

CSE211: Compiler Design

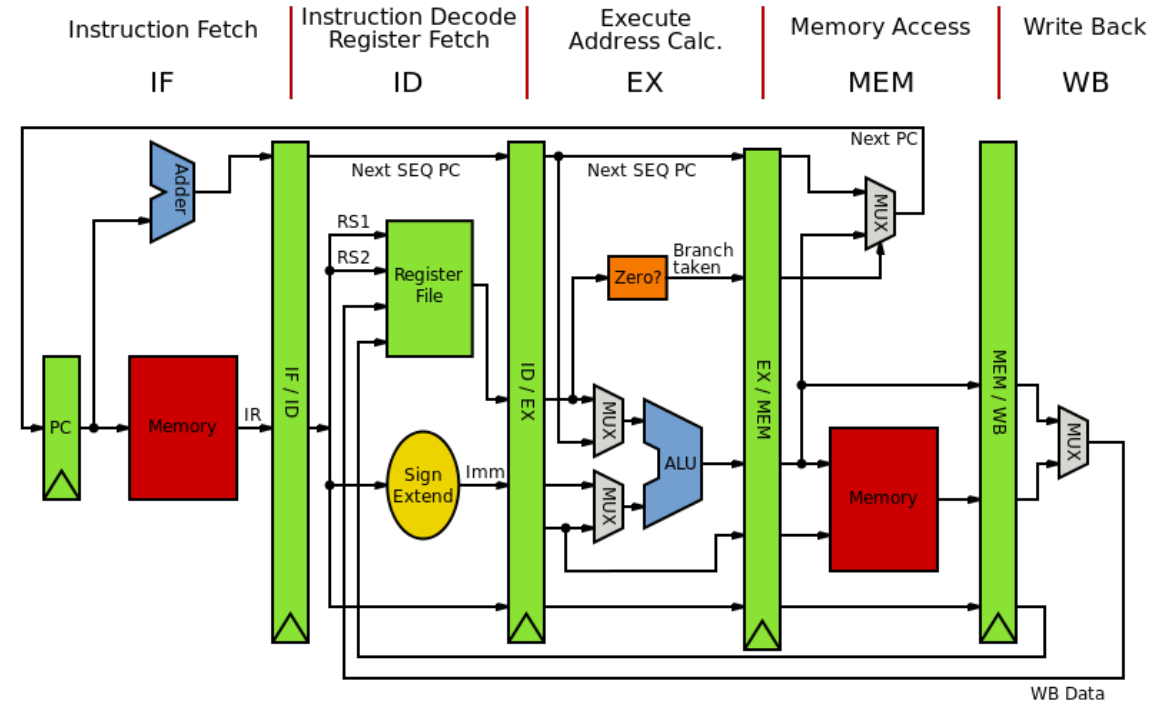
Nov. 10, 2020

- **Topic:** instruction-level parallelism (ILP)

- dependency graphs/chains
- loop unrolling
- reductions

- **Discussion questions:**

- What is instruction level parallelism?
- How can modern processors exploit ILP?



MIPS pipeline image from:

[https://commons.wikimedia.org/wiki/Pipeline_\(computer_hardware\)](https://commons.wikimedia.org/wiki/Pipeline_(computer_hardware))

Announcements

- Homework 2 is due Nov. 12
 - One more office hour session on Wednesday
 - I will do the open door sessions
 - API for pycfg can be a little confusing. Please discuss on Canvas
- Midterm will be released Nov. 12 and due on Nov. 19
 - Email me clarification questions, not technical questions
 - Please don't discuss with classmates
- Paper/Project proposals submitted (latest possible, preferably earlier):
Nov. 24

Resources: The Dragon Book

- This week: sections from chapter 10
- Next week: sections from chapter 11
- Book link posted on Canvas

CSE211: Compiler Design

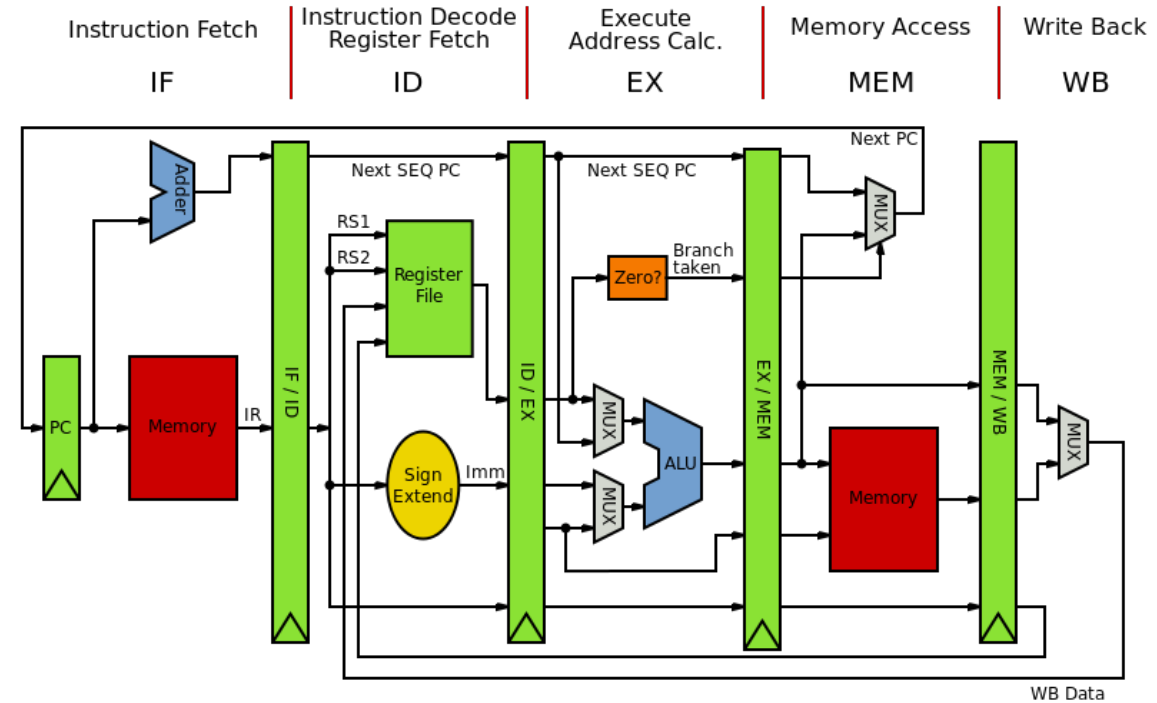
Nov. 10, 2020

- **Topic:** instruction-level parallelism (ILP)

- dependency graphs/chains
- loop unrolling
- reductions

- **Discussion questions:**

- What is instruction level parallelism?
- How can modern processors exploit ILP?



MIPS pipeline image from:

[https://commons.wikimedia.org/wiki/Pipeline_\(computer_hardware\)](https://commons.wikimedia.org/wiki/Pipeline_(computer_hardware))

Instruction-level Parallelism (ILP)

- Parallelism from a single stream of instructions.
 - Output of program must match exactly a sequential execution!
- Widely applicable:
 - most mainstream programming languages are sequential
 - most deployed hardware has components to execute ILP
- Can benefit from a combination of hardware and software scheduling
- While it can be done by hand, its better to implement in a compiler

Finding dependencies in the compiler

- What type of instructions can be done in parallel?

Finding dependencies in the compiler

- What type of instructions can be done in parallel?

two instructions can be executed in parallel if they are independent

Finding dependencies in the compiler

- What type of instructions can be done in parallel?

two instructions can be executed in parallel if they are independent

```
x = z + w;  
a = b + c;
```

Two instructions are independent if the operand registers are disjoint from the result registers

Finding dependencies in the compiler

- What type of instructions can be done in parallel?

two instructions can be executed in parallel if they are independent

```
x = z + w;  
a = b + c;
```

Two instructions are independent if the operand registers are disjoint from the result registers

instructions that are not independent cannot be executed in parallel

```
x = z + w;  
a = b + x;
```

Finding dependencies in the compiler

- What type of instructions can be done in parallel?

two instructions can be executed in parallel if they are independent

```
x = z + w;  
a = b + c;
```

Two instructions are independent if the operand registers are disjoint from the result registers

instructions that are not independent cannot be executed in parallel

```
x = z + w;  
a = b + x;
```

Many times, dependencies can be easily tracked in the compiler:

*Easier with:
+ within a basic block
+ using SSA form*

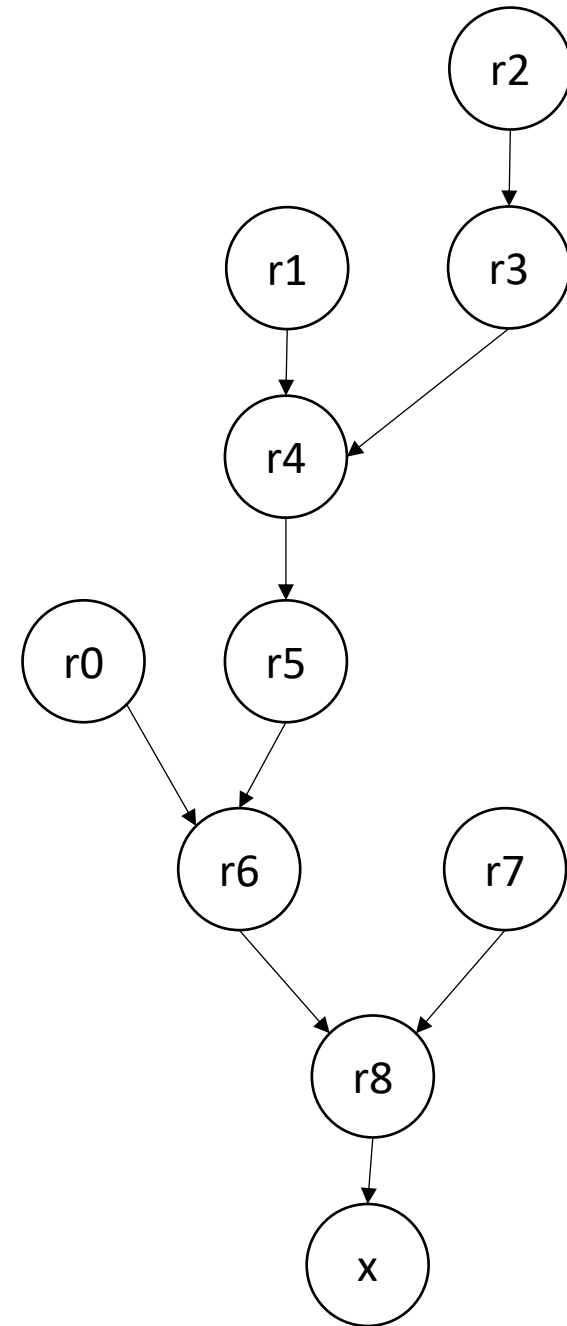
*Harder with:
- memory locations*

Different types of dependencies

- Data Dependence
- Control Dependence
- Memory Dependence

Data Dependency analysis from Oct. 15 lecture:

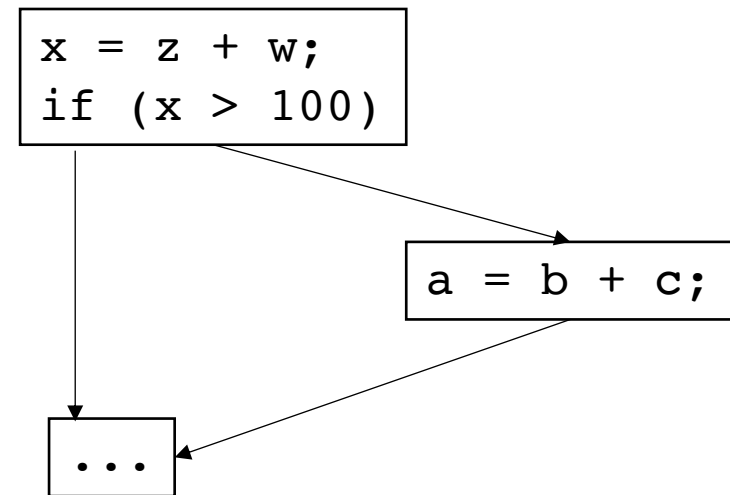
```
r0 = neg(b);  
r1 = b * b;  
r2 = 4 * a;  
r3 = r2 * c;  
r4 = r1 - r3;  
r5 = sqrt(r4);  
r6 = r0 - r5;  
r7 = 2 * a;  
r8 = r6 / r7;  
x  = r8;
```



Control dependencies

```
x = z + w;  
if (x > 100)  
    a = b + c;
```

*Instructions in different CFG nodes
have control-dependencies*



Memory dependencies

True dependence:

Read-after-write

```
a[i] = z + w;
```

```
x = a[i]
```

Memory dependencies

True dependence:
Read-after-write

```
a[i] = z + w;  
x = a[i]
```

Output dependence:
Write-after-write

```
a[i] = z + w;  
a[i] = a + b;
```

Memory dependencies

True dependence:
Read-after-write

```
a[i] = z + w;  
x = a[i]
```

Output dependence:
Write-after-write

```
a[i] = z + w;  
a[i] = a + b;
```

anti-dependence:
Write-after-read

```
x = a[i]  
a[i] = z + w;
```


Memory dependencies

True dependence:
Read-after-write

```
a[i] = z + w;  
x = a[i]
```

Output dependence:
Write-after-write

```
a[i] = z + w;  
a[i] = a + b;
```

anti-dependence:
Write-after-read

```
x = a[i]  
a[i] = z + w;
```

Dependencies can be
removed

```
reg_a_i = z + w;  
a[i] = a + b;
```

Dependencies can be
delayed

```
x = a[i]  
reg_a_i = z + w;  
...  
a[i] = reg_a_i;
```

