependencies

- No Dependency
- Data Dependency (i1 uses data computed by i0)
- Procedural Dependency
- Resource Conflict (i0 and i1 use the same functional unit)

[Sta09 Fig 14.3]

Computer Organization II, Spring 2010, Tiina Niklander 16.2.2010 6

Dependencies specific to out-of-order completion

- Output Dependency (Kirjoitusriippuvuus)
 - write-after-write (WAW)
 - Two instructions alter the same register or memory location, the latter in the original code must stay
- Antidependency (Antiriippuvuus)
 - Write-after-read (WAR)
 - The former read instruction must be able to fetch the register content, before the latter write stores new value there
- Alias?
 - Two registers use indirect references to the same memory location?
 - Different virtual address, same physical address?
 - What is visible on instruction level (before MMU)?

```
load r1, X
add r2, r1, r3
add r1, r4, r5
```

```
move r2, r1
add r1, r4, r5
```

```
store R5, 40(R1)
load R6, 0(R2)
```

Computer Organization II, Spring 2010, Tiina Niklander 16.2.2010 7

Dependencies

NOTE: Newest 8th ed. only Ch14 - Incorrect definitions! SEE ERRATA Ch 12 - Correct definitions!

- i: "write" R1
.....
j: "read" R1
- i: "read" R1
.....
j: "write" R1
- i: "write" R1
.....
j: "write" R1

data dependency (RAW)

In data dependency instruction j cannot be executed before instr. i!

antidependency (WAR)

Anti- and output dependency allow change in execution order for instructions i and j, but afterwards must be checked that the right value and result remains

output dependency (WAW)

Computer Organization II, Spring 2010, Tiina Niklander 16.2.2010 8

How to handle dependencies?

- Starting point
 - All dependences must be handled one way or other
- Simple solution (as before)
 - Special hardware detects dependency and force the pipeline to wait (bubble)
- Alternative solution
 - Compiler generates instructions in such a way that there will be NO dependencies
 - No special hardware
 - simpler CPU that need not detect dependencies
 - Compiler must have very detailed and specific information about the target processor's functionality

Computer Organization II, Spring 2010, Tiina Niklander 16.2.2010 9

Parallelism (rinnakkaisuus)

- Instruction-level parallelism (käskyntason rinnakkaisuus)
 - Independent instructions of a sequence can be executed in parallel by overlapping
 - Theoretical upper limit for parallel execution of instructions
 - Depends on the code itself
- Machine parallelism (konetason rinnakkaisuus)
 - Ability of the processor to execute instructions parallel
 - How many instructions can be fetched and executed at the same time?
 - How many pipelines can be used
 - Always smaller than instruction-level parallelism
 - Cannot exceed what instructions allow, but can limit the true parallelism
 - Dependencies, bad optimization?

```
load r1 ← r2
add r3 ← r3+1
add r4 ← r4, r2
```

```
add r3 ← r3+1
add r4 ← r3, r2
load r0 ← r4
```

Computer Organization II, Spring 2010, Tiina Niklander 16.2.2010 10

Superscalar execution

issue - laukaisu, liikkeellelaskeminen
dispatch - vuorottaminen, lähettää suorittamaan

Computer Organization II, Spring 2010, Tiina Niklander 16.2.2010 11

Superscalar execution

- Instruction fetch (käskyjen nouto)
 - Branch prediction (hyppyjen ennustus)
 - prefetch (ennaltanouto) from memory to CPU
 - Instruction window (valintaikkuna)
 - set of fetched instructions
- Instruction dispatch/issue (käskyn päästäminen hihnalle)
 - Check (and remove) data, control and resource dependencies
 - Reorder; dispatch the suitable instructions to pipelines
 - Pipelines proceed without waits
 - If no suitable instruction, wait here
- Instruction complete, retire (suoritus valmistuu)
 - Commit or abort (hyväksy tai hylkää)
 - Check and remove write and antidependencies
 - wait / reorder (järjestä uudelleen)

Computer Organization II, Spring 2010, Tiina Niklander 16.2.2010 12

In-order issue, in-order complete

- Traditional sequential execution order
- No need for instruction window
- Instructions dispatched to pipelines in original order
 - Compiler handles most of the dependencies
 - Still need to check dependencies, if needed add bubbles
 - Can allow overlapping on multiple pipelines
- Instructions complete and commit in original order
 - Cannot pass, overtake (*ohittaa*) on other pipeline
 - Several instructions can complete at same time
 - Commit/Abort

