NORTHWEST INSTITUTE *for* ADVANCED COMPUTING

Pacific Northwest
NATIONAL LABORATORY
*Proudly Operated by Battelle
for the U.S. Department of Energy*

**W**
UNIVERSITY *of*
WASHINGTON

# AMATH 483/583
# High Performance Scientific Computing
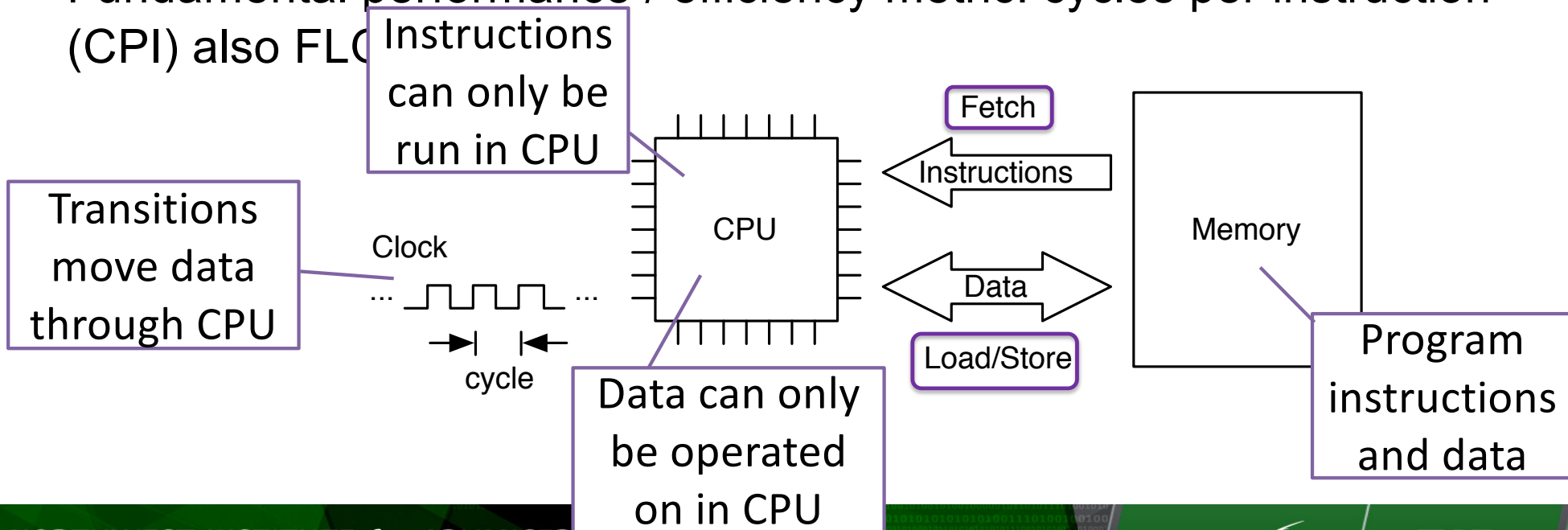
# Lecture 6:
# High Performance in Hierarchical Memory

Andrew Lumsdaine

Northwest Institute for Advanced Computing

Pacific Northwest National Laboratory

University of Washington

Seattle, WA

# Overview

- "PDP-11" machine model

- Pipelining, pipeline stalls

- Hierarchical memory

- Timing and benchmarking

- Compiler optimizations

- Tiling

- Blocking

NORTHWEST INSTITUTE for ADVANCED COMPUTING

2

AMATH 483/583 High-Performance Scientific Computing Spring 2019
University of Washington by Andrew Lumsdaine

Pacific Northwest
NATIONAL LABORATORY

Proudly Operated by Battelle
for the U.S. Department of Energy
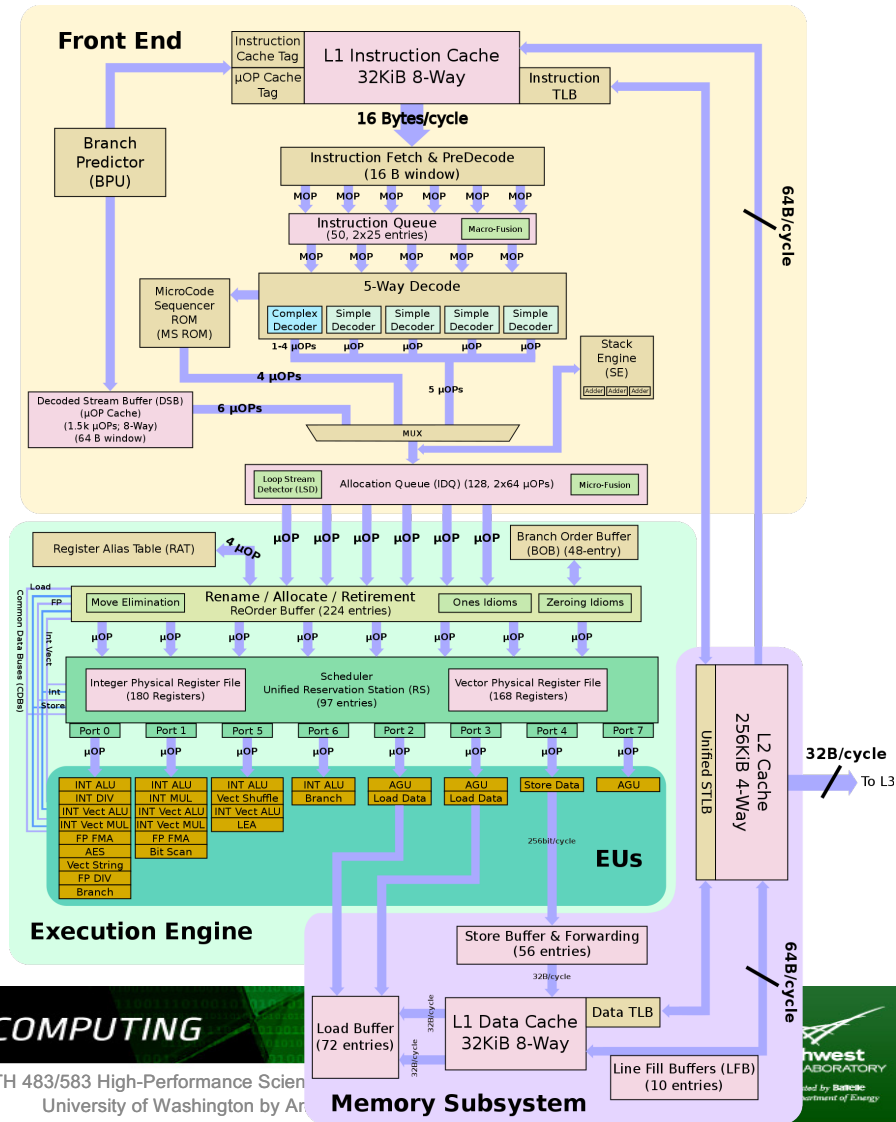
UNIVERSITY of
WASHINGTON