Memory dependencies

*All of this depends on
accurate pointer analysis!*

True dependence:
Read-after-write

```
a[i] = z + w;  
x = a[i]
```

Output dependence:
Write-after-write

```
a[i] = z + w;  
a[i] = a + b;
```

anti-dependence:
Write-after-read

```
x = a[i]  
a[i] = z + w;
```

Dependencies can be
removed

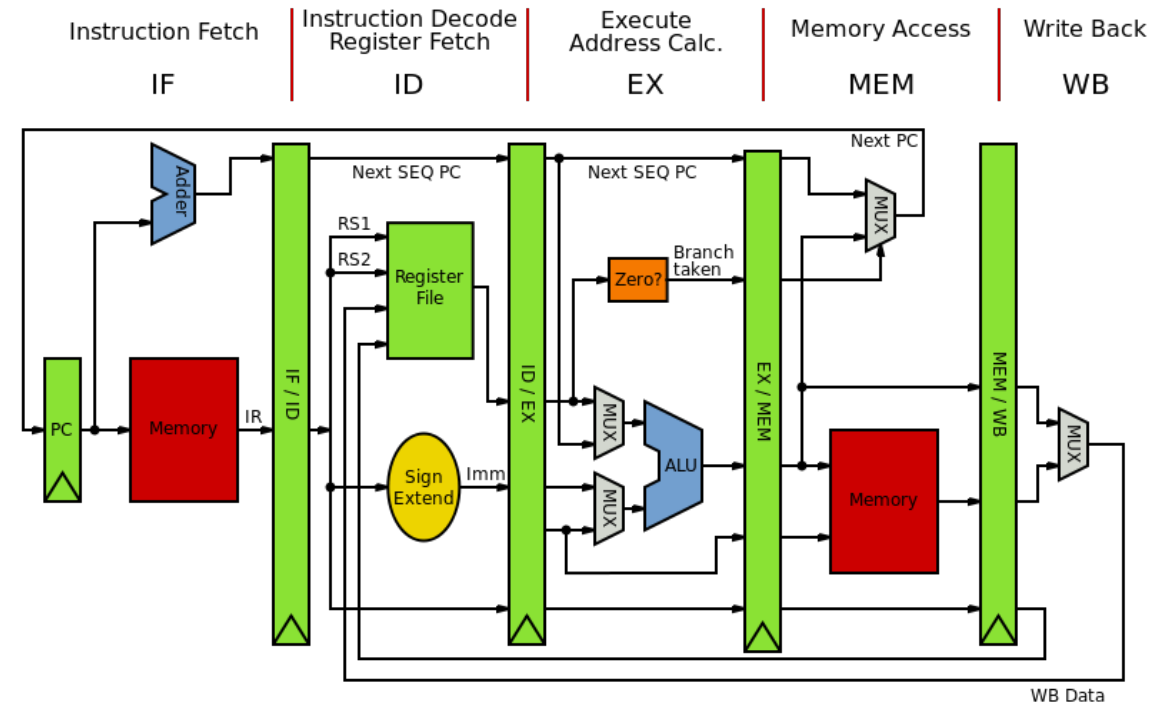
```
reg_a_i = z + w;  
a[i] = a + b;
```

Dependencies can be
delayed

```
x = a[i]  
reg_a_i = z + w;  
...  
a[i] = reg_a_i;
```

How can hardware execute ILP?

- Pipeline parallelism
- Abstract mental model for compiler:
 - N-stage pipeline
 - N instructions can be in-flight
 - Dependencies stall pipeline

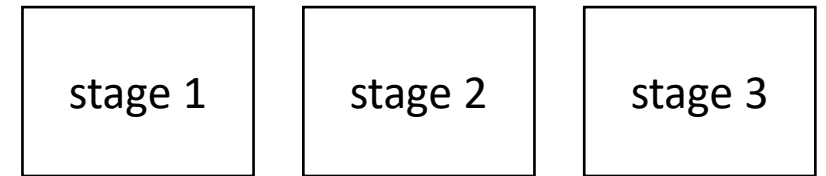


MIPS pipeline image from:
[https://commons.wikimedia.org/wiki/Pipeline_\(computer_hardware\)](https://commons.wikimedia.org/wiki/Pipeline_(computer_hardware))

How can hardware execute ILP?

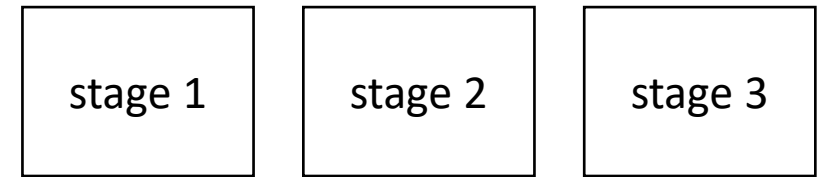
- Pipeline parallelism
- Abstract mental model for compiler:
 - N-stage pipeline
 - N instructions can be in-flight
 - Dependencies stall pipeline

```
instr1;  
instr2;  
instr3;
```



How can hardware execute ILP?

- Pipeline parallelism
- Abstract mental model for compiler:
 - N-stage pipeline
 - N instructions can be in-flight
 - Dependencies stall pipeline



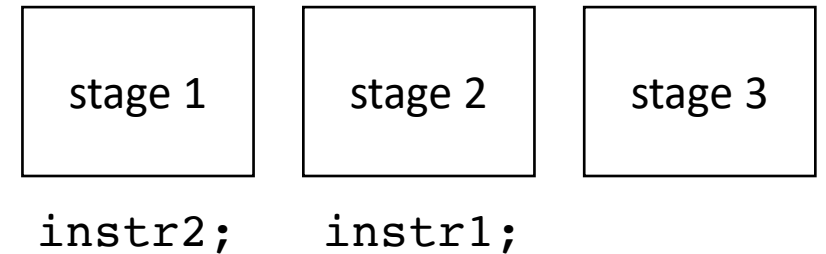
`instr1;`

`instr2;`

`instr3;`

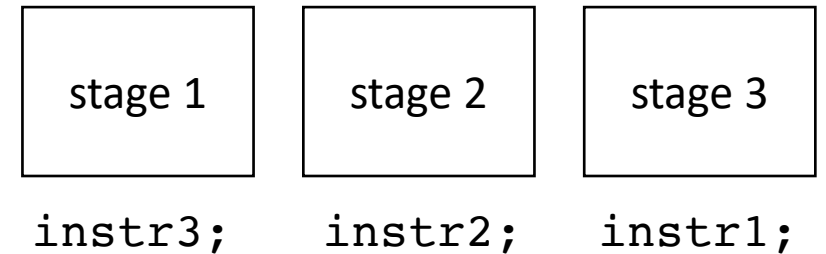
How can hardware execute ILP?

- Pipeline parallelism
- Abstract mental model for compiler:
 - N-stage pipeline
 - N instructions can be in-flight
 - Dependencies stall pipeline



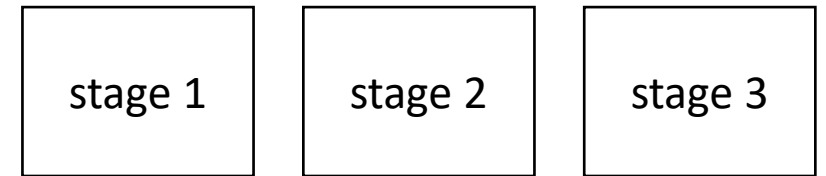
How can hardware execute ILP?

- Pipeline parallelism
- Abstract mental model for compiler:
 - N-stage pipeline
 - N instructions can be in-flight
 - Dependencies stall pipeline



How can hardware execute ILP?

- Pipeline parallelism
- Abstract mental model for compiler:
 - N-stage pipeline
 - N instructions can be in-flight
 - Dependencies stall pipeline

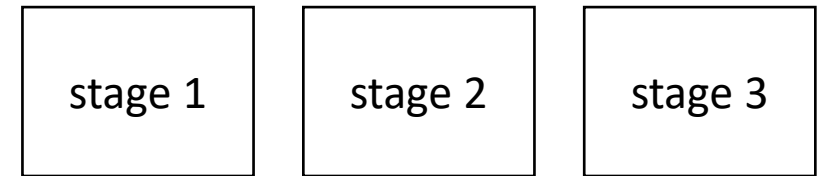


6 cycles for 3 independent instructions

Converges to 1 instruction per cycle

How can hardware execute ILP?

- Pipeline parallelism
- Abstract mental model for compiler:
 - N-stage pipeline
 - N instructions can be in-flight
 - Dependencies stall pipeline

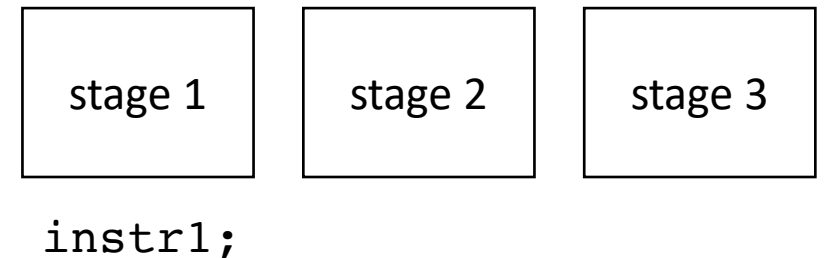


```
instr1;  
instr2;  
instr3;
```

What if the instructions depend on each other?

How can hardware execute ILP?

- Pipeline parallelism
- Abstract mental model for compiler:
 - N-stage pipeline
 - N instructions can be in-flight
 - Dependencies stall pipeline

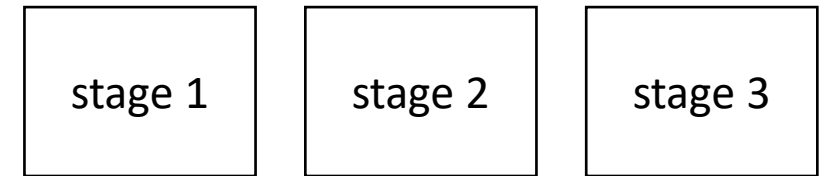


`instr2;`
`instr3;`

What if the instructions depend on each other?

How can hardware execute ILP?

- Pipeline parallelism
- Abstract mental model for compiler:
 - N-stage pipeline
 - N instructions can be in-flight
 - Dependencies stall pipeline



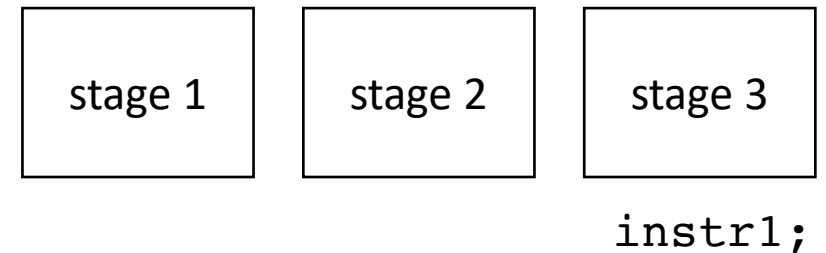
`instr1;`

`instr2;`
`instr3;`

What if the instructions depend on each other?

How can hardware execute ILP?

- Pipeline parallelism
- Abstract mental model for compiler:
 - N-stage pipeline
 - N instructions can be in-flight
 - Dependencies stall pipeline

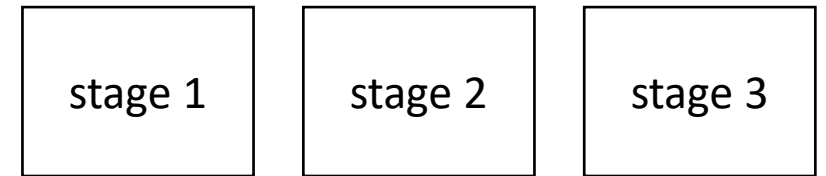


instr2;
instr3;

What if the instructions depend on each other?

How can hardware execute ILP?

- Pipeline parallelism
- Abstract mental model for compiler:
 - N-stage pipeline
 - N instructions can be in-flight
 - Dependencies stall pipeline

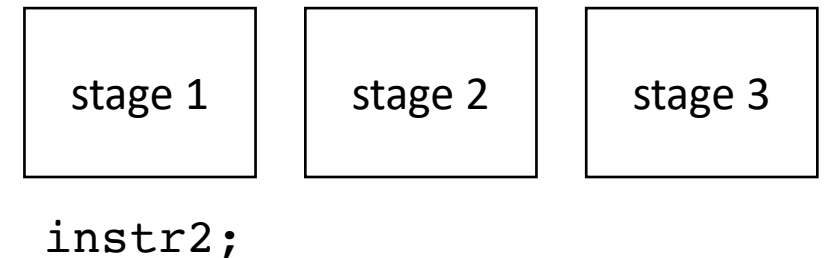


```
instr2;  
instr3;
```

What if the instructions depend on each other?

How can hardware execute ILP?

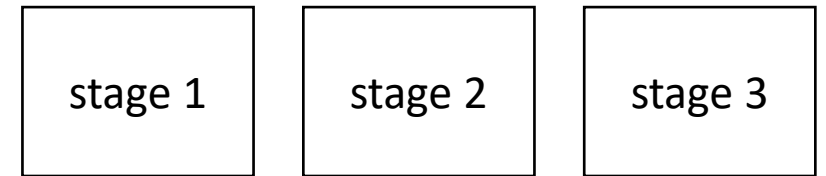
- Pipeline parallelism
- Abstract mental model for compiler:
 - N-stage pipeline
 - N instructions can be in-flight
 - Dependencies stall pipeline



What if the instructions depend on each other?