Computer Organization II, Spring 2010, Tiina Niklander 16.2.2010 13

In-order issue, in-order complete

Fetch 2 instructions at the same time
 I1 needs 2 cycles for execution
 I3 and I4: resource dependency
 I5 (consume) and I4 (produce): data dependency
 I5 and I6: resource dependency

| Decode | Execute | Write | Cycle |
|--------|---------|-------|-------|
| I1 I2 | | | 1 |
| I3 I4 | I1 I2 | | 2 |
| I5 I6 | | | 3 |
| | I3 I4 | I1 I2 | 4 |
| | | I3 I4 | 5 |
| | I5 I6 | | 6 |
| | | I5 I6 | 7 |
| | | | 8 |

Decode clean before fetching next two instructions
 Instructions queue for execution in decode unit
 Writes delayed to maintain in-order completion

Computer Organization II, Spring 2010, Tiina Niklander 16.2.2010 14

In-order issue, out-of-order complete

- Like previous, but
 - Allow commit in different order than issued order (allow passing)
 - Clear write and antidep. before writing the results

Fetch 2 instructions at the same time
 I1 needs 2 cycles for execution
 I3 and I4: resource dependency
 I5 (consume) and I4 (produce): data dep.
 I5 and I6: resource dependency

(WAW)

| Decode | Execute | Write | Cycle |
|--------|---------|-------|-------|
| I1 I2 | | | 1 |
| I3 I4 | I1 I2 | | 2 |
| I5 I6 | | | 3 |
| | I3 I4 | I1 I2 | 4 |
| | | I3 I4 | 5 |
| | I5 I6 | | 6 |
| | | I5 I6 | 7 |

Output dependency

| |
|------------------|
| 1: R3 ← R3 op R5 |
| 2: R4 ← R3 + 1 |
| 3: R3 ← R5 + 1 |
| 4: R7 ← R3 op R4 |

Computer Organization II, Spring 2010, Tiina Niklander 16.2.2010 15

Out-of-order issue, out-of-order complete

- Dispatch instructions for execution in any suitable order
 - Need instruction window
 - Processor looks ahead (at the future instructions)
 - Must consider the dependencies during dispatch
- Allow instructions to complete and commit in any suitable order
 - Check and clear write and antidependencies

Anti dependency (WAR)

| |
|------------------|
| 1: R3 ← R3 op R5 |
| 2: R4 ← R3 + 1 |
| 3: R3 ← R5 + 1 |
| 4: R7 ← R3 op R4 |

I3 must not write to R3, before I1 has read the content

True superscalar processor

Computer Organization II, Spring 2010, Tiina Niklander 16.2.2010 16

Out-of-order issue, Out-of-order complete

Fetch 2 instructions at the same time
 I1 needs 2 cycles for execution
 I3 and I4: resource dependency
 I5 (consume) and I4 (produce): data dependency
 I5 and I6: resource dependency

| Decode | Window | Execute | Write | Cycle |
|--------|--------|---------|-------|-------|
| I1 I2 | | | | 1 |
| I3 I4 | I1 I2 | | | 2 |
| I5 I6 | I3 I4 | | | 3 |
| | I5 | I1 I2 | I1 I2 | 4 |
| | | I3 I4 | I3 I4 | 5 |
| | | I5 I6 | I5 I6 | 6 |

Instruction window, (just a buffer, not a pipeline stage)

Computer Organization II, Spring 2010, Tiina Niklander 16.2.2010 17

Register renaming (rekistereiden uudelleennimeäminen)

- One cause for some of the dependencies is the usage of names
 - The same name could be used for several independent elements
 - Thus, instructions have unneeded write and antidependencies
 - Causing unnecessary waits
- Solution: Register renaming
 - Hardware must have more registers (than visible to the programmer and compiler)
 - Hardware allocates new real registers during execution in order to avoid name-based dependencies (nimiriippuvuus)
- Need
 - More internal registers (register files, register set), e.g. Pentium II has 40 working registers
 - Hardware that is capable of allocating and managing registers and performing the needed mapping

| |
|--------------|
| R3 ← R3 + R5 |
| R4 ← R3 + 1 |
| R3 ← R5 + 1 |
| R7 ← R3 + R4 |

Computer Organization II, Spring 2010, Tiina Niklander 16.2.2010 18

Register renaming

Output dependency (WAW):
(*Kirjoitusriippuvuus*)
i3 must not finish before i1

Anti dependency (RAW):
(*antiriippuvuus*)
i3 must not finish before i2 has read the value from R3

Rename R3 use work registers R3a, R3b, R3c
Other registers similarly: R4b, R5a, R7b

No more dependencies based on names!