# Microprocessors

- Basic operation: read and execute program instructions stored in memory

- Fundamental performance / efficiency metric: cycles per instruction (CPI) also FL...

**Instructions can only be run in CPU**

**Transitions move data through CPU**

Clock

... ⊓⊔⊓⊔ ...

cycle

CPU

**Data can only be operated on in CPU**

Fetch

Instructions

Data

Load/Store

Memory

**Program instructions and data**

NORTHWEST INSTITUTE for ADVANCED COMPUTING

Pacific Northwest
NATIONAL LABORATORY
Proudly Operated by Battelle
for the U.S. Department of Energy

UNIVERSITY of WASHINGTON

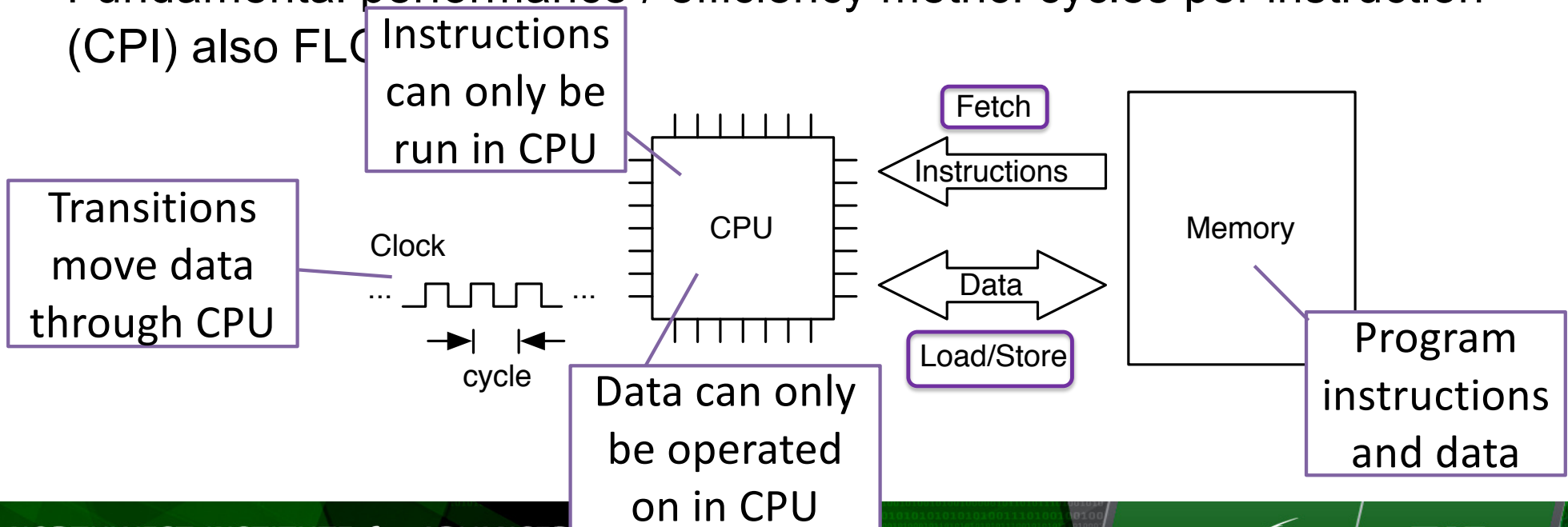# Performance-Oriented Architecture Features

- Execution Pipeline
  - Stages of functionality to process issued instructions
  - Hazards are conflicts with continued execution
  - Forwarding supports closely associated operations exhibiting precedence constraints
- Out of Order Execution
  - Uses reservation stations
  - Hides some core latencies and provide fine grain asynchronous operation supporting concurrency
- Branch Prediction
  - Permits computation to proceed at a conditional branch point prior to resolving predicate value
  - Overlaps follow-on computation with predicate resolution
  - Requires roll-back or equivalent to correct false guesses
  - Sometimes follows both paths, and several deep

NORTHWEST INSTITUTE for ADVANCED COMPUTING

4

AMATH 483/583 High-Performance Scientific Computing Spring 2019
University of Washington by Andrew Lumsdaine

Pacific Northwest
NATIONAL LABORATORY

Proudly Operated by Battelle
for the U.S. Department of Energy

UNIVERSITY of
WASHINGTON

# Skylake



**Front End**

| Instruction Cache Tag | L1 Instruction Cache 32KiB 8-Way | |
| µOP Cache Tag | | Instruction TLB |

**16 Bytes/cycle**

Branch Predictor (BPU)

Instruction Fetch & PreDecode (16 B window)