How can hardware execute ILP?

- Pipeline parallelism
- Abstract mental model for compiler:
 - N-stage pipeline
 - N instructions can be in-flight
 - Dependencies stall pipeline



What if the instructions depend on each other?

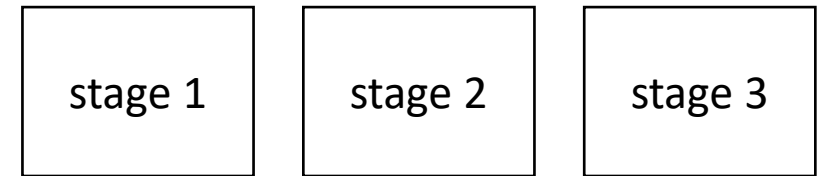
9 cycles for 3 instructions

converges to 3 cycles per instruction

How can hardware execute ILP?

- Pipeline parallelism
- Abstract mental model for compiler:
 - N-stage pipeline
 - N instructions can be in-flight
 - Dependencies stall pipeline

```
instr1;  
instrX0;  
instrX1;  
instr2;  
instrX2;  
instrX3;  
instr3;
```

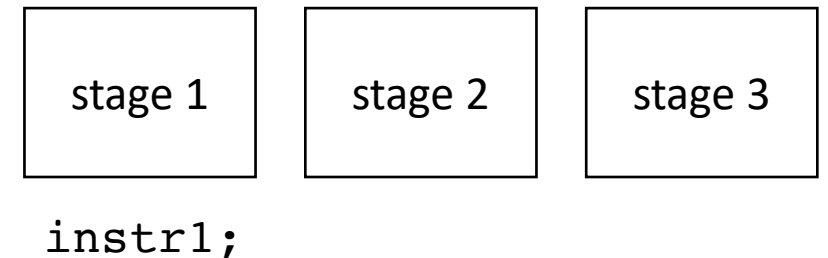


If there are non-dependent instructions from other places in the program that we can interleave then we can get back performance!

How can hardware execute ILP?

- Pipeline parallelism
- Abstract mental model for compiler:
 - N-stage pipeline
 - N instructions can be in-flight
 - Dependencies stall pipeline

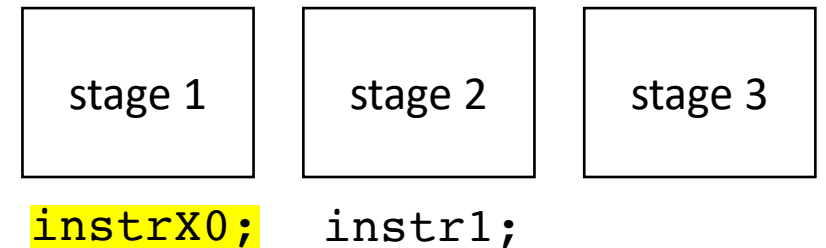
```
instrX0;  
instrX1;  
instr2;  
instrX2;  
instrX3;  
instr3;
```



If there are non-dependent instructions from other places in the program that we can interleave then we can get back performance!

How can hardware execute ILP?

- Pipeline parallelism
- Abstract mental model for compiler:
 - N-stage pipeline
 - N instructions can be in-flight
 - Dependencies stall pipeline

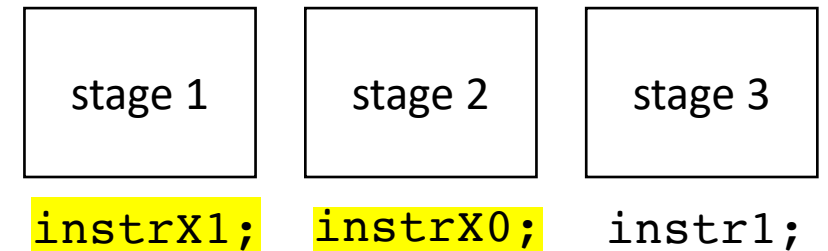


```
instrX1;  
instr2;  
instrX2;  
instrX3;  
instr3;
```

If there are non-dependent instructions from other places in the program that we can interleave then we can get back performance!

How can hardware execute ILP?

- Pipeline parallelism
- Abstract mental model for compiler:
 - N-stage pipeline
 - N instructions can be in-flight
 - Dependencies stall pipeline

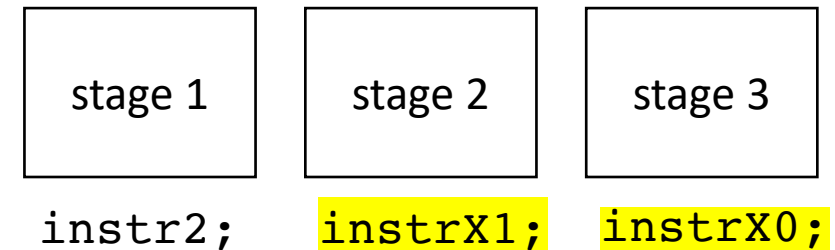


```
instr2;  
instrX2;  
instrX3;  
instr3;
```

If there are non-dependent instructions from other places in the program that we can interleave then we can get back performance!

How can hardware execute ILP?

- Pipeline parallelism
- Abstract mental model for compiler:
 - N-stage pipeline
 - N instructions can be in-flight
 - Dependencies stall pipeline

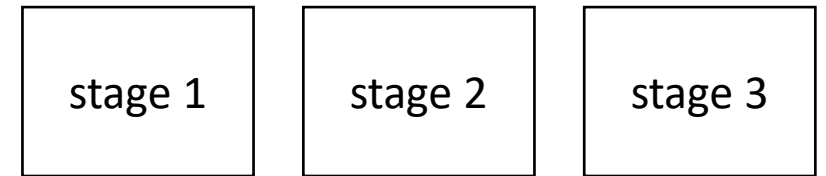


instrX2;
instrX3;
instr3;

If there are non-dependent instructions from other places in the program that we can interleave then we can get back performance!

How can hardware execute ILP?

- Pipeline parallelism
- Abstract mental model for compiler:
 - N-stage pipeline
 - N instructions can be in-flight
 - Dependencies stall pipeline

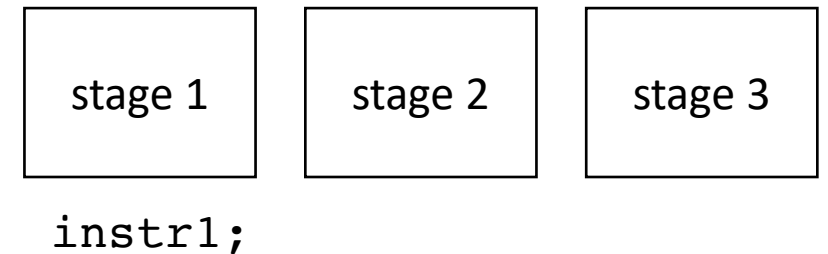


```
instr1;  
instr2;  
instr3;
```

*Say instr2; and instr3;
have a control
dependence on instr1;*

How can hardware execute ILP?

- Pipeline parallelism
- Abstract mental model for compiler:
 - N-stage pipeline
 - N instructions can be in-flight
 - Dependencies stall pipeline

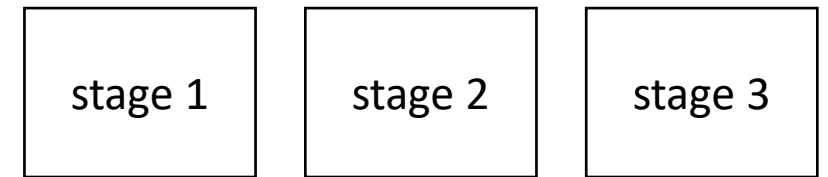


`instr2;`
`instr3;`

*Say `instr2;` and `instr3;`
have a control
dependence on `instr1;`*

How can hardware execute ILP?

- Pipeline parallelism
- Abstract mental model for compiler:
 - N-stage pipeline
 - N instructions can be in-flight
 - Dependencies stall pipeline



`instr2;`

`instr1;`

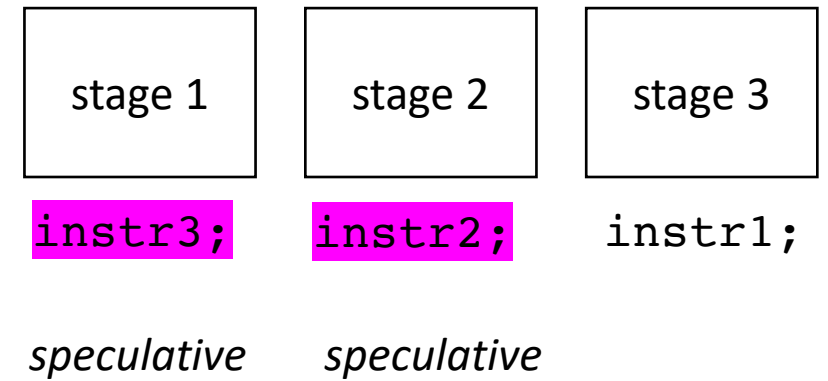
speculative

`instr3;`

*Say `instr2;` and `instr3;`
have a control
dependence on `instr1;`*

How can hardware execute ILP?

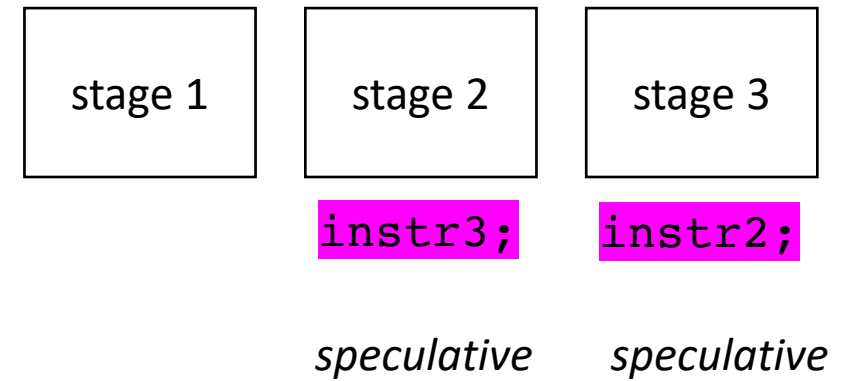
- Pipeline parallelism
- Abstract mental model for compiler:
 - N-stage pipeline
 - N instructions can be in-flight
 - Dependencies stall pipeline



*Say instr2; and instr3;
have a control
dependence on instr1;*

How can hardware execute ILP?

- Pipeline parallelism
- Abstract mental model for compiler:
 - N-stage pipeline
 - N instructions can be in-flight
 - Dependencies stall pipeline



*Say instr2; and instr3;
have a control
dependence on instr1;*

*before we commit
the speculative instructions,
we check if the control
dependence was satisfied.*

How can hardware execute ILP?

- Executing multiple instructions at once:
- Very Long Instruction Word (VLIW) architecture
 - Multiple instructions are combined into one by the compiler
- Superscalar architecture:
 - Several sequential operations are issued in parallel

How can hardware execute ILP?

- Executing multiple instructions at once:
- Superscalar architecture:
 - Several sequential operations are issued in parallel
 - hardware detects dependencies

```
instr0;  
instr1;  
instr2;
```

issue-width is maximum number of instructions that can be issued in parallel

How can hardware execute ILP?

- Executing multiple instructions at once:
- Superscalar architecture:
 - Several sequential operations are issued in parallel
 - hardware detects dependencies

```
instr0;
```

```
instr1;
```

```
instr2;
```

issue-width is maximum number of instructions that can be issued in parallel

if instr0 and instr1 are independent, they will be issued in parallel

It's even more complicated

- Out-of-order execution delays dependent instructions
 - Reorder buffers (RoB) track dependencies
 - Load-Store Queues (LSQ) hold outstanding memory requests

What does this look like in the real world?

- Intel Haswell (2013):
 - Issue width of 4
 - 14-19 stage pipeline
 - OoO execution
- Intel Nehalem (2008)
 - 20-24 stage pipeline
 - Issue width of 2-4
 - OoO execution
- ARM
 - V7 has 3 stage pipeline; Cortex V8 has 13
 - Cortex V8 has issue width of 2
 - OoO execution
- RISC-V
 - Ariane and Rocket are In-Order
 - 3-6 stage pipelines
 - some super scaler implementations (BOOM)

What does this mean for compiler writers?