R3 ← R3 + R5 (i1)
R4 ← R3 + 1 (i2)
R3 ← R5 + 1 (i3)
R7 ← R3 + R4 (i4)

R3b ← R3a + R5a (i1)
R4b ← R3b + 1 (i2)
R3c ← R5a + 1 (i3)
R7b ← R3c + R4b (i4)

Why R3a and R3b?

Computer Organization II, Spring 2010, Tina Niklander 16.2.2010 19

Impact of additional hardware

base: out-of-order issue
+ld/st: base and duplicate load/store unit for data cache
+alu: base and duplicate ALU

Window size (construction) 8 16 32

Computer Organization II, Spring 2010, Tina Niklander 16.2.2010 20

Superscalar – conclusion

- Several functionally independent units
- Efficient use of memory hierarchy
 - Allows parallel memory fetch and store
- Instruction prefetch (*käskeyjen ennaltanouto*)
 - Branch prediction (*hyppyjen ennustaminen*)
- Hardware-level logic for dependency detections
 - Circuits to pass information for other functional unit at the same time as storing to register or memory
- Hardware-level logic to dispatch several independent instructions
 - Dependencies → dispatching order
- Hardware-level logic to maintain correct completion order (*valmistumisjärjestys*)
 - Dependencies → commit-order

Computer Organization II, Spring 2010, Tina Niklander 16.2.2010 21

Computer Organization II

Pentium 4

Computer Organization II, Spring 2010, Tina Niklander 16.2.2010 22

Pentium 4

AGU = address generation unit
BTB = branch target buffer
D-TLB = data translation lookaside buffer
I-TLB = instruction translation lookaside buffer

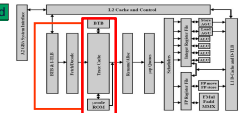
Computer Organization II, Spring 2010, Tina Niklander 16.2.2010 23

Pipeline

- Outside CISC (IA-32)
- Execution in micro-operations as RISC
 - Fetch CISC instruction and translate it to one or more micro-operations (μops) to L1-level cache (*trace cache*)
 - Rest of the superscalar pipeline operates with these fixed-length micro-operations (118b)
- Long pipeline
 - Extra stages (5 and 20) for propagation delays

TC Next IP = trace cache next instruction pointer
TC Fetch = trace cache fetch
Alloc = allocate
Rename = register renaming
Que = micro-op queuing
Sch = micro-op scheduling
Disp = Dispatch
RF = register file
Ex = execute
Flgs = flags
Br Ck = branch check

Computer Organization II, Spring 2010, Tina Niklander 16.2.2010 24



Generation of μ ops

a) Fetch IA-32 instruction from L2 cache and generate μ ops to L1

- Uses Instruction Lookaside Buffer (I-TLB) and Branch Target Buffer (BTB)
- four-way set-associative cache, 512 lines
- 1-4 μ ops (\approx 118 bit RISC) per instruction (most cases), if more then stored to microcode ROM

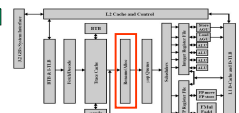
b) Trace Cache Next Instruction Pointer - instruction selection

- Dynamic branch prediction based on history (4-bit)
- If no history available, Static branch prediction
 - backward, predict "taken"
 - forward, predict "not taken"

c) Fetch instruction from L1-level trace cache

d) Drive – wait (instruction from trace cache to rename/allocator)

Computer Organization II, Spring 2010, Tiina Niklander 16.2.2010 25



Resource allocation

e) Allocate resources

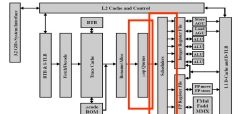
- 3 micro-operations per cycle
- Allocate an entry from Reorder Buffer (ROB) for the μ ops (126 entries available)
- Allocate one of the 128 internal work registers for the result
- And, possibly, one load (of 48) OR store (of 24) buffer

f) Register renaming

- Clear 'name dependencies' by remapping registers (16 architectural regs to 128 physical registers)
- If no free resource, wait (\rightarrow out-of-order)

- ROB-entry contains bookkeeping of the instruction progress
 - Micro-operation and the address of the original IA-32 instr.
 - State: scheduled, dispatched, completed, ready
 - Register Alias Table (RAT): which IA-32 register \rightarrow which physical register