MOP MOP MOP MOP MOP MOP

Instruction Queue (50, 2x25 entries) — Macro-Fusion

MOP MOP MOP MOP MOP

MicroCode Sequencer ROM (MS ROM)

5-Way Decode

| Complex Decoder | Simple Decoder | Simple Decoder | Simple Decoder | Simple Decoder |

1-4 µOPs   µOP   µOP   µOP   µOP

Stack Engine (SE)
Adder Adder Adder

**4 µOPs**        **5 µOPs**

Decoded Stream Buffer (DSB) (µOP Cache) (1.5k µOPs; 8-Way) (64 B window)

**6 µOPs**

MUX

Loop Stream Detector (LSD) — Allocation Queue (IDQ) (128, 2x64 µOPs) — Micro-Fusion

µOP µOP µOP µOP µOP µOP

**64B/cycle**

Register Alias Table (RAT)

**4 µOP**

Branch Order Buffer (BOB) (48-entry)

Load
FP
Int Vect.
Int
Store

Move Elimination

Rename / Allocate / Retirement
ReOrder Buffer (224 entries)

Ones Idioms   Zeroing Idioms

µOP µOP µOP µOP µOP µOP µOP µOP

Common Data Buses (CDBs)

Integer Physical Register File (180 Registers)

Scheduler
Unified Reservation Station (RS) (97 entries)

Vector Physical Register File (168 Registers)

| Port 0 | Port 1 | Port 5 | Port 6 | Port 2 | Port 3 | Port 4 | Port 7 |
| µOP | µOP | µOP | µOP | µOP | µOP | µOP | µOP |

| INT ALU | INT ALU | INT ALU | INT ALU | AGU | AGU | Store Data | AGU |
| INT DIV | INT MUL | Vect Shuffle | Branch | Load Data | Load Data | | |
| INT Vect ALU | INT Vect ALU | INT Vect ALU | | | | | |
| INT Vect MUL | INT Vect MUL | LEA | | | | | |
| FP FMA | FP FMA | | | | | | |
| AES | Bit Scan | | | | | | |
| Vect String | | | | | | | |
| FP DIV | | | | | | | |
| Branch | | | | | | | |

256bit/cycle

**EUs**

**Execution Engine**

Unified STLB

L2 Cache 256KiB 4-Way

**32B/cycle** → To L3

Store Buffer & Forwarding (56 entries)

32B/cycle

**64B/cycle**

Load Buffer (72 entries)

L1 Data Cache 32KiB 8-Way

Data TLB

Line Fill Buffers (LFB) (10 entries)