- We should have an abstract and parametrized performance model for instruction scheduling (the order of instructions)
- Try not to place dependent instructions in sequence
- *Above all, instructions must respect sequential semantics!*

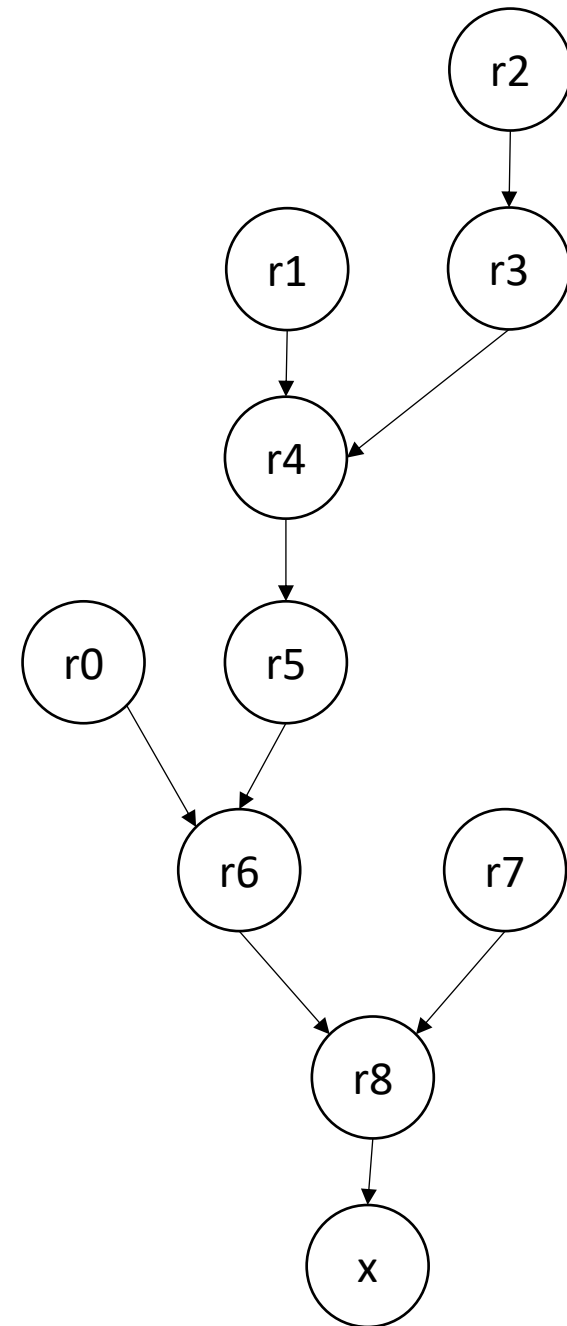
Four compiler techniques for better ILP

- Priority topological ordering
- Anticipatable expressions
- Independent for loops
- Reduction for loops

Priority Topological Ordering of DDGs for Superscalar

```
r0 = neg(b);  
r1 = b * b;  
r2 = 4 * a;  
r3 = r2 * c;  
r4 = r1 - r3;  
r5 = sqrt(r4);  
r6 = r0 - r5;  
r7 = 2 * a;  
r8 = r6 / r7;  
x = r8;
```

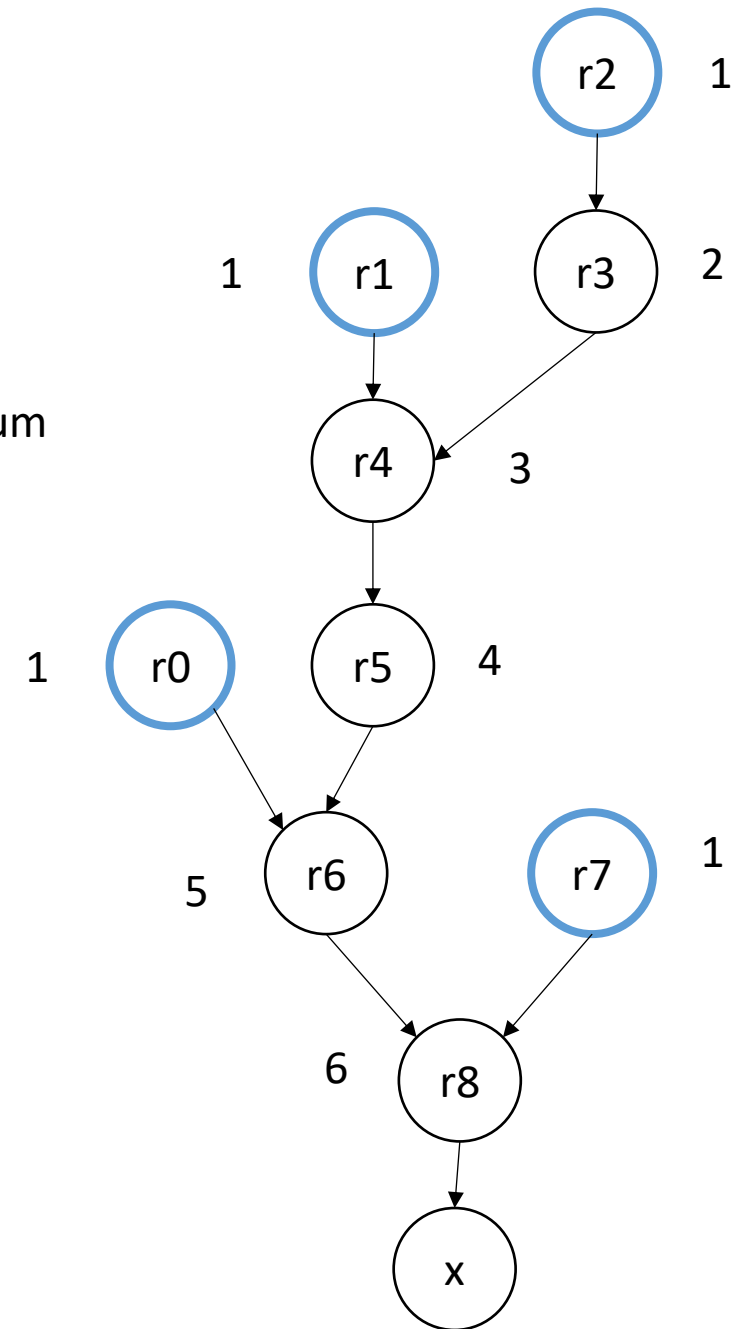
First, consider optimizing for superscalar



Priority Topological Ordering of DDGs for Superscalar

```
r0 = neg(b);  
r1 = b * b;  
r2 = 4 * a;  
r3 = r2 * c;  
r4 = r1 - r3;  
r5 = sqrt(r4);  
r6 = r0 - r5;  
r7 = 2 * a;  
r8 = r6 / r7;  
x = r8;
```

Label nodes with the maximum distance to a source

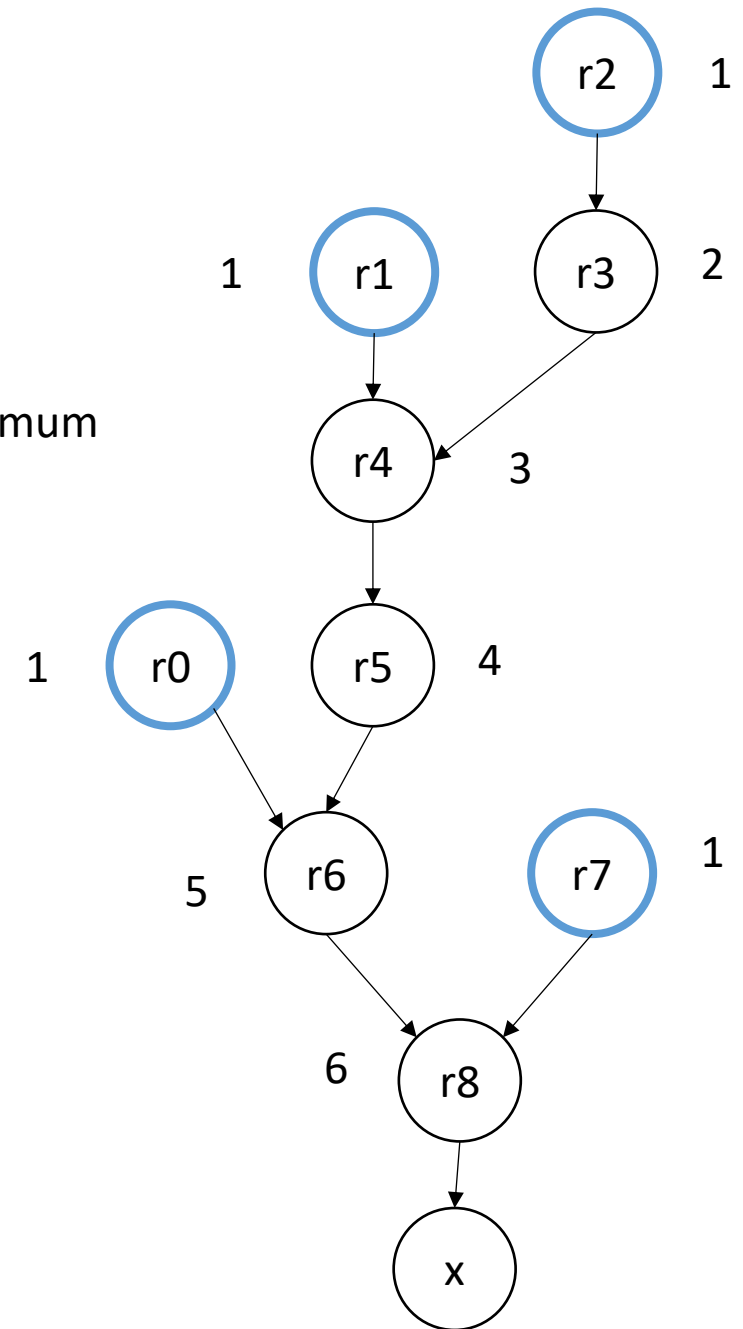


Priority Topological Ordering of DDGs for Superscalar

```
r0 = neg(b);  
r1 = b * b;  
r2 = 4 * a;  
r3 = r2 * c;  
r4 = r1 - r3;  
r5 = sqrt(r4);  
r6 = r0 - r5;  
r7 = 2 * a;  
r8 = r6 / r7;  
x = r8;
```

number nodes with the maximum distance to a source

Break ties in topological order using this number

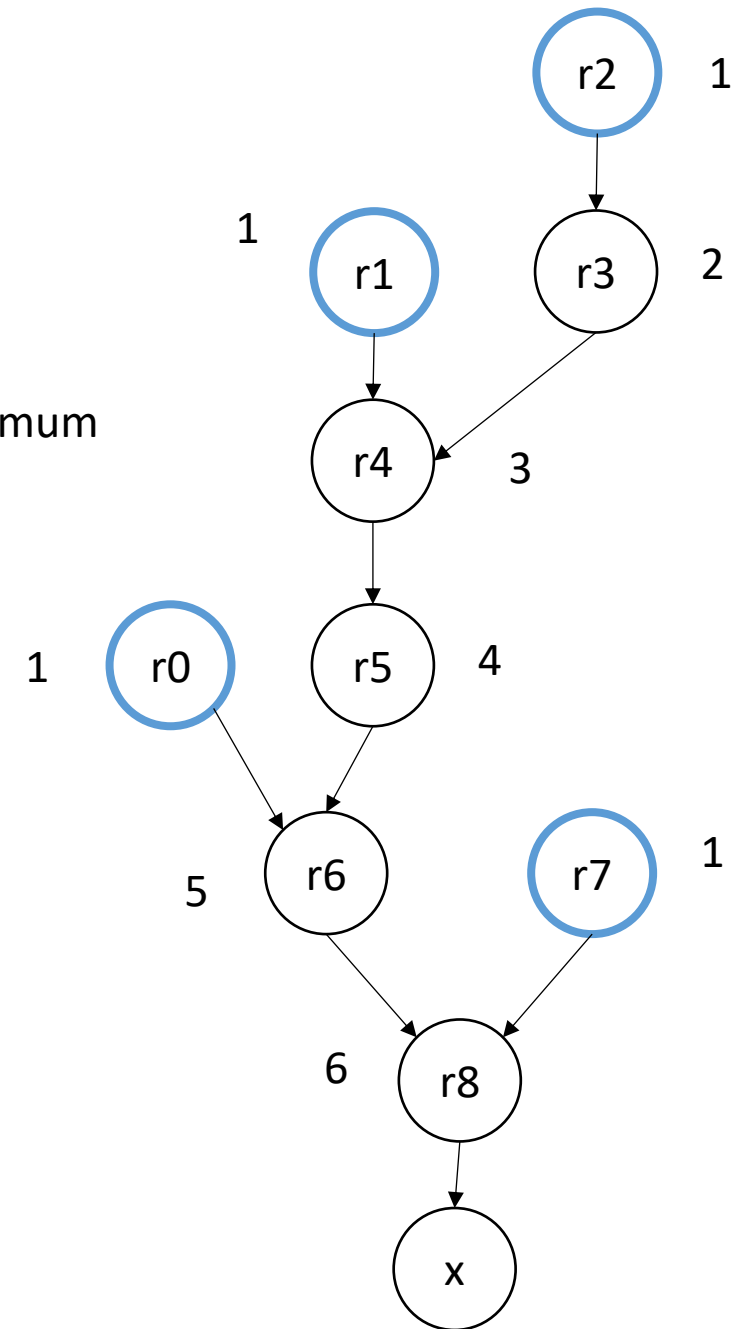


Priority Topological Ordering of DDGs for Superscalar

```
r0 = neg(b);  
r1 = b * b;  
r2 = 4 * a;  
r3 = r2 * c;  
r4 = r1 - r3;  
r5 = sqrt(r4);  
r6 = r0 - r5;  
r7 = 2 * a;  
r8 = r6 / r7;  
x = r8;
```

