#### Instruction ordering

# Instruction scheduling

Advanced Compiler Construction Michel Schinz – 2021-04-22 When emitting the instructions of a program, a compiler imposes a total order on them, but:

- there is (usually) more than one valid order,

- some orders might be better than others.

Example:

Two independent instructions appearing in sequence can be swapped.

### Instruction scheduling

#### Goal of instruction scheduling:

Find, among all valid permutations of the instructions of a program, one that is better than the others. (Usually, better = executes faster.)

#### **Pipeline stalls**

Modern, pipelined architectures can usually issue at least one instruction per clock cycle.

However, an instruction can execute only if its arguments are ready, otherwise the pipeline **stalls** until it is the case.

Causes for stalls:

- the instruction producing the argument has not finished executing yet,
- the argument must be fetched from memory,
- etc.

# Scheduling example

The following example will illustrate how proper scheduling can reduce the time required to execute a piece of RTL code. We assume the following delays for instructions:

Instruction kind	RTL notation	Delay
Memory load or store	R <sub>a</sub> ← Mem[R <sub>b</sub> +c] Mem[R <sub>b</sub> +c] ← R <sub>a</sub>	3
Multiplication	$R_a \leftarrow R_b \star R_c$	2
Addition	$R_a \leftarrow R_b + R_c$	1

## Scheduling example



### Instruction dependences

An instruction  $i_2$  **depends** on an instruction  $i_1$  when it is not possible to execute  $i_2$  before  $i_1$  without changing the behavior of the program. We distinguish three kinds of dependencies:

- 1. true dependency  $i_2$  reads a value written by  $i_1$  (**read after write** or **RAW**),
- 2. anti-dependency i2 writes a value read by i1 (write after read or WAR),
- 3. anti-dependency  $i_2$  writes a value written by  $i_1$  (write after write or WAW).

#### Anti-dependencies

Anti-dependencies do not arise from the flow of data. They are due to a single location being reused.

hey are due to a single location being reused.

Often, they can be removed by "renaming" locations, e.g. using different registers.

In the example below, the program (left) contains a WAW anti-dependency that can be removed by "renaming" the second use of  $R_1$ .



#### **Computing dependencies**

# Dependence graph

Identifying dependencies is:

- easy if they only access registers,

- impossible (in general) if they access memory.

For memory, conservative approximations have to be used. We won't cover them here.

# The **dependence graph** is a directed graph representing dependencies among instructions:

- the nodes are the instructions to schedule,

- there is an edge from  $n_1$  to  $n_2$  iff  $n_2$  depends on  $n_1. \label{eq:n_1}$ 

Any topological sort of the nodes of this graph is a valid schedule of the instructions.

# Dependence graph example





#### ···▶ antidependence

# List scheduling

Optimal instruction scheduling is NP-complete. **List scheduling** is:

- a heuristic scheduling technique,
- that works on a single basic block.

Basic idea:

- simulate the execution of the instructions, and
- schedule them only when their operands are ready.

# List scheduling algorithm

The list scheduling algorithm maintains two lists:

- **ready**: the instructions that could be scheduled without stall, ordered by priority,
- active: the instructions that are being executed.
- At each step:
- the highest-priority instruction from ready is scheduled,
- it gets moved to active,
- it stays there for a time equal to its delay.

Before scheduling is performed, renaming is done to remove all antidependencies that can be removed.

# Instruction priority

Nodes (i.e. instructions) are sorted by priority in the ready list. The priority of a node can be defined as:

- the length of the longest latency-weighted path from it to a root of the dependence graph,
- the number of its immediate successors,
- the number of its descendants,
- its latency,
- etc.

Unfortunately, none of these is better for all cases.

#### List scheduling example

a <sup>13</sup>	Cycle	ready	active
b <sup>10</sup> c <sup>12</sup> priority	1 [a	<sup>13</sup> ,C <sup>12</sup> ,e <sup>10</sup> ,g <sup>8</sup> ]	[a]
	2 [ <mark>c</mark>	<sup>12</sup> ,e <sup>10</sup> ,g <sup>8</sup> ]	[a,c]
	3 [ <mark>e</mark>	<sup>10</sup> ,g <sup>8</sup> ]	[a,c,e]
	4 [ <mark>b</mark>	<sup>10</sup> ,g <sup>8</sup> ]	[b,c,e]
f <sup>7</sup> , g <sup>8</sup>	5 [ <mark>d</mark>	<sup>9</sup> ,g <sup>8</sup> ]	[d,e]
	6 [ <mark>g</mark>	8]	[d,g]
n° 🔨	7 [ <mark>f</mark>	7]	[f,g]
i <sup>3</sup>	8 []		[f,g]
	9 [ <mark>h</mark>	5]	[h]
A node's priority is the length	10 []		[h]
of the longest latency-	11 [ <mark>i</mark>	3]	[i]
weighted path from it to a root	12 []		[i]
of the dependence graph	13 []		[i]
	14 []		[]

# Scheduling conflicts

Should scheduling be done before or after register allocation?

- If it is done first, register allocation can introduce spilling code that destroys the schedule.
- If it is done second, register allocation can introduce anti-dependencies when reusing registers.

Solution:

- schedule first,
- allocate registers
- schedule once more if spilling was necessary.