**Memory Subsystem**

NORTHWEST INSTITUTE for ADVANCED COMPUTING

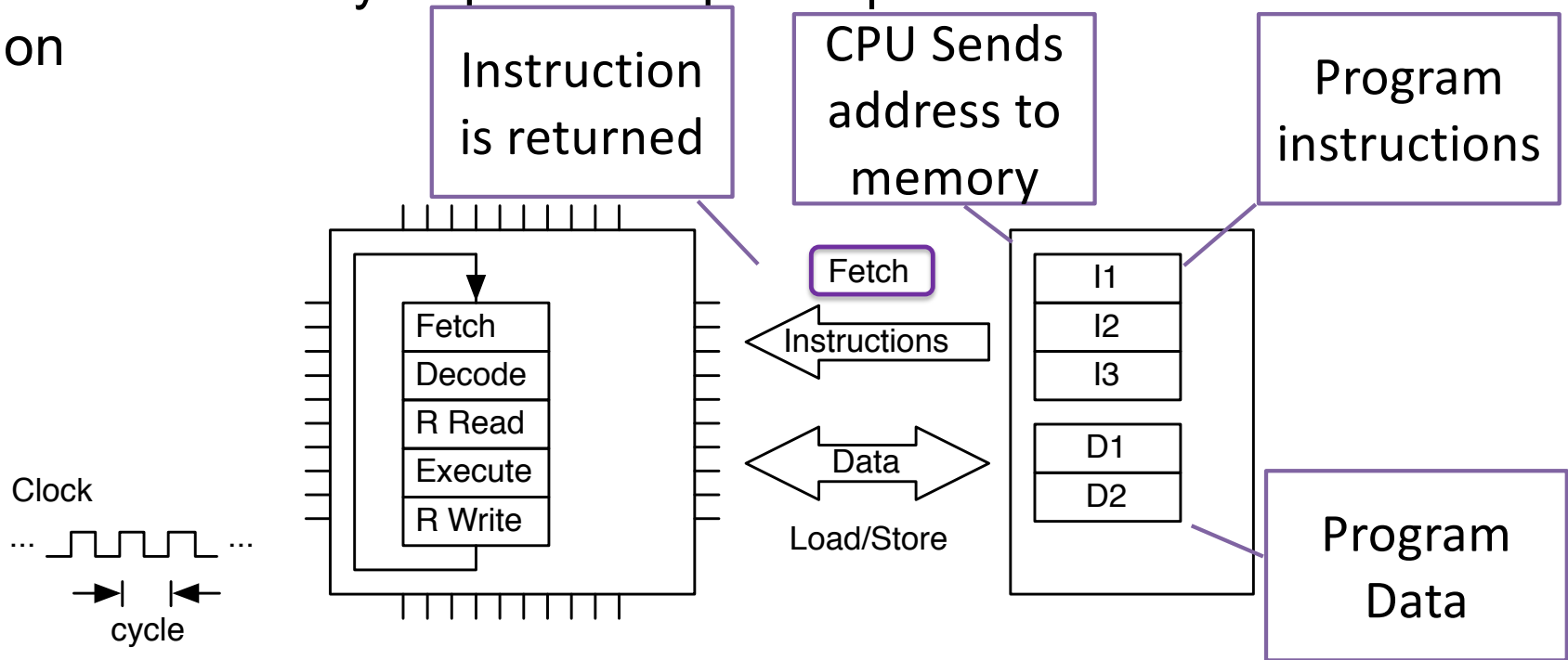...hwest ...ABORATORY

UNIVERSITY of WASHINGTON

# Microprocessors

- Basic operation: read and execute program instructions stored in memory
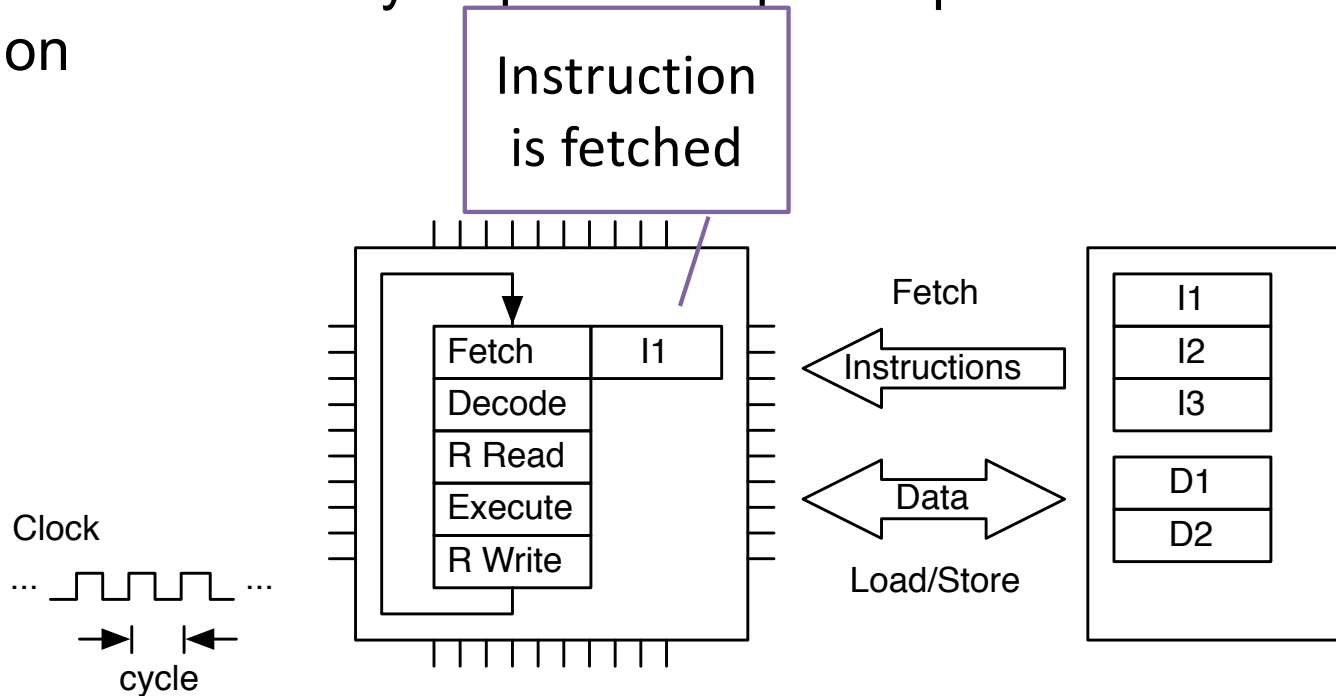- Fundamental performance / efficiency metric: cycles per instruction (CPI) also FLO...

Instructions can only be run in CPU

Transitions move data through CPU

Clock

... cycle

CPU

Fetch

Instructions

Data

Load/Store

Memory

Data can only be operated on in CPU

Program instructions and data

NORTHWEST INSTITUTE for ADVANCED COMPUTING

Pacific Northwest
NATIONAL LABORATORY
Proudly Operated by Battelle
for the U.S. Department of Energy

UNIVERSITY of WASHINGTON
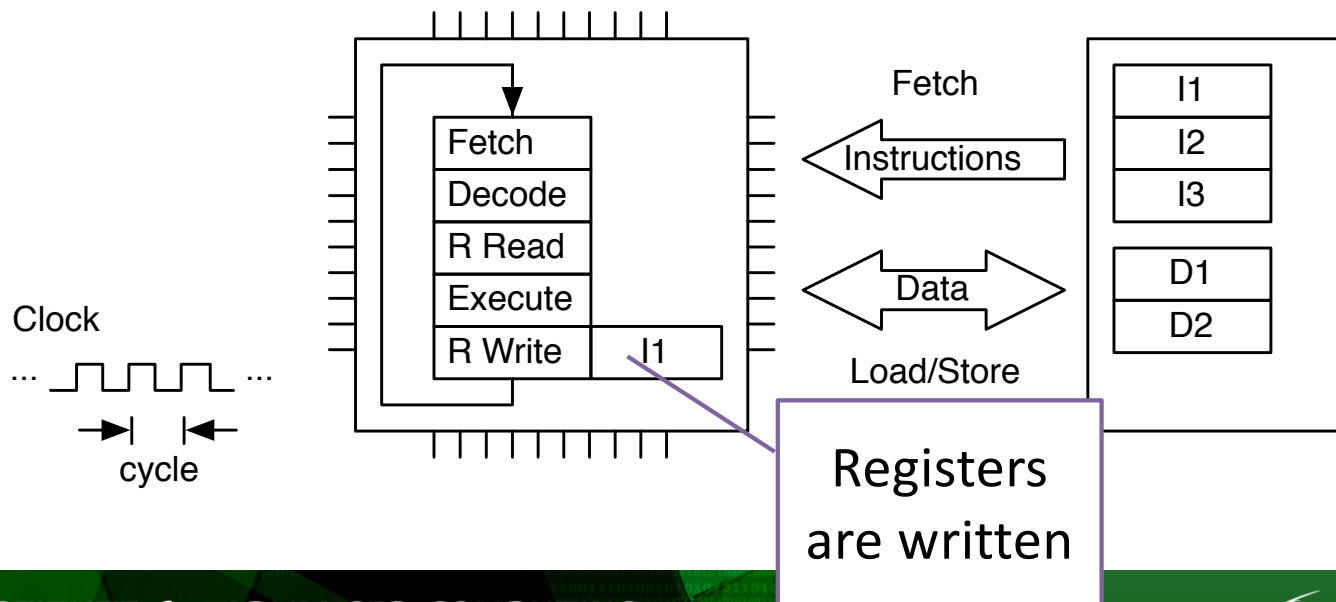
# Processor Core Instruction Handling