number nodes with the maximum distance to a source

Break ties in topological order using this number

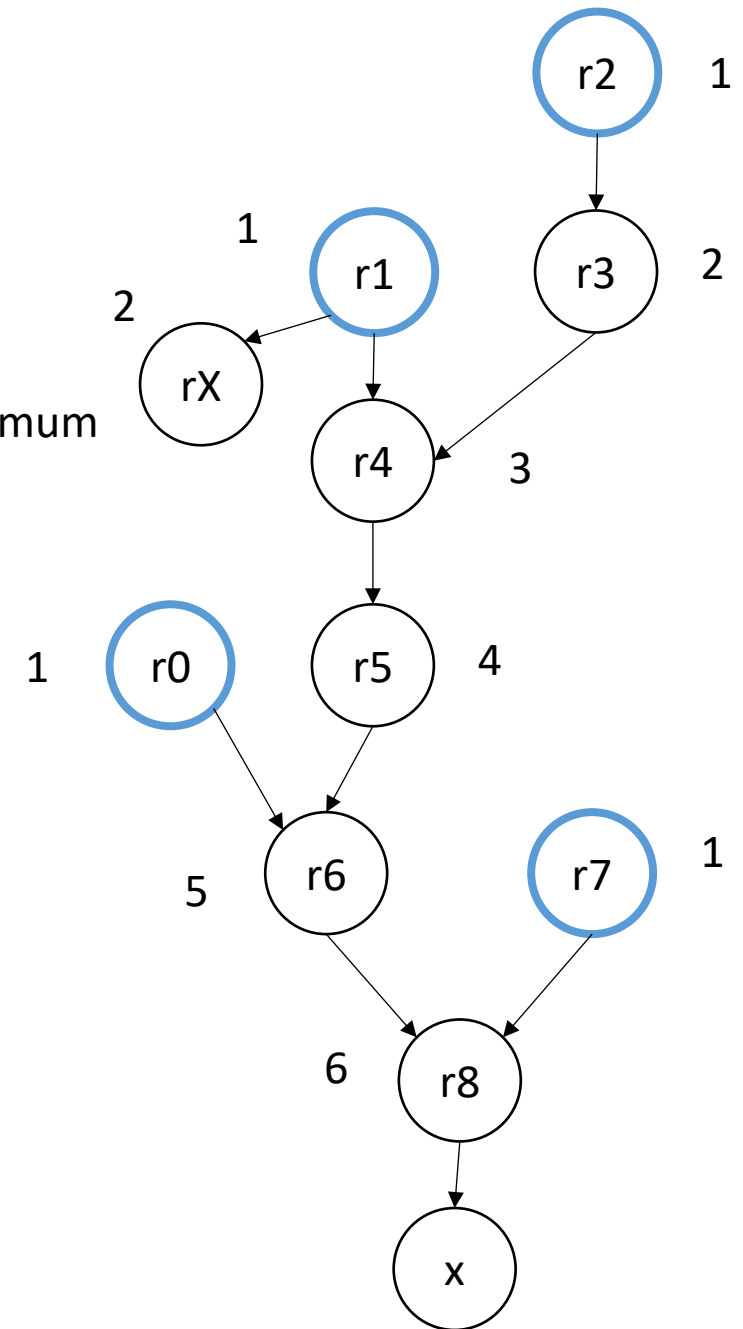


Priority Topological Ordering of DDGs for Superscalar

```
r0 = neg(b);  
r1 = b * b;  
r2 = 4 * a;  
r7 = 2 * a;  
r3 = r2 * c;  
r4 = r1 - r3;  
r5 = sqrt(r4);  
r6 = r0 - r5;  
r8 = r6 / r7;  
x = r8;
```

number nodes with the maximum distance to a source

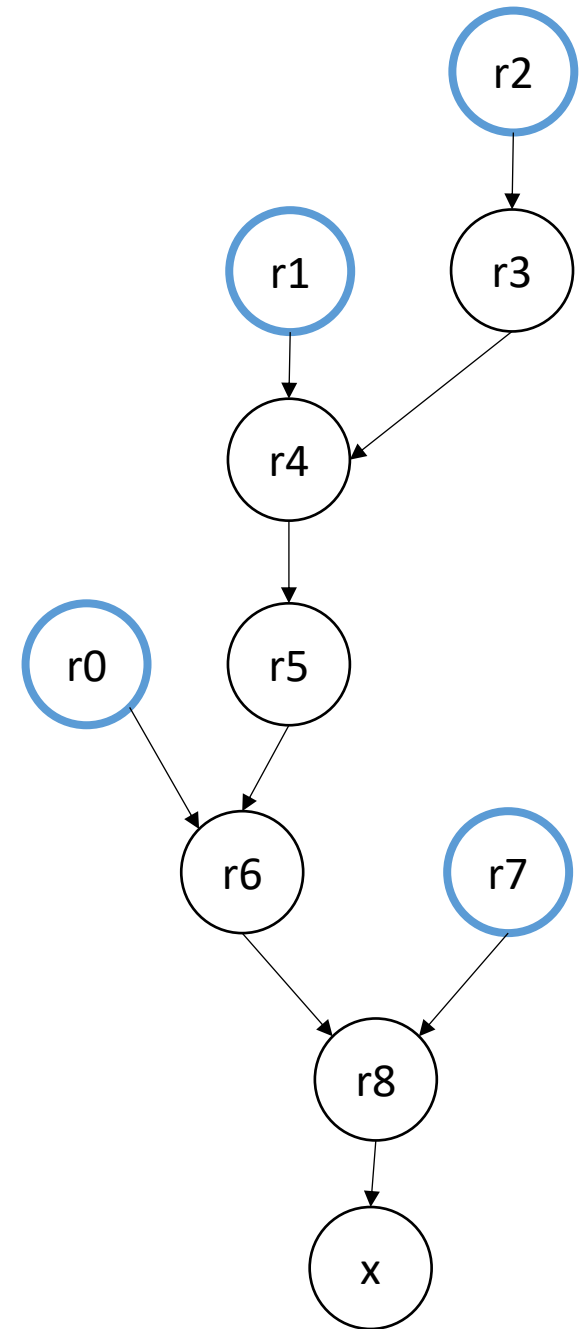
Break ties in topological order using this number



Priority Topological Ordering of DDGs for Pipelining

```
r0 = neg(b);  
r1 = b * b;  
r2 = 4 * a;  
r3 = r2 * c;  
r4 = r1 - r3;  
r5 = sqrt(r4);  
r6 = r0 - r5;  
r7 = 2 * a;  
r8 = r6 / r7;  
x = r8;
```

superscalar should move independent instructions as high as possible. What about pipelining?

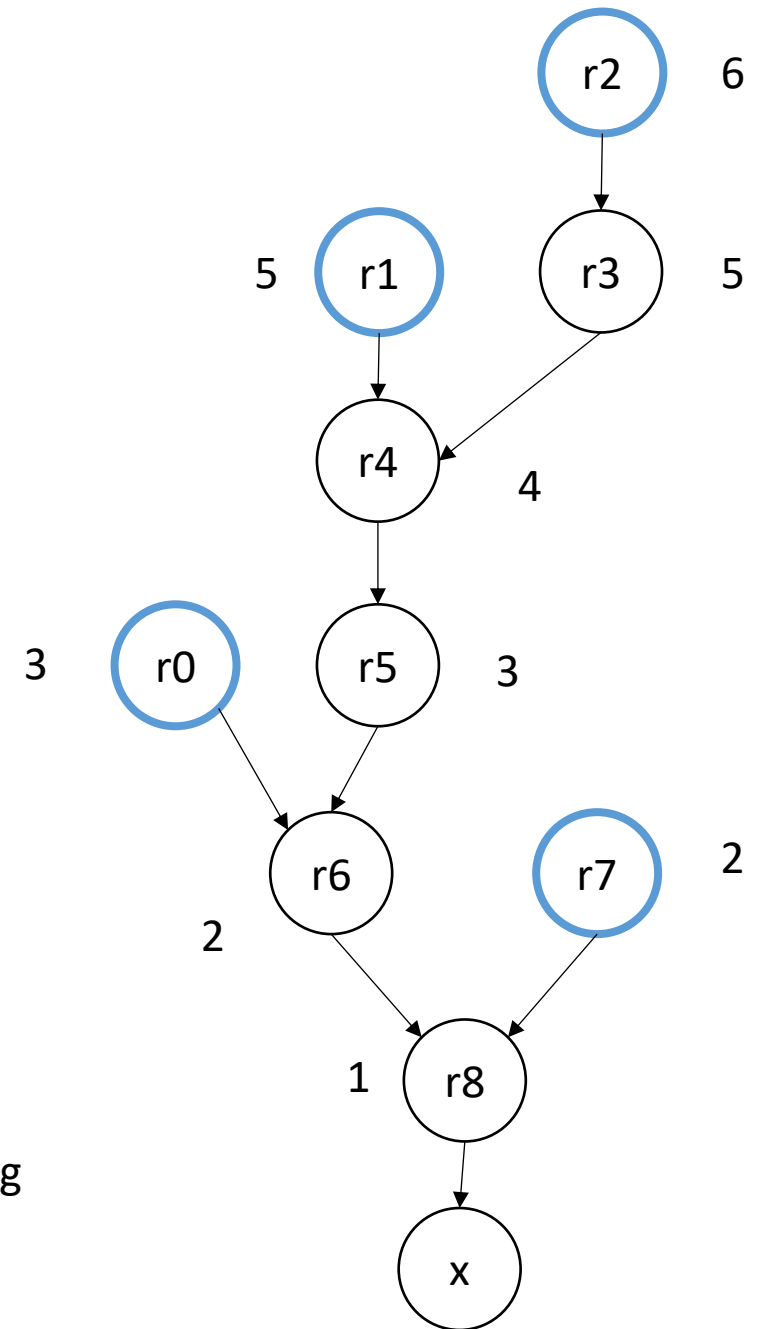


Priority Topological Ordering of DDGs for Pipelining

```
r0 = neg(b);  
r1 = b * b;  
r2 = 4 * a;  
r3 = r2 * c;  
r4 = r1 - r3;  
r5 = sqrt(r4);  
r6 = r0 - r5;  
r7 = 2 * a;  
r8 = r6 / r7;  
x = r8;
```

superscalar should move independent instructions as high as possible. What about pipelining?

label each node with a distance from the root. Schedule each node according to the level

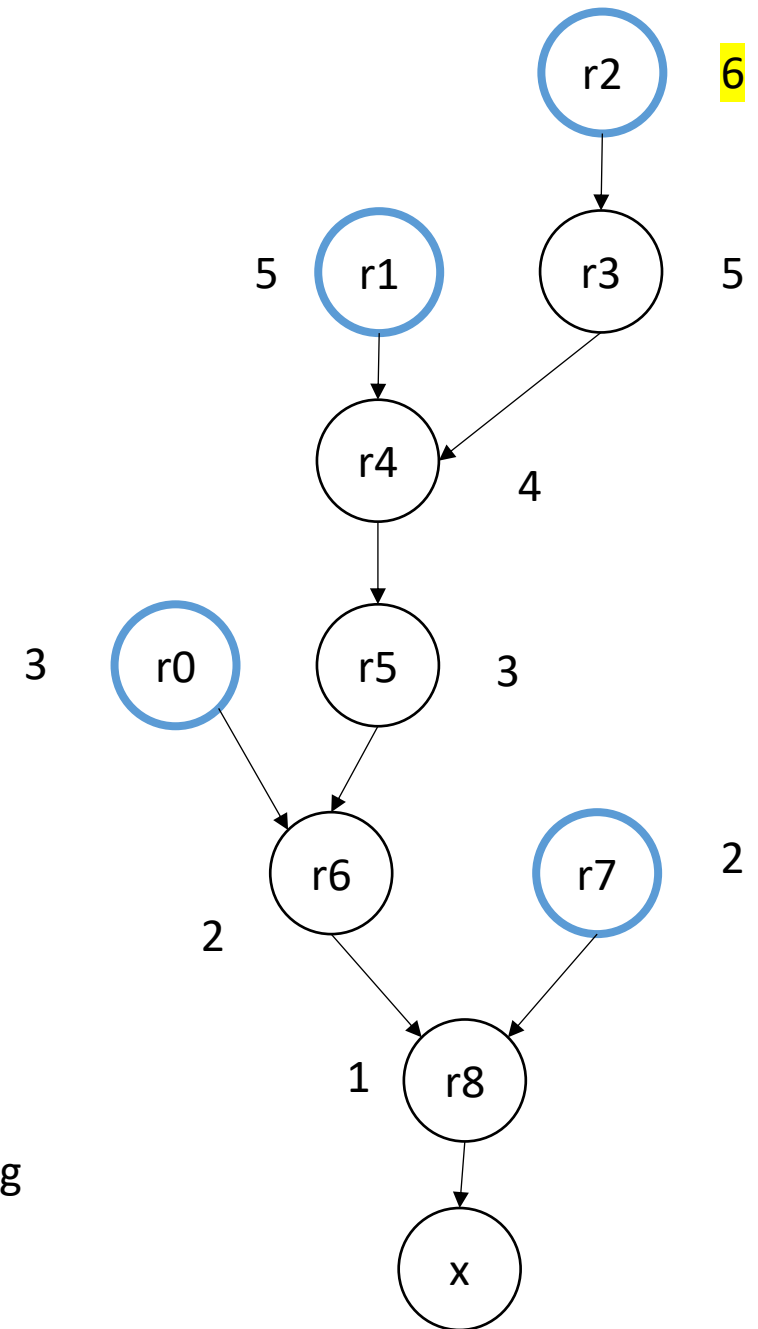


Priority Topological Ordering of DDGs for Pipelining

```
r0 = neg(b);  
r1 = b * b;  
r2 = 4 * a;  
r3 = r2 * c;  
r4 = r1 - r3;  
r5 = sqrt(r4);  
r6 = r0 - r5;  
r7 = 2 * a;  
r8 = r6 / r7;  
x = r8;
```

superscalar should move independent instructions as high as possible. What about pipelining?

