

TDTS08: whattolearn

Name	PIN	Email
Pål Kastman	851212-7575	palka285@student.liu.se

Contents

1	Introduction	4
2	Lecture 1: Memory Systems	4
2.1	Cache memory	4
2.1.1	Locality of reference	4
2.1.2	Replacement methods	4
2.1.3	Cache Design	4
2.1.4	Split Caches	4
2.1.5	Unified Caches	4
2.1.6	Direct mapping cache	4
2.1.7	Associative mapping cache	4
2.1.8	Fully Associative Organization	4
2.1.9	Write Policy	4
2.2	Virtual memory	5
2.2.1	Paging	5
2.2.2	Page fault	5
2.2.3	Page replacement	5
2.2.4	Replacement algorithms	5
3	Lecture 3: Instruction pipelining	5
3.1	Pipeline hazards	5
3.1.1	Structural (resource) hazards	5
3.1.2	Data hazards	6
3.1.3	Control hazards	6
3.2	Branch handling	6
3.2.1	Pre-fetch branch target	6
3.2.2	Loop buffer	6
3.2.3	Delayed branch	6
3.3	Branch prediction	6
3.3.1	Static Branch Prediction	6
3.3.2	Dynamic Branch Prediction	6
3.3.3	Bimodal Prediction	6
4	Lecture 4: RISC Computers	6
4.1	Program execution features	6
4.2	RISC characteristics	6
4.2.1	Register Windows	7
4.2.2	Main advantages of RISC	7
4.2.3	Criticism of RISC	7
4.3	RISC vs. CISC	7
5	Lecture 5: Superscalar Architecture (SSA)	7
5.1	Instruction-level parallelism	8
5.2	Superpipelining	8
5.3	Difference between Superpipelined and Superscalar Designs	8
5.4	Superscalar Superpipeline Design	9
5.5	Dependency issues	9
5.5.1	Resource Conflicts	10

5.5.2	Procedural Dependency	11
5.5.3	Data Conflicts	11
5.5.4	Window of execution	11
5.5.5	Data Dependency	11
5.5.6	True Data Dependency	11
5.5.7	Output Dependency	11
5.5.8	Anti Dependency	11
5.5.9	Effect of Dependencies	11
5.6	Parallel instruction execution	11
5.6.1	Instruction vs Machine Parallelism	11
5.6.2	Division and Decoupling	11
5.6.3	SSA instruction Execution Policies	11
5.6.4	In-Order with In-Order Completion	11
5.6.5	In-Order issue with Out-of-Order Completion	11
6	Lecture 6: VLIW Processors	11
6.1	Very Long Instruction Word Processors	12
6.1.1	Explicit Parallelism	12
6.1.2	Main issues	12
6.1.3	Software issues	13
6.2	Loop unrolling	13
6.3	IA-64 architecture	13
6.3.1	Predicated execution	13
6.3.2	Instruction format	13
6.3.3	Branch Predication	13
6.3.4	Placement of Loading	13
6.3.5	Speculative Loading	13
7	Parallel Processing	13
7.0.1	Why Parallel Processing?	13
7.0.2	Parallel Computer	13
7.0.3	Parallel Program	13
7.1	Architecture classification	13
7.1.1	Flynn's Classification of Architectures	13
7.2	Performance evaluation	14
7.3	Interconnection network	14
7.3.1	Single Bus	14
7.3.2	Completely Connected Network	14
7.3.3	Crossbar Network	14
7.3.4	Mesh Network	14
7.3.5	Hypercube Network	14
8	Lecture 8: SIMD Architectures	14

1 Introduction

This documents only purpose is to document what I will need to learn in the course tdts08

2 Lecture 1: Memory Systems

Memory hierarchy

registers, cache, main memory etc.

2.1 Cache memory

how does it work? why do we use it? positive and negative things about it!

2.1.1 Locality of reference

Temporal- and spatial locality, differences and how we use these in caches.

2.1.2 Replacement methods

yes, there's different kinds of them.

2.1.3 Cache Design

One, two, three level caches, how does it work. up- and downsides of them.

2.1.4 Split Caches

2.1.5 Unified Caches

2.1.6 Direct mapping cache

2.1.7 Associative mapping cache

2.1.8 Fully Associative Organization

2.1.9 Write Policy

Write through

Write through with buffered write

Write back

2.2 Virtual memory

2.2.1 Paging

2.2.2 Page fault

2.2.3 Page replacement

2.2.4 Replacement algorithms

3 Lecture 3: Instruction pipelining

how does it work, and why do we need it?

3.1 Pipeline hazards

what problems do we have with pipelining

3.1.1 Structural (resource) hazards

what are they and how do we solve it?

3.1.2 Data hazards**3.1.3 Control hazards****3.2 Branch handling****3.2.1 Pre-fetch branch target****3.2.2 Loop buffer****3.2.3 Delayed branch****3.3 Branch prediction****3.3.1 Static Branch Prediction****3.3.2 Dynamic Branch Prediction****3.3.3 Bimodal Prediction****4 Lecture 4: RISC Computers**

What are they, how do they work, and why do we need them?