Computer Organization II, Spring 2010, Tiina Niklander 16.2.2010 26



Window of Execution

g) Micro-Op Queueing

- 2 FIFO queues for μ ops
 - One for memory operations (load, store)
 - One for everything else
- No dependencies, proceed when room in scheduling

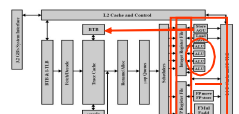
h) Micro-Op Scheduling

- Retrieve μ ops from queue and dispatch for execution
- Only when operands ready (check from ROB-entry)

i) Dispatching

- Check the first instructions of FIFO-queues (their ROB-entries)
- If execution unit needed is free, dispatch to that unit
- Two queues \rightarrow out-of-order issue
- max 6 micro-ops dispatched in one cycle
 - ALU and FPU can handle 2 per cycle
 - Load and store each can handle 1 per cycle

Computer Organization II, Spring 2010, Tiina Niklander 16.2.2010 27



Integer and FP Units

j) Get data from register or L1 cache

k) Execute instruction, set flags (lipuke)


- Several pipelined execution units
 - 4 * ALU, 2 * FPU, 2 * load/store
 - E.g. fast ALU for simple ops, own ALU for multiplications
- Result storing: in-order complete
- Update ROB, allow next instruction to the unit

l) Branch check

- What happen in the jump /branch instruction
- Was the prediction correct?
- Abort incorrect instruction from the pipeline (no result storing)

m) Drive – update BTB with the branch result


Computer Organization II, Spring 2010, Tiina Niklander 16.2.2010 28



Pentium 4 Hyperthreading

- One physical IA-32 CPU, but 2 logical CPUs
- OS sees as 2 CPU SMP (symmetric multiprocessing)
 - Processors execute different processes or threads
 - No code-level issues
 - OS must be capable to handle more processors (like scheduling, locks)
- Uses CPU wait cycles
 - Cache miss, dependences, wrong branch prediction
- If one logical CPU uses FP unit the other one can use INT unit
 - Benefits depend on the applications

Computer Organization II, Spring 2010, Tiina Niklander 16.2.2010 29



Pentium 4 Hyperthreading

- Duplicated (kahdennettu)
 - IP, EFLAGS and other control registers
 - Instruction TLB
 - Register renaming logic
- Split (puolitettu)
 - No monopoly, non-even split allowed
 - Reordering buffers (ROB)
 - Micro-op queues
 - Load/store buffers
- Shared (jaettu)
 - Register files (128 GPRs, 128 FPRs)
 - Caches: trace cache, L1, L2, L3
 - Registers needed during μ ops execution
 - Functional units: 2 ALU, 2 FPU, 2 ld/st-units

Intel Nehalem arch.:
8 cores on one chip,
1-16 threads (820 million transistors)
First lauched processor
Core i7 (Nov 2008)

Computer Organization II, Spring 2010, Tiina Niklander 16.2.2010 30

Computer Organization II

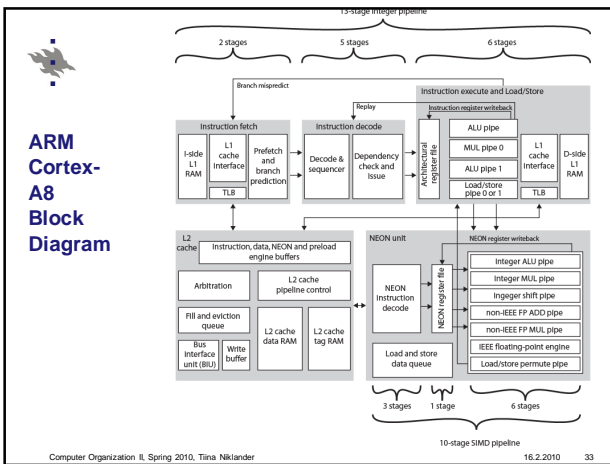
ARM Cortex-A8

Computer Organization II, Spring 2010, Tina Niklander 16.2.2010 31

ARM CORTEX-A8

- ARM refers to Cortex-A8 as application processors
- Embedded processor running complex operating system
 - Wireless, consumer and imaging applications
 - Mobile phones, set-top boxes, gaming consoles, automotive navigation/entertainment systems
- Three functional units
- Dual, in-order-issue, 13-stage pipeline
 - Keep power required to a minimum
 - Out-of-order issue needs extra logic consuming extra power
- Separate SIMD (single-instruction-multiple-data) unit called NEON
 - 10-stage pipeline