label each node with a distance from the root. Schedule each node according to the level



Priority Topological Ordering of DDGs for Pipelining

```
r2 = 4 * a;
```

```
r0 = neg(b);
```

```
r1 = b * b;
```

```
r3 = r2 * c;
```

```
r4 = r1 - r3;
```

```
r5 = sqrt(r4);
```

```
r6 = r0 - r5;
```

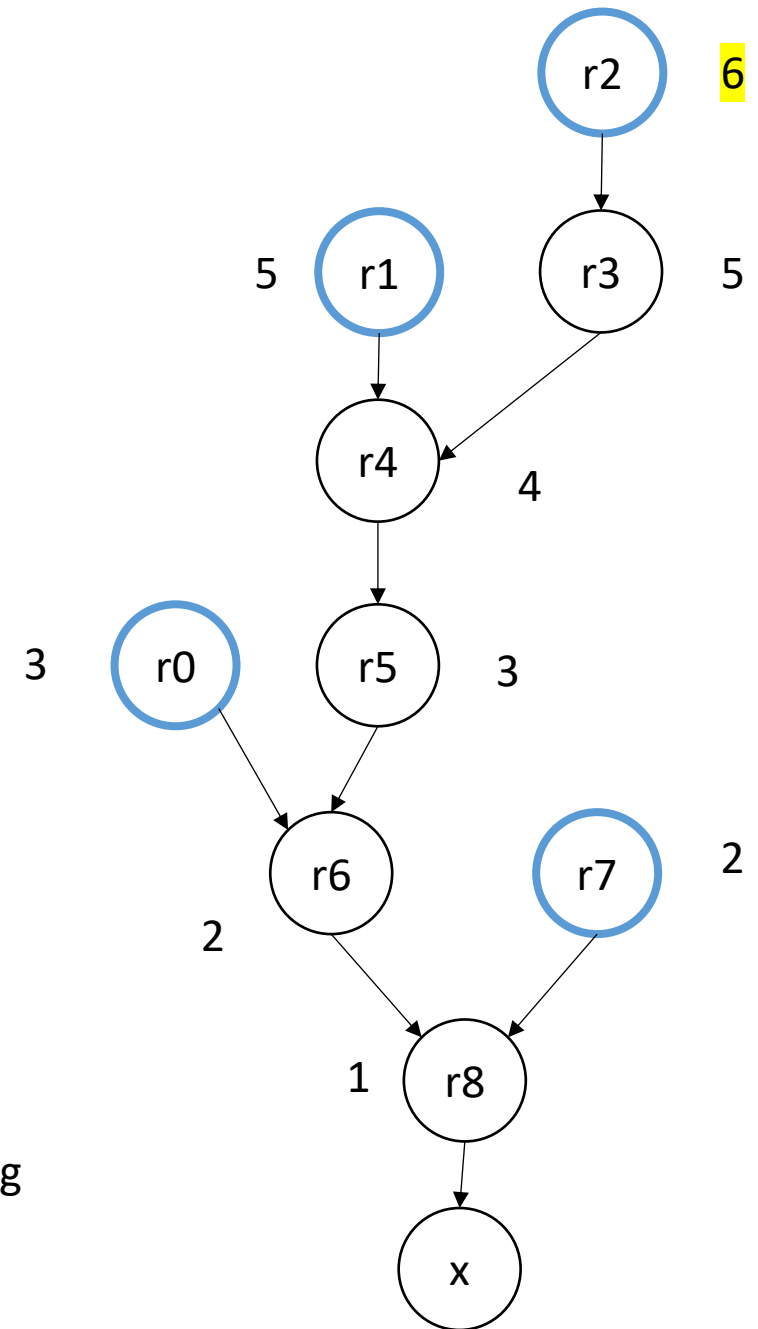
```
r7 = 2 * a;
```

```
r8 = r6 / r7;
```

```
x = r8;
```

superscalar should move independent instructions as high as possible. What about pipelining?

label each node with a distance from the root. Schedule each node according to the level

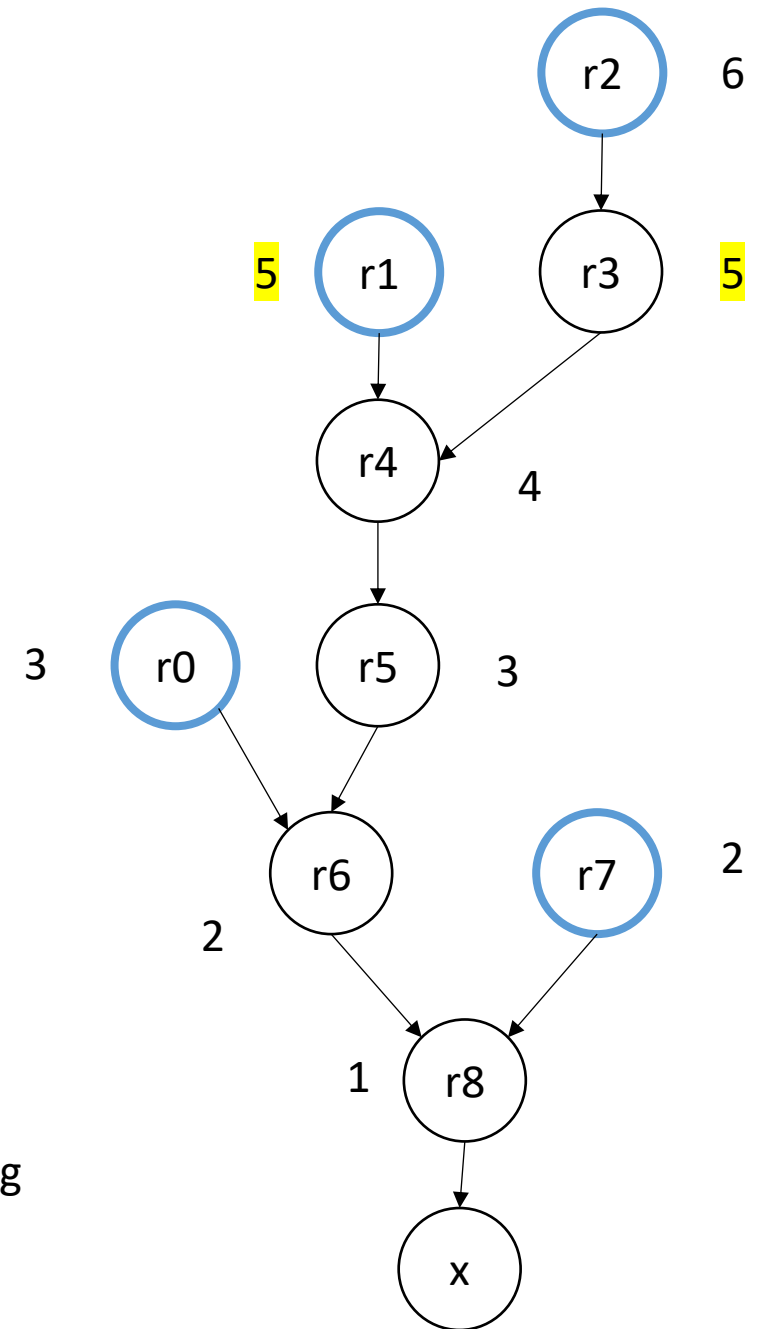


Priority Topological Ordering of DDGs for Pipelining

```
r2 = 4 * a;  
r0 = neg(b);  
r1 = b * b;  
r3 = r2 * c;  
r4 = r1 - r3;  
r5 = sqrt(r4);  
r6 = r0 - r5;  
r7 = 2 * a;  
r8 = r6 / r7;  
x = r8;
```

Ties are broken with the node that has the least parents

label each node with a distance from the root.
Schedule each node according to the level

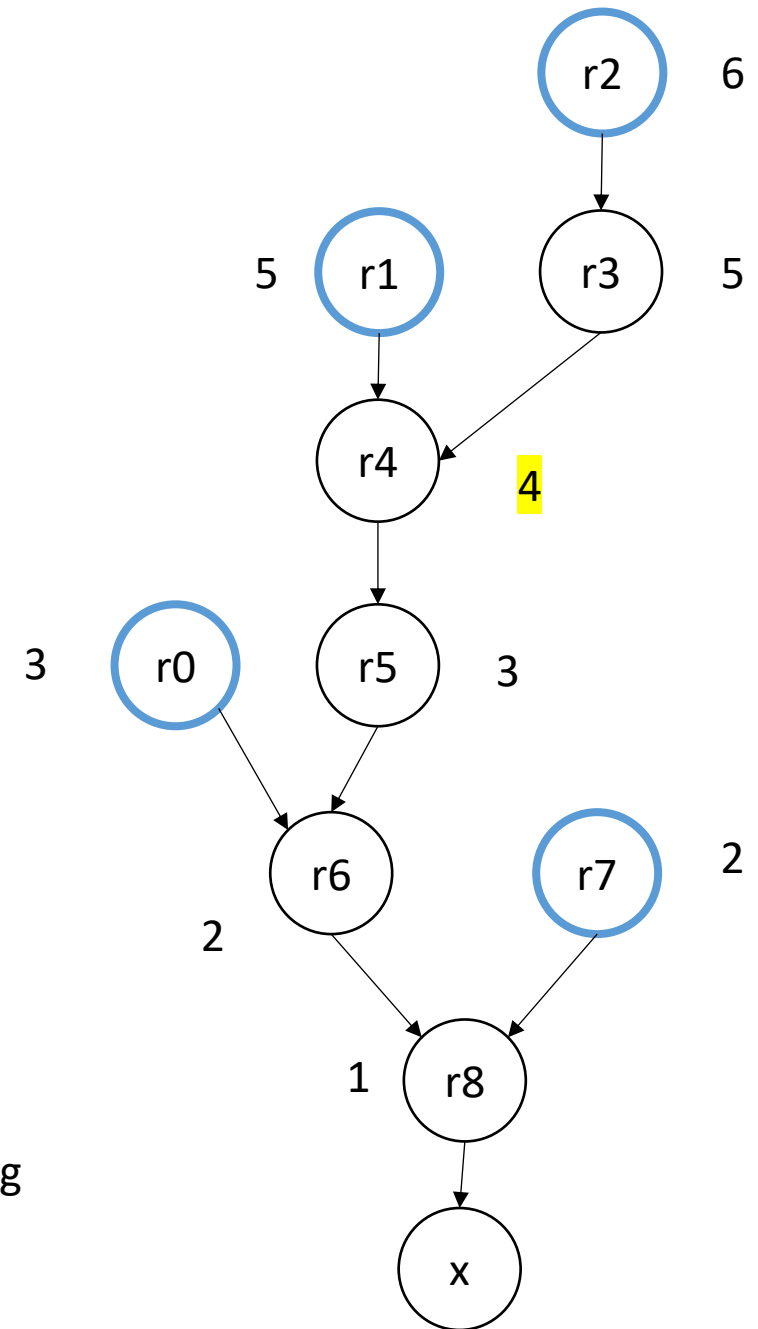


Priority Topological Ordering of DDGs for Pipelining

```
r2 = 4 * a;  
r1 = b * b;  
r3 = r2 * c;  
r0 = neg(b);  
r4 = r1 - r3;  
r5 = sqrt(r4);  
r6 = r0 - r5;  
r7 = 2 * a;  
r8 = r6 / r7;  
x = r8;
```

Ties are broken with the node that has the least parents

label each node with a distance from the root.
Schedule each node according to the level



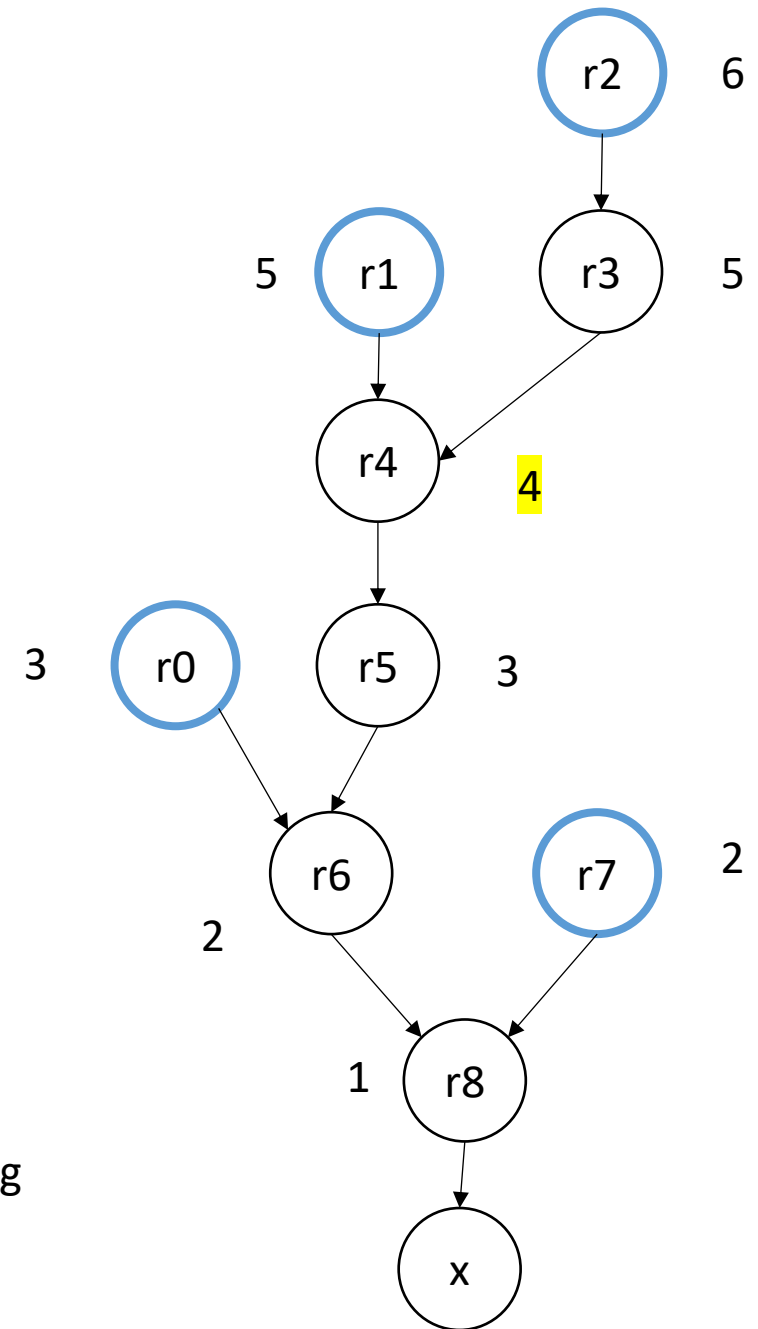
Priority Topological Ordering of DDGs for Pipelining

final

```
r2 = 4 * a;  
r1 = b * b;  
r3 = r2 * c;  
r4 = r1 - r3;  
r0 = neg(b);  
r5 = sqrt(r4);  
r7 = 2 * a;  
r6 = r0 - r5;  
r8 = r6 / r7;  
x = r8;
```

Ties are broken with the node that has the least parents

label each node with a distance from the root.
Schedule each node according to the level



In practice

- Some hybrid approach depending on the target
- Some approaches use a resource model that explicitly encode the issue-width and pipeline

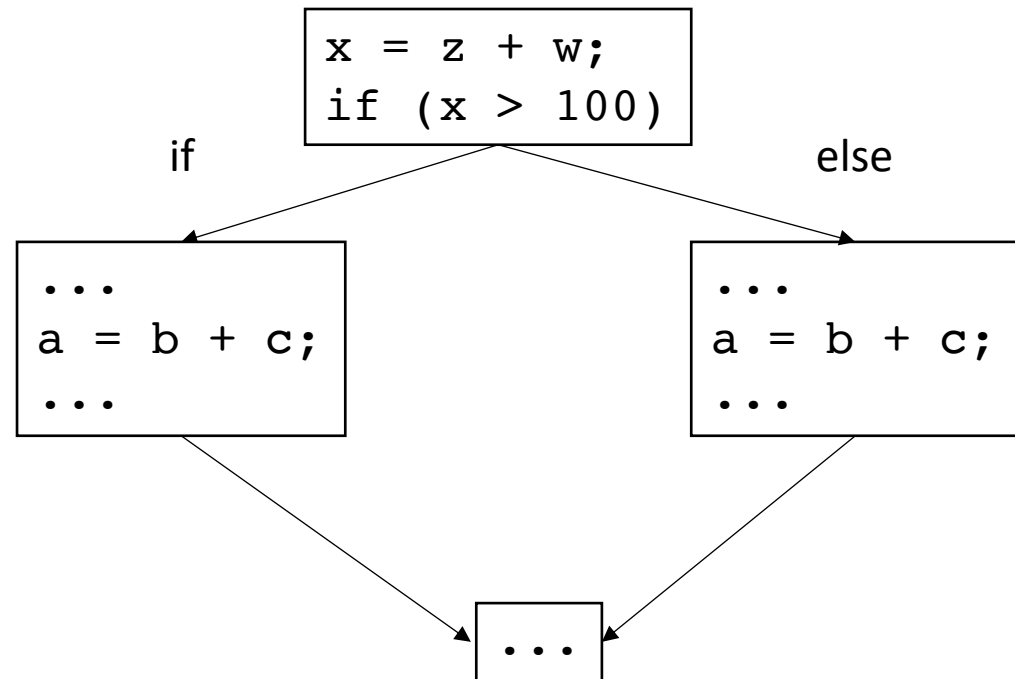
Anticipable Expressions

$$AntOut(n) = \bigcap_{s \text{ in succ}} UExpr(s) \cup (AntOut(s) \cap \overline{ExprKill(s)})$$

An expression e is “anticipable” at a basic block b_x if for all paths that leave b_x , e is evaluated

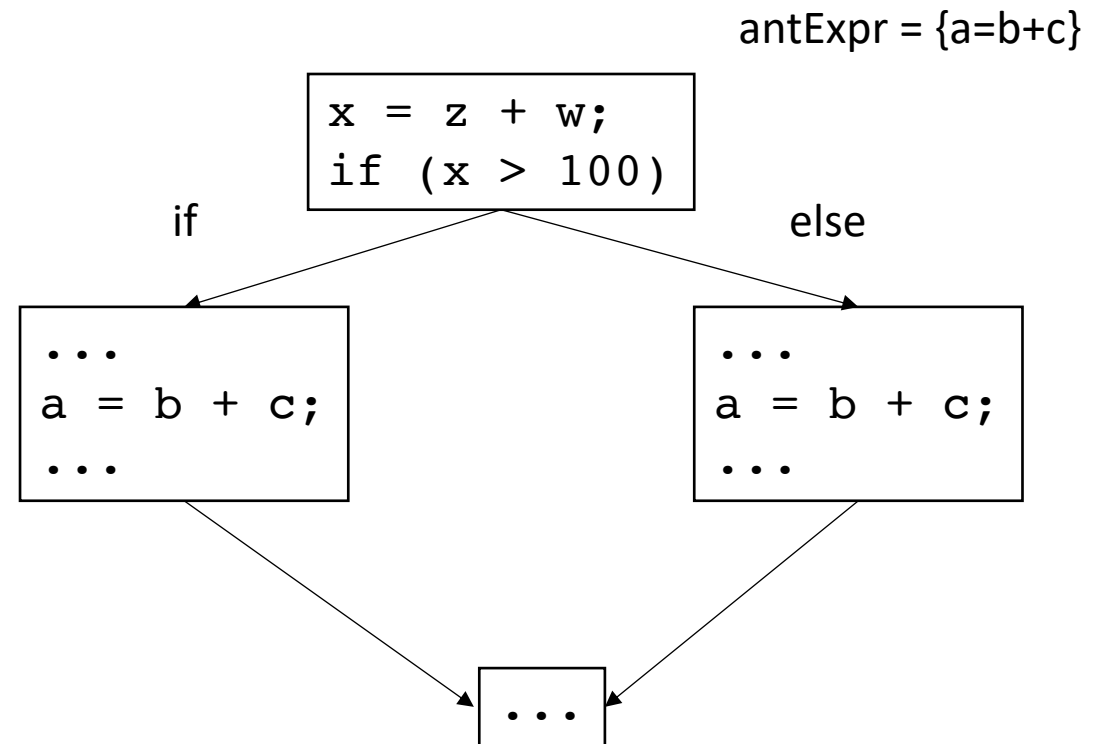
Anticipable Expressions

```
x = z + w;  
if (x > 100) {  
    ...  
    a = b + c;  
    ...  
}  
else {  
    ...  
    a = b + c;  
    ...  
}
```



Anticipable Expressions

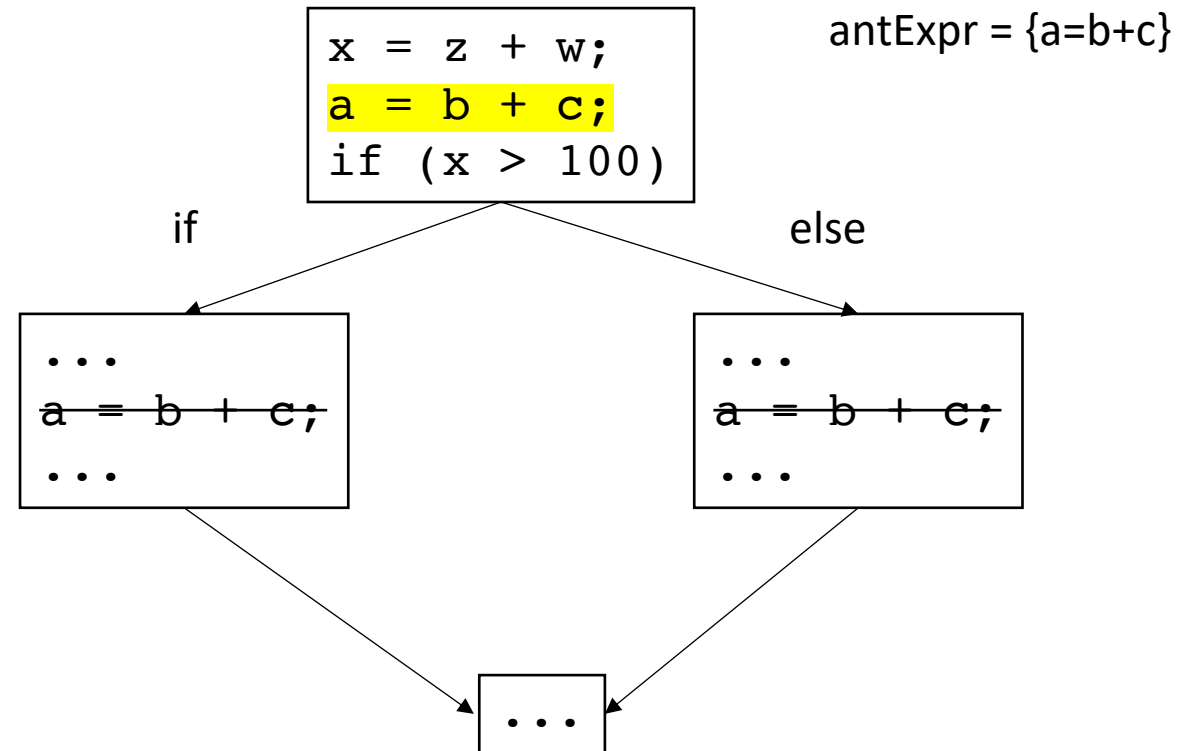
```
x = z + w;  
if (x > 100) {  
    ...  
    a = b + c;  
    ...  
}  
else {  
    ...  
    a = b + c;  
    ...  
}
```



Anticipable Expressions

also called "Upward code motion"

```
x = z + w;  
a = b + c;  
if (x > 100) {  
  ...  
  a = b + c;  
  ...  
}  
else {  
  ...  
  a = b + c;  
  ...  
}
```



Using Loop Unrolling to Exploit ILP

- for loops with independent chains of computation

```
for (int i = 0; i < SIZE; i++) {  
    SEQ(i);  
}
```

where: `SEQ(i) = instr1;`

`instr2;`

`...`

`a[i] = instrN;`

and let `instr(N)` depends on `instr(N-1)`

loops only write to memory
addressed by the loop variable

Using Loop Unrolling to Exploit ILP

- Simple loop unrolling:

```
for (int i = 0; i < SIZE; i+=2) {  
    SEQ(i);  
    SEQ(i+1);  
}
```

Saves one addition and one comparison per loop, but doesn't help with ILP

Using Loop Unrolling to Exploit ILP

- Simple loop unrolling:

```
for (int i = 0; i < SIZE; i+=2) {  
    SEQ(i);  
    SEQ(i+1);  
}
```

Let $SEQ(i, j)$ be the j th instruction of $SEQ(i)$.

Let each instruction chain have N instructions

Using Loop Unrolling to Exploit ILP

- Simple loop unrolling:

```
for (int i = 0; i < SIZE; i+=2) {  
    SEQ(i,1);  
    SEQ(i+1,1);  
    SEQ(i,2);  
    SEQ(i+1,2);  
    ...  
    SEQ(i,N);  
    SEQ(i+1, N);  
}
```

Let $SEQ(i, j)$ be the j th instruction of $SEQ(i)$.

Let each instruction chain have N instructions

Using Loop Unrolling to Exploit ILP

- Simple loop unrolling:

```
for (int i = 0; i < SIZE; i+=2) {  
    SEQ(i, 1);  
    SEQ(i+1, 1);  
    SEQ(i, 2);  
    SEQ(i+1, 1);  
    ...  
    SEQ(i, N);  
    SEQ(i+1, N);  
}
```

two instructions can be pipelined, or executed on a superscalar processor

Let $SEQ(i, j)$ be the j th instruction of $SEQ(i)$.

Let each instruction chain have N instructions

Using Loop Unrolling to Exploit ILP

- Simple loop unrolling:

```
for (int i = 0; i < SIZE; i+=2) {  
    SEQ(i, 1);  
    SEQ(i+1, 1);  
    SEQ(i, 2);  
    SEQ(i+1, 1);  
    ...  
    SEQ(i, N);  
    SEQ(i+1, N);  
}
```

two instructions can be pipelined, or executed on a superscalar processor

Compiler should:

- * be in charge of unrolling factor
- * detect such loops

Loop Unrolling for Reduction Loops

- Prior approach examined loops with independent iterations and chains of dependent computations
- Now we will look at reduction loops:
 - Entire computation is dependent
 - Typically short bodies (addition, multiplication, max, min)

1	2	3	4	5	6
---	---	---	---	---	---

addition: 21

max: 6

min: 1

Loop Unrolling for Reduction Loops

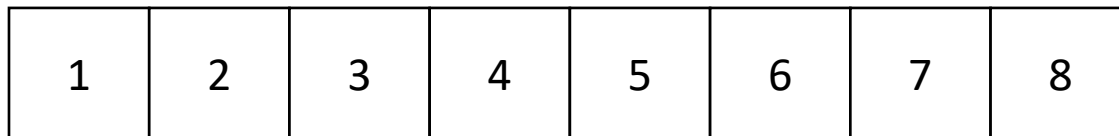
- Simple implementation:

```
for (int i = 1; i < SIZE; i++) {  
    a[0] = REDUCE(a[0], a[i]);  
}
```

If the reduction operator is associative, we can do better!

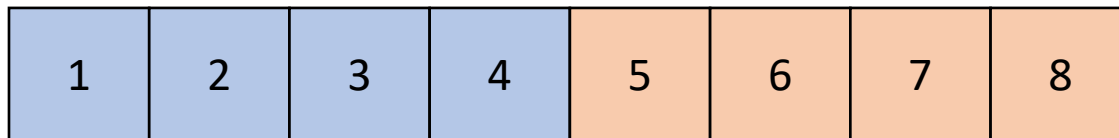
Loop Unrolling for Reduction Loops

- chunk array in equal sized partitions and do local reductions
- Consider size 2:



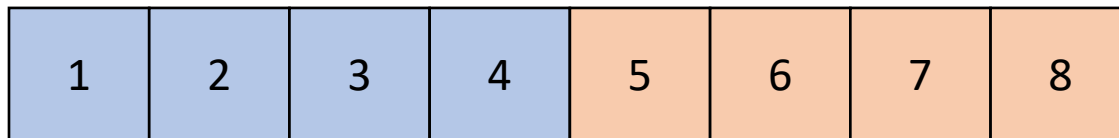
Loop Unrolling for Reduction Loops

- chunk array in equal sized partitions and do local reductions
- Consider size 2:



Loop Unrolling for Reduction Loops

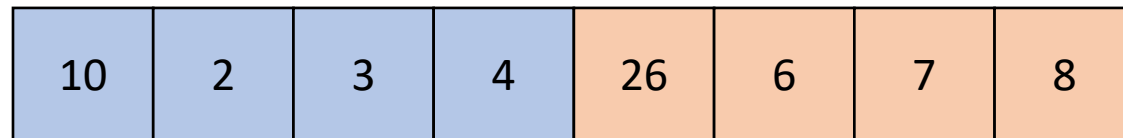
- chunk array in equal sized partitions and do local reductions
- Consider size 2:



Do addition reduction in base memory location

Loop Unrolling for Reduction Loops

- chunk array in equal sized partitions and do local reductions
- Consider size 2:



Do addition reduction in base memory location

Loop Unrolling for Reduction Loops

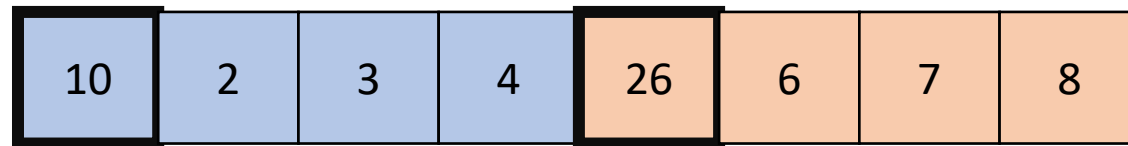
- chunk array in equal sized partitions and do local reductions
- Consider size 2:



Add together base locations

Loop Unrolling for Reduction Loops

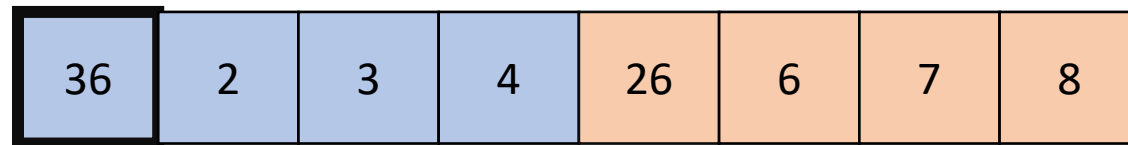
- chunk array in equal sized partitions and do local reductions
- Consider size 2:



Add together base locations

Loop Unrolling for Reduction Loops

- chunk array in equal sized partitions and do local reductions
- Consider size 2:



Add together base locations

Loop Unrolling for Reduction Loops

- Simple implementation:

```
for (int i = 1; i < SIZE/2; i++) {  
    a[0] = REDUCE(a[0], a[i]);  
    a[SIZE/2] = REDUCE(a[SIZE/2], a[(SIZE/2)+i]);  
}
```

```
a[0] = REDUCE(a[0], a[SIZE/2])
```

Loop Unrolling for Reduction Loops

- Simple implementation:

```
for (int i = 1; i < SIZE/2; i++) {  
    a[0] = REDUCE(a[0], a[i]);  
    a[SIZE/2] = REDUCE(a[SIZE/2], a[(SIZE/2)+i]);  
}
```

```
a[0] = REDUCE(a[0], a[SIZE/2])
```

Loop Unrolling for Reduction Loops

- Simple implementation:

```
for (int i = 1; i < SIZE/2; i++) {  
    a[0] = REDUCE(a[0], a[i]);  
    a[SIZE/2] = REDUCE(a[SIZE/2], a[(SIZE/2)+i]);  
}
```

```
a[0] = REDUCE(a[0], a[SIZE/2])
```

*independent
instructions
can be done
in parallel!*

Watch out!

- Our abstraction: separate dependent instructions as far as possible
- Pros:
 - Simple
- Cons:
 - Can lead to register spilling, causing expensive loads

consider `instr1` and `instr2` have a data dependence, and `instrX`'s are independent

`instr1;`

`instrX0;`

`instrX1;`

...

`instr2;`

independent instructions. If they overwrite the register storing `instr1`'s result, then it will have to be stored to memory and retrieved before `instr2`

Watch out!

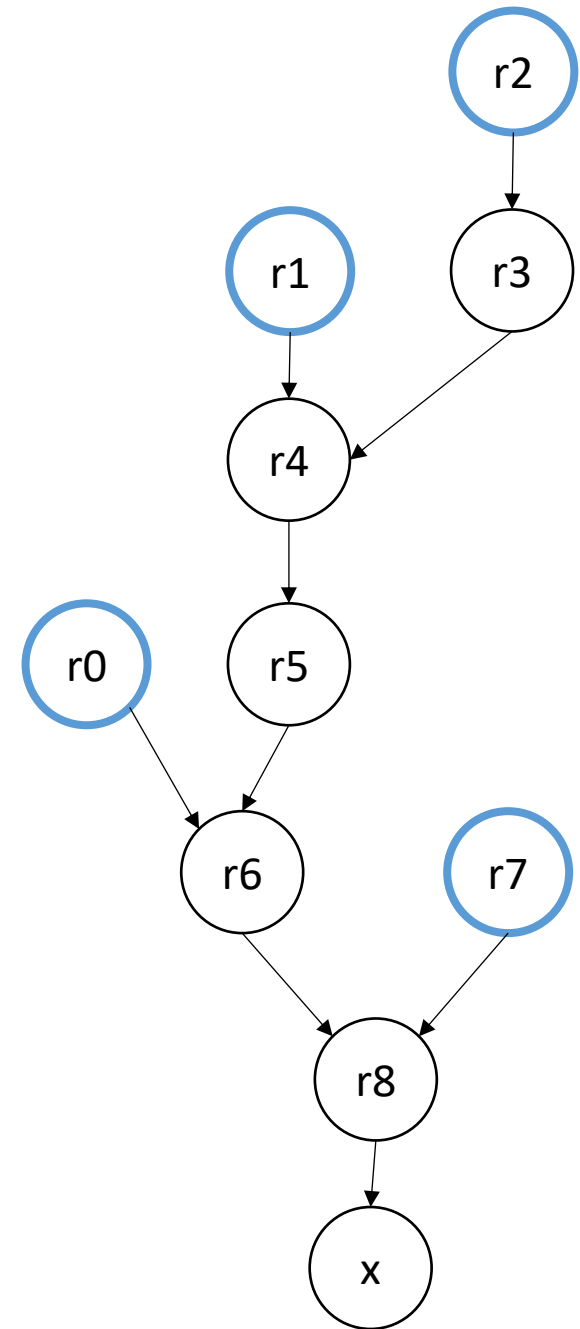
- Our abstraction: separate dependent instructions as far as possible
- Pros:
 - Simple
- Cons:
 - Can lead to register spilling, causing expensive loads

Solutions include using a **resource model** to guide the topological ordering. Highly architecture dependent. Algorithms become more expensive

Typically doesn't show up in basic block analysis. In loop unrolling, it will influence the number of unrolls you do.

Priority Topological Ordering of DDGs

```
r7 = 2 * a;  
r0 = neg(b);  
r1 = b * b;  
r2 = 4 * a;  
r3 = r2 * c;  
r4 = r1 - r3;  
r5 = sqrt(r4);  
r6 = r0 - r5;  
r8 = r6 / r7;  
x = r8;
```



Discussion

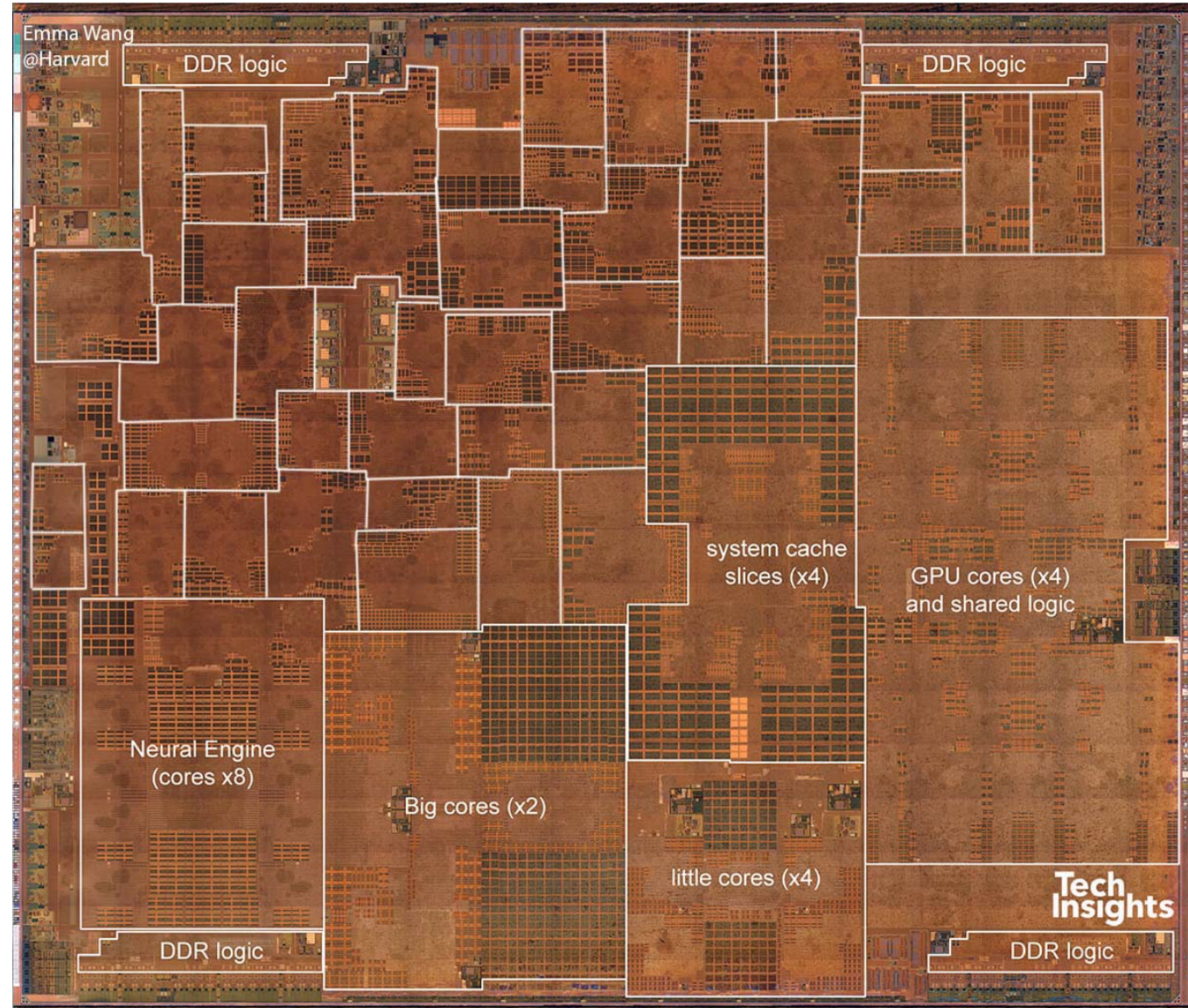
- Where is parallelism most commonly found?
 - Non-numeric applications are thought to have very little. lots of:
 - I/O (file, network, user),
 - events,
 - *source needed*
 - numeric applications have more:
 - media processing (image, video, sound)
 - machine-learning (esp. inference)
- More and more, numeric applications are moving to accelerators

Modern SoC

- From David Brooks lab at Harvard:

<http://vlsiarch.eecs.harvard.edu/research/accelerators/die-photo-analysis/>

- Compilers will need to be able to map software efficiently to a range of different accelerators



Current tensions

- Simple cores with accelerators/GPUs?
 - Less need for pipelines, OoO, and superscalar
 - Hard to port legacy code
- Complicated cores
 - area/power hungry
 - great for legacy code
- Where do compilers fit in?

Next week

- SPMD parallelism for loops