Semantic gap

Problems with RISC

4.1 Program execution features

what are programs doing most of the time?

4.2 RISC characteristics

Small number of simple instructions

Execution of one instruction per clock cycle

Complex operations are executed as a sequence of simple instructions

only LOAD and STORE instructions reference data in memory

Only a few simple addressing modes are used

Instructions are of fixed length and uniform format.

Large number of registers, this is because the reduced complexity of the processor leaves silicon space on the chip to implement them (opposite of CISC).

4.2.1 Register Windows

Large number of registers is usually very useful.

However, if contents of all registers must be saved at every procedure call, more registers mean longer delay.

A solution to this problem is to divide the register file into a set of fixed-size windows.

- Each window is assigned to a procedure.
- Windows for adjacent procedures are overlapped to allow parameter passing

4.2.2 Main advantages of RISC

Best support is given by optimizing most used and most time consuming architecture aspects.

- Frequently executed instructions.
- Simple memory reference.
- Procedure call/return.
- Pipeline design.

Consequently, we have:

- High performance for many applications;
- Less design complexity;
- Reduced power consumption:
- reducing design cost and time-to-market (newer technology)

4.2.3 Criticism of RISC

Operation might take several instructions to accomplish

more memory access might be needed

Execution speed may be reduced for certain applications.

It usually leads to longer programs, which needs larger memory space to store.

It makes it more difficult to program machine codes and assembly programs.

4.3 RISC vs. CISC

In embedded processors RISC is the better choice

5 Lecture 5: Superscalar Architecture (SSA)

Computer designed to improve computation on scalars instructions.

A scalar is a variable that can hold only one atomic value at a time, e.g., an integer or a real.

A scalar architecture processes one data item at a time – the computers we discussed up till now.

Examples of non-scalar variables:

- Arrays – Vector Processor
- Matrices – Graphics Processing Unit (GPU)
- Records

In a superscalar architecture (SSA), several scalar instructions can be initiated simultaneously and executed independently.

5.1 Instruction-level parallelism

Most operations are on scalar quantities, speed up these operations will lead to large performance improvement.

5.2 Superpipelining

Divide pipelining stages into different several sub-stages, and hence increase the number of instructions which are handled by the pipeline at the same time.

A picture would be nice here

- For example by diving each stage into two sub stages, we will be able (in the ideal situation) to perform each stage at twice the stage.
- No duplication of hardware is needed.
- Not all stages can be divided into (equal length) sub stages.
- Hazards more difficult to resolve.
- More complex hardware.
- Interrupt handling and testing will be more complicated.

5.3 Difference between Superpipelined and Superscalar Designs

The difference is that Superpipeline can perform several instructions in one clock cycle, whereas Superscalar performs several instructions in parallel