- A single instruction may require multiple steps from fetch to completion

Instruction is returned

CPU Sends address to memory

Program instructions

Fetch

Fetch
Decode
R Read
Execute
R Write

Instructions

Data

Load/Store

Clock

··· ··· 

cycle

I1
I2
I3

D1
D2

Program Data

NORTHWEST INSTITUTE for ADVANCED COMPUTING

Pacific Northwest
NATIONAL LABORATORY
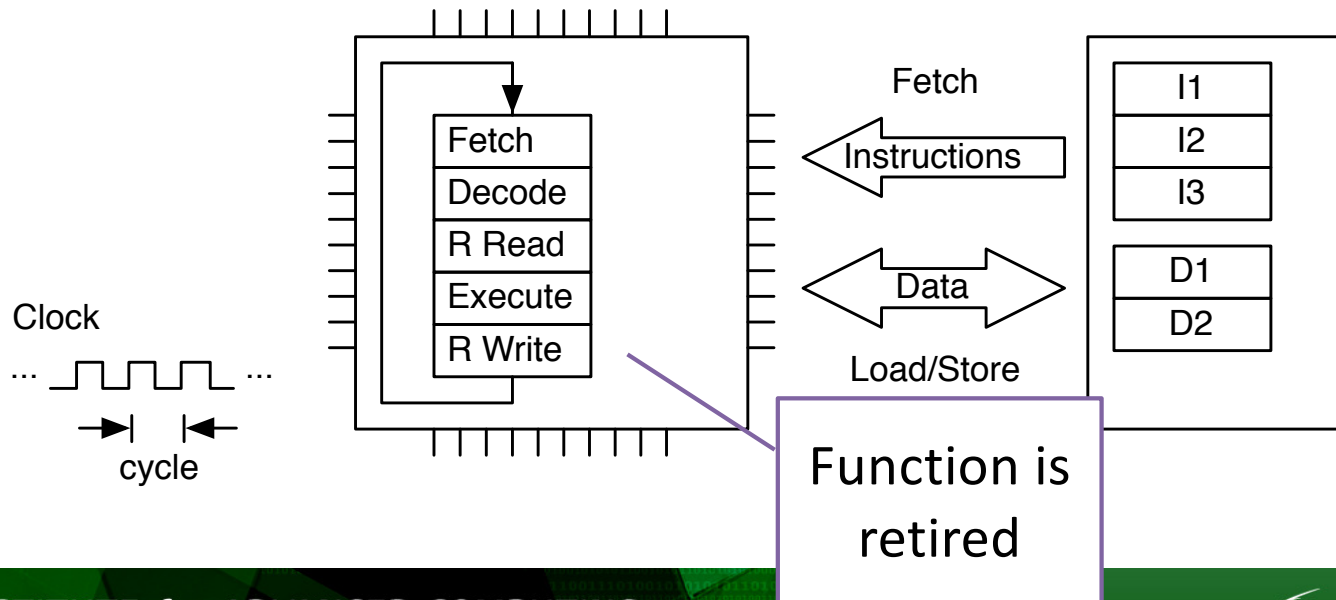Proudly Operated by Battelle
for the U.S. Department of Energy

UNIVERSITY of WASHINGTON

# Processor Core Instruction Handling

- A single instruction may require multiple steps from fetch to completion

Instruction is fetched

Fetch

Instructions

Data

Load/Store

| Fetch | I1 |
|-------|----|
| Decode | |
| R Read | |
| Execute | |
| R Write | |

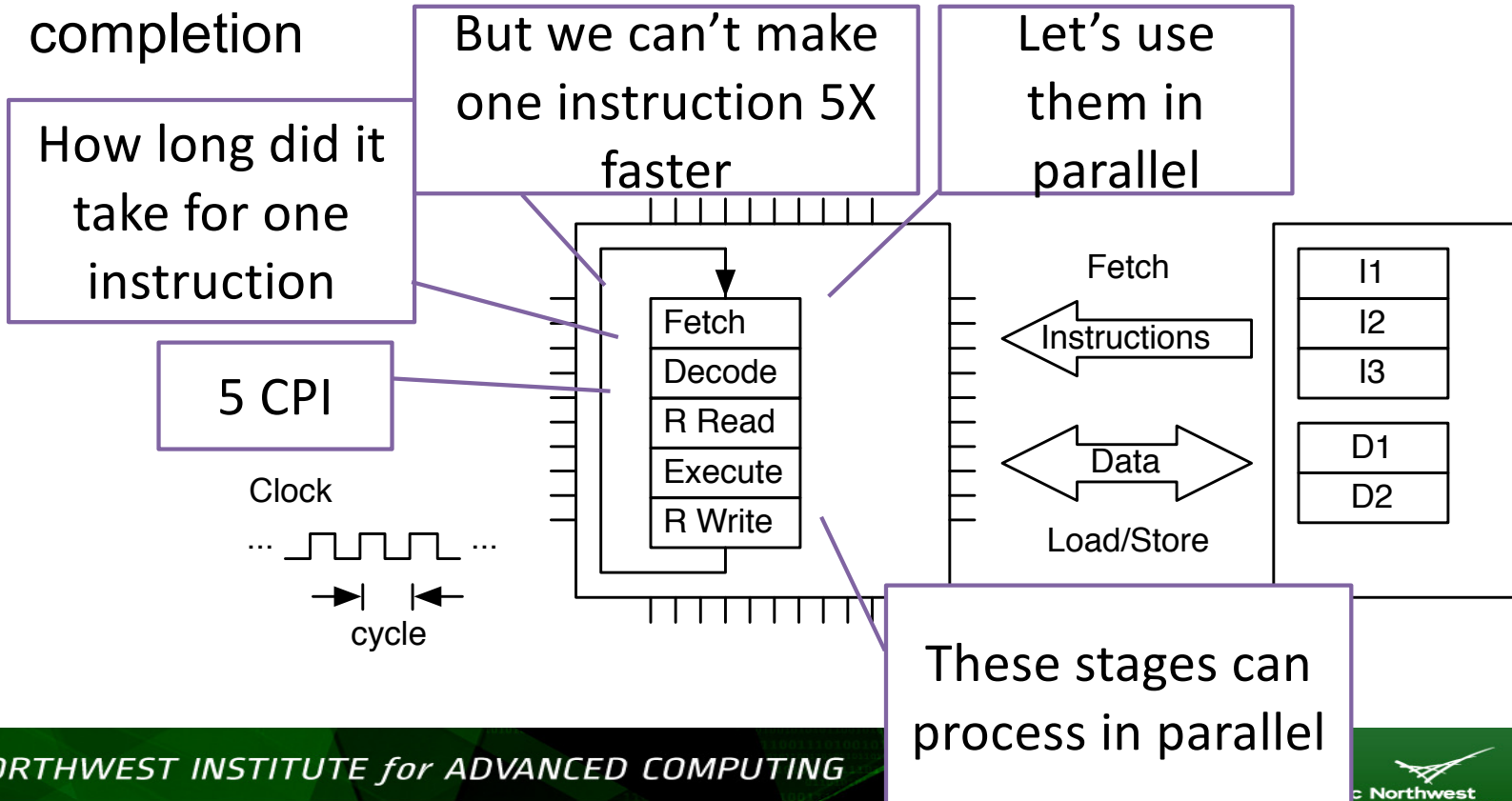| I1 |
|----|
| I2 |
| I3 |

| D1 |
|----|
| D2 |

Clock

... cycle

# Processor Core Instruction Handling

- A single instruction may require multiple steps from fetch to completion

# Processor Core Instruction Handling

- A single instruction may require multiple steps from fetch to completion

# Processor Core Instruction Handling

- A single instruction may require multiple steps from fetch to completion



Instruction is executed

NORTHWEST INSTITUTE for ADVANCED COMPUTING

# Processor Core Instruction Handling

- A single instruction may require multiple steps from fetch to completion



Clock

... cycle

Fetch | Decode | R Read | Execute | R Write | I1

Fetch
Instructions

Data

Load/Store

I1
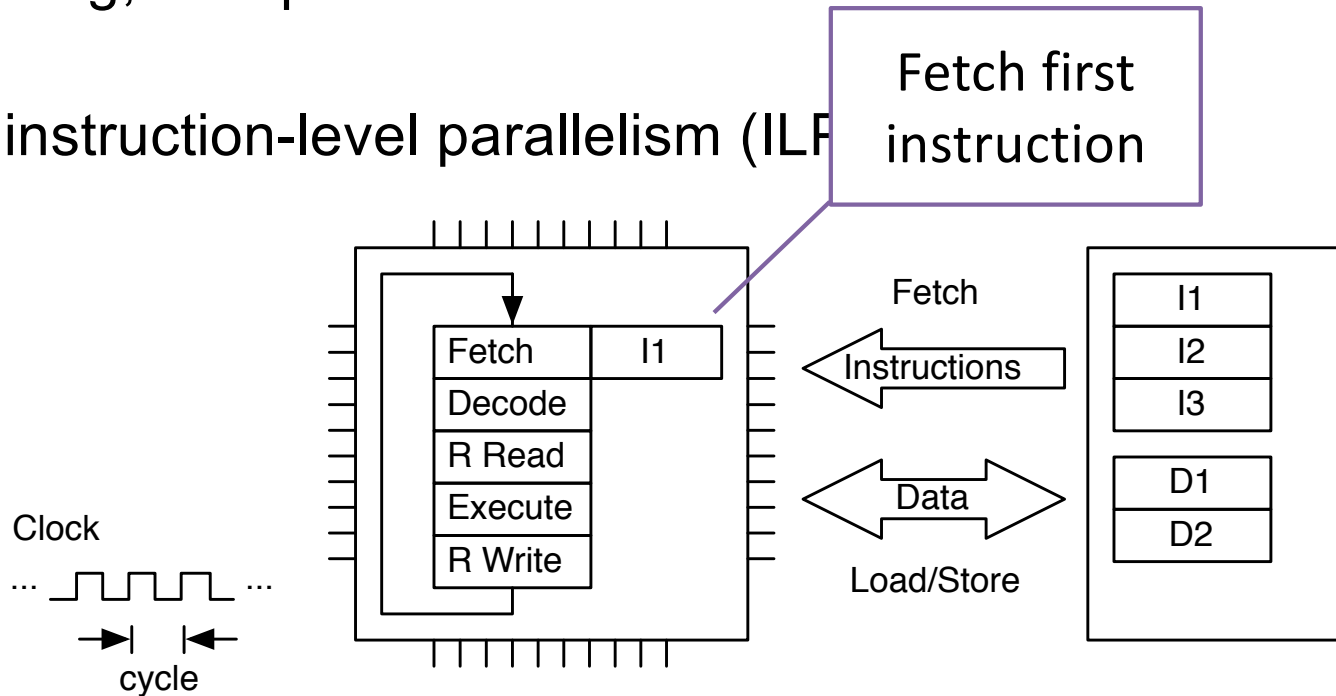I2
I3

D1
D2

Registers are written

# Processor Core Instruction Handling

- A single instruction may require multiple steps from fetch to completion

# Processor Core Instruction Handling

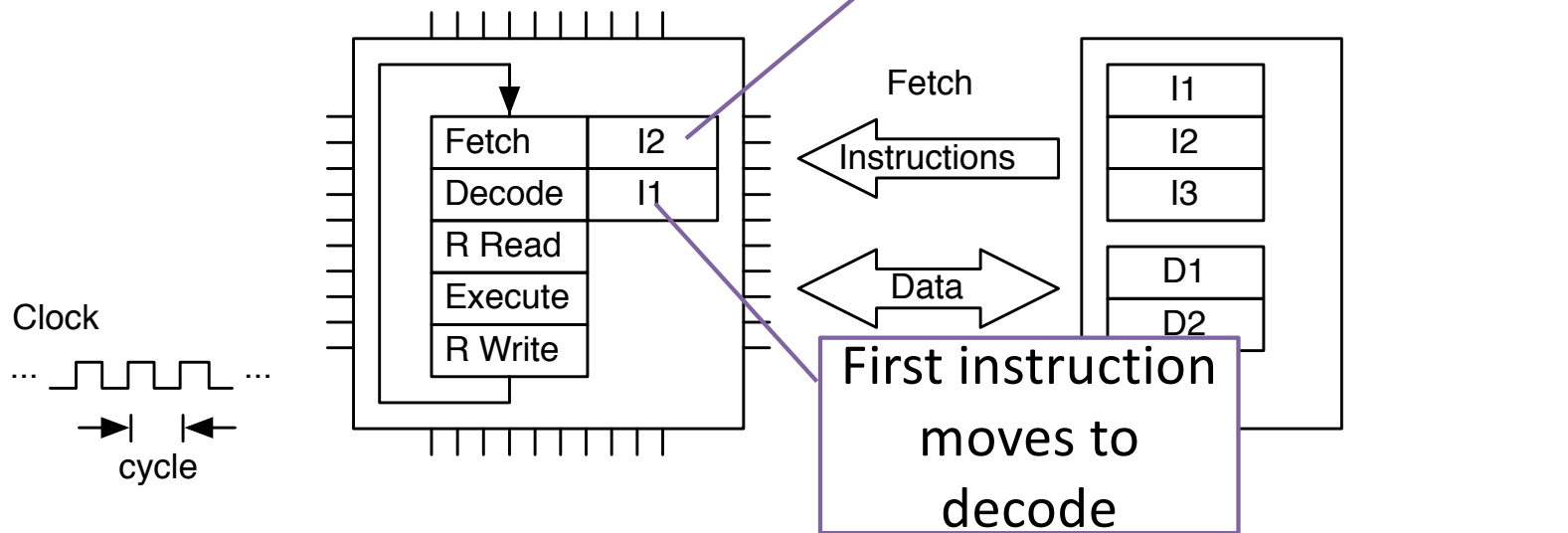- A single instruction may require multiple steps from fetch to completion

But we can't make one instruction 5X faster

Let's use them in parallel

How long did it take for one instruction

5 CPI

Clock

... ⊓_⊓_⊓ ...

|←→| |←→| cycle

Fetch
Decode
R Read
Execute
R Write

Fetch
Instructions

Data

Load/Store

I1
I2
I3

D1
D2

These stages can process in parallel

AMATH 483/583 High-Performance Scientific Computing Spring 2019
University of Washington by Andrew Lumsdaine

# Processor Core Instruction Handling

- By pipelining, multiple instructions can be executed at each clock cycle
- Form of instruction-level parallelism (ILP)

Fetch first instruction

# Processor Core Instruction Handling

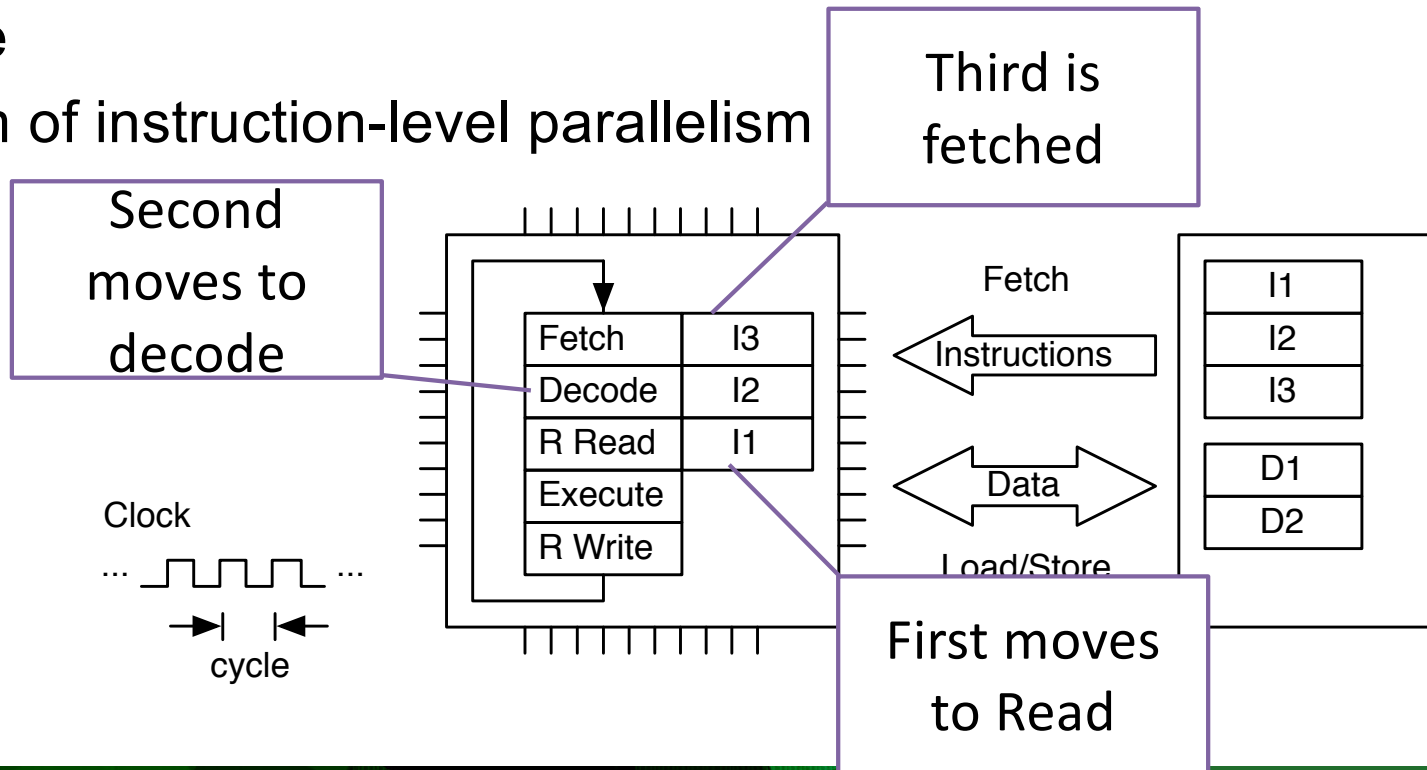- By pipelining, multiple instructions can be executed at each clock cycle
- Form of instruction-level parallelism (ILP)

When first instruction is in decode, fetch second

First instruction moves to decode

| Fetch | I2 |
| Decode | I1 |
| R Read | |
| Execute | |
| R Write | |

Clock

… ⊓⊔⊓⊔ …

cycle

Fetch

Instructions

Data

| I1 |
| I2 |
| I3 |

| D1 |
| D2 |

AMATH 483/583 High-Performance Scientific Computing Spring 2019
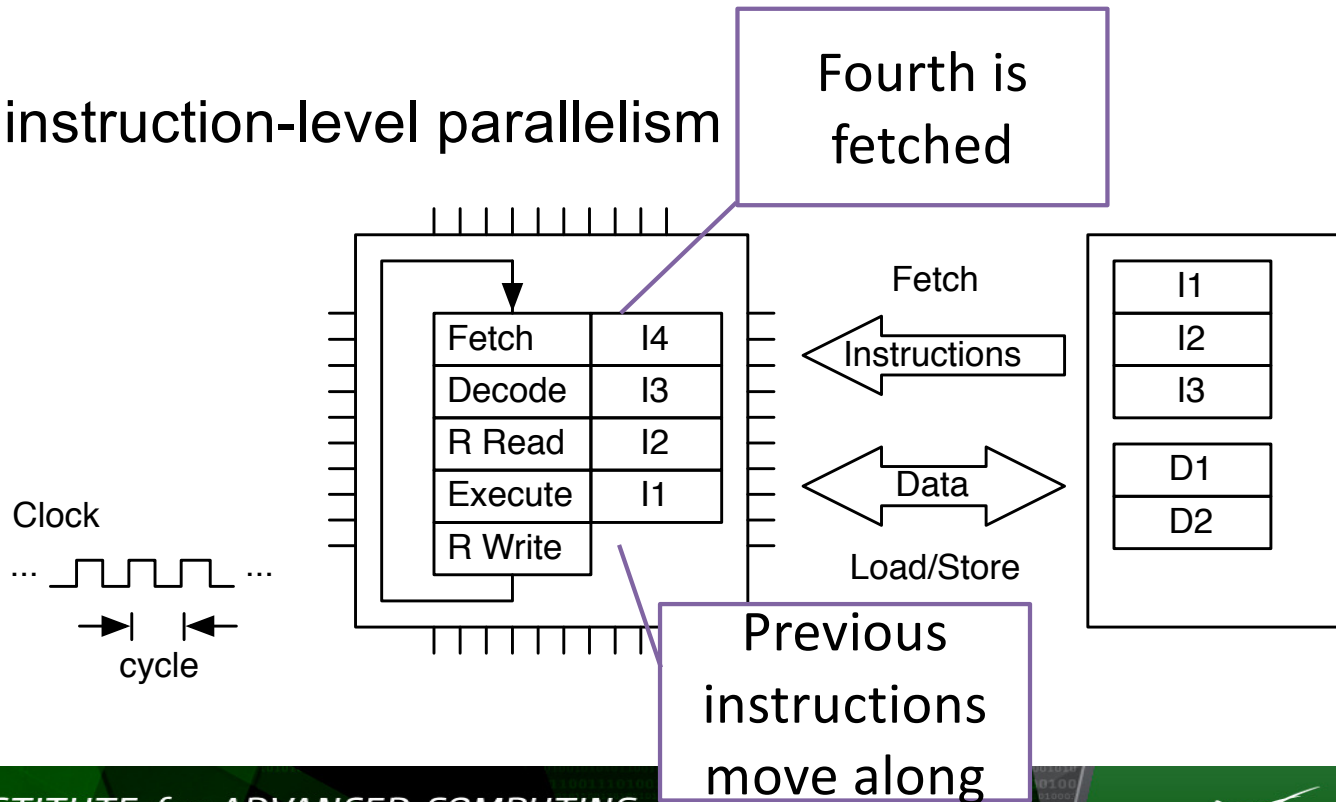