Computer Organization II, Spring 2010, Tina Niklander 16.2.2010 32



Instruction Fetch Unit

- Predicts instruction stream
- Fetches instructions from the (included) L1 instruction cache
 - Into buffer for decode pipeline
 - Up to four instructions per cycle
- Speculative instruction fetches
- Branch or exceptional instruction cause pipeline flush
- Two-level global history branch predictor
 - Branch Target Buffer (BTB) and Global History Buffer (GHB)
- Return stack to predict subroutine return addresses
- Can fetch and queue up to 12 instructions

Computer Organization II, Spring 2010, Tina Niklander 16.2.2010 34

Instruction Fetch Unit: Processing Stages

- F0 address generation unit (AGU)
 - Next address sequentially
 - Or branch target address (from branch prediction of previous address)
- F1 fetch instructions from L1
 - In parallel, check the branch prediction for the next address
- F2 Place instruction to instr. queue
 - If branch prediction, new target address sent to AGU
- Issues instructions to decode two at a time

Computer Organization II, Spring 2010, Tina Niklander 16.2.2010 35

Instruction Decode Unit

- Dual pipeline structure, *pipe0* and *pipe1*
 - Two instructions at a time
 - Pipe0 contains older instruction in program order
 - If instruction in pipe0 cannot issue, pipe1 will not issue
- In-order instruction issue and retire
 - Results written back to register file at end of execution pipeline
 - no WAR hazards
 - tracks WAW hazards and straightforward recovery from flush
 - Decode pipeline to prevent RAW hazards

Computer Organization II, Spring 2010, Tina Niklander 16.2.2010 36

Instruction Decode Unit: Processing Stages

- D0 Decompress Thumbs and do preliminary decode
- D1 Instruction decode completed
- D2 Write instruction to and read instructions from pending/replay queue
- D3 instruction scheduling logic
 - Scoreboard predicts register availability using static scheduling
 - Hazard checking
- D4 Final decode - control signals for integer execute load/store units

Computer Organization II, Spring 2010, Tina Niklander 16.2.2010 37

Integer Execution Unit

- Two symmetric (ALU) pipelines, an address generator for load and store instructions, and multiply pipeline
- Multiply unit instructions routed to pipe0
 - Performed in stages E1 through E3
 - Multiply accumulate operation in E4
- E0 Access register file
 - Up to six registers for two instructions
- E1 Barrel shifter if needed.
- E2 ALU function
- E3 If needed, completes saturation arithmetic
- E4 Change in control flow prioritized and processed
- E5 Results written back to register file

Computer Organization II, Spring 2010, Tina Niklander 16.2.2010 38

Integer Execution Unit

Computer Organization II, Spring 2010, Tina Niklander 16.2.2010 39

Load/store pipeline

- Parallel to integer pipeline
- E1 Memory address generated from base and index register
- E2 address applied to cache arrays
- E3 load, data returned and formatted
- E3 store, data are formatted and ready to be written to cache
- E4 Updates L2 cache, if required
- E5 Results are written to register file

Computer Organization II, Spring 2010, Tina Niklander 16.2.2010 40

SIMD and Floating-Point Pipeline

- SIMD and floating-point instructions pass through integer pipeline
- Processed in separate 10-stage pipeline
 - NEON unit
 - Handles packed SIMD instructions
 - Provides two types of floating-point support
- If implemented, vector floating-point (VFP) coprocessor performs IEEE 754 floating-point operations
 - If not, separate multiply and add pipelines implement floating-point operations

Computer Organization II, Spring 2010, Tina Niklander 16.2.2010 41

ARM Cortex-A8 NEON & Floating Point Pipeline

Computer Organization II, Spring 2010, Tina Niklander 16.2.2010 42



Review Questions / Kertauskysymyksiä

- Differences / similarities of superscalar and trad. pipeline?
- What new problems must be solved?
- How to solve those?
- What is register renaming and why it is used?

- Miten superskalaaritoteutus eroaa tavallisesta liukuhinnoitetusta toteutuksesta?
- Mitä uusia rakenteesta johtuvia ongelmia tulee ratkottavaksi?
- Miten niitä ongelmia ratkotaan?
- Mitä tarkoittaa rekistereiden uudelleennimeäminen ja mitä hyötyä siitä on?