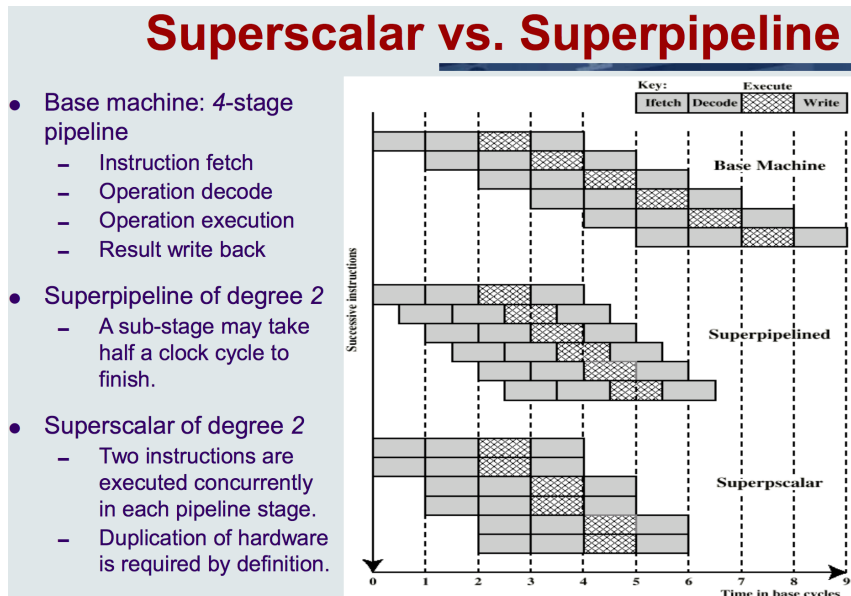


Figure 1 – Difference between a Superscalar and a Superpipeline architecture

5.4 Superscalar Superpipeline Design

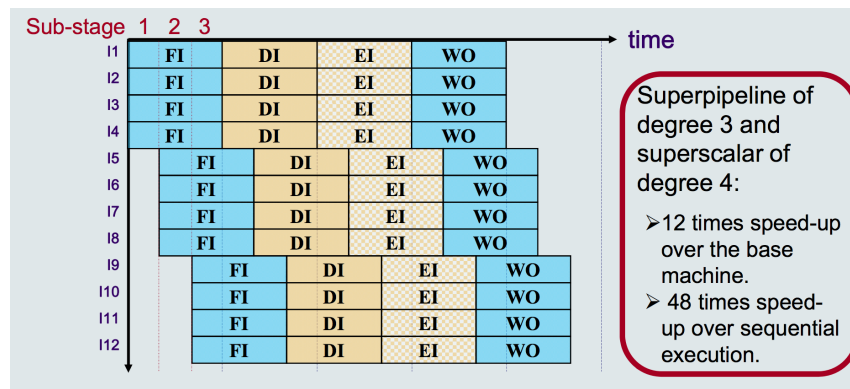


Figure 2 – Superscalar Superpipeline architecture

The new trend is to combined the two ideas. In Figure 2 it says that this would give us a 48 times speedup, this is not the case though because of data dependencies.

5.5 Dependency issues

The main problems are: – Resource conflicts.
– Control (procedural) dependency.

- Data dependencies.

These are very similar to the cases in normal pipelining (data hazards). But the consequences are more severe here because the parallelism are greater and thus larger amount of performance will be lost.

5.5.1 Resource Conflicts

Several instructions compete for the same hardware resource at the same time.

- For instance, two arithmetic instructions need the same floating-point unit for execution.
- similar to structural hazards in pipeline.

They can be solved partly by introducing several hardware units for the same functions.

- e.g. have two floating point units.
- the hardware units can also be pipelined to support several operations at the same time.
- however, memory units **can't be duplicated**.

5.5.2 Procedural Dependency**5.5.3 Data Conflicts****5.5.4 Window of execution****5.5.5 Data Dependency****5.5.6 True Data Dependency****5.5.7 Output Dependency****5.5.8 Anti Dependency****5.5.9 Effect of Dependencies****5.6 Parallel instruction execution****5.6.1 Instruction vs Machine Parallelism****5.6.2 Division and Decoupling****5.6.3 SSA instruction Execution Policies****5.6.4 In-Order with In-Order Completion**

This seems important to understand

5.6.5 In-Order issue with Out-of-Order Completion

This seems important to understand

6 Lecture 6: VLIW Processors

The difference from Superscalar is that in that architecture the hardware solves everything for us, if improvements are made on the design we don't need to change the programs.

The down side though is that the architecture is very complex.

- A lot of hardware is needed for run-time detection of parallelism.
- It consumes a lot of power.
- There is, therefore a limit in how far we can go with this technique.

The instruction window for execution is limited in size.

- This limits the capacity to detect large number of parallel instructions.

6.1 Very Long Instruction Word Processors

Several operations that can be executed in parallel are placed in a single instruction word.

- VLIW rely on compile-time detection of parallelism.
- The compiler analyzes the program and detects operations to be executed in parallel.

After one instruction has been fetched, all the corresponding operations are issued in parallel. The instruction window limitation disappears: the compiler can potentially analyze the whole program to detect parallel operations.

6.1.1 Explicit Parallelism

Instruction parallelism scheduled at compile time.

- Included within the machine instructions explicitly.

An EPIC (Explicitly Parallel Instruction Computing) processor

- uses this information to perform parallel execution.

The hardware is very much simplified

- The controller is similar to a simple scalar computer.
- The number of FUs can be increased without needing additional sophisticated hardware to detect parallelism, as in SSA.

Compiler has much more time to determine parallel operations.

- This analysis is only done once off-line, while run-time detection is carried out by SSA hardware for each execution of the code.
- Good compilers can detect parallelism based on global analysis of the whole program.

6.1.2 Main issues

- Need a large number of registers.
- Large data transport capacity is needed between FUs and the register files and between register files and memory.
- High bandwidth between instruction cache and fetch unit is also needed due to long instructions.

6.1.3 Software issues

6.2 Loop unrolling

6.3 IA-64 architecture

6.3.1 Predicated execution

6.3.2 Instruction format

6.3.3 Branch Predication

NOT THE SAME AS BRANCH PREDICTION

6.3.4 Placement of Loading

6.3.5 Speculative Loading

7 Parallel Processing

7.0.1 Why Parallel Processing?

7.0.2 Parallel Computer

7.0.3 Parallel Program

7.1 Architecture classification

7.1.1 Flynn's Classification of Architectures

- Single instruction, single data stream - **SISD**
- Single instruction, multiple data stream - **SIMD**
- Multiple instruction, single data stream - **MISD**
- Multiple instruction, multiple data stream- **MIMD**

GÅ IGENOM DESSA

7.2 Performance evaluation

7.3 Interconnection network

7.3.1 Single Bus

7.3.2 Completely Connected Network

7.3.3 Crossbar Network

7.3.4 Mesh Network

7.3.5 Hypercube Network

8 Lecture 8: SIMD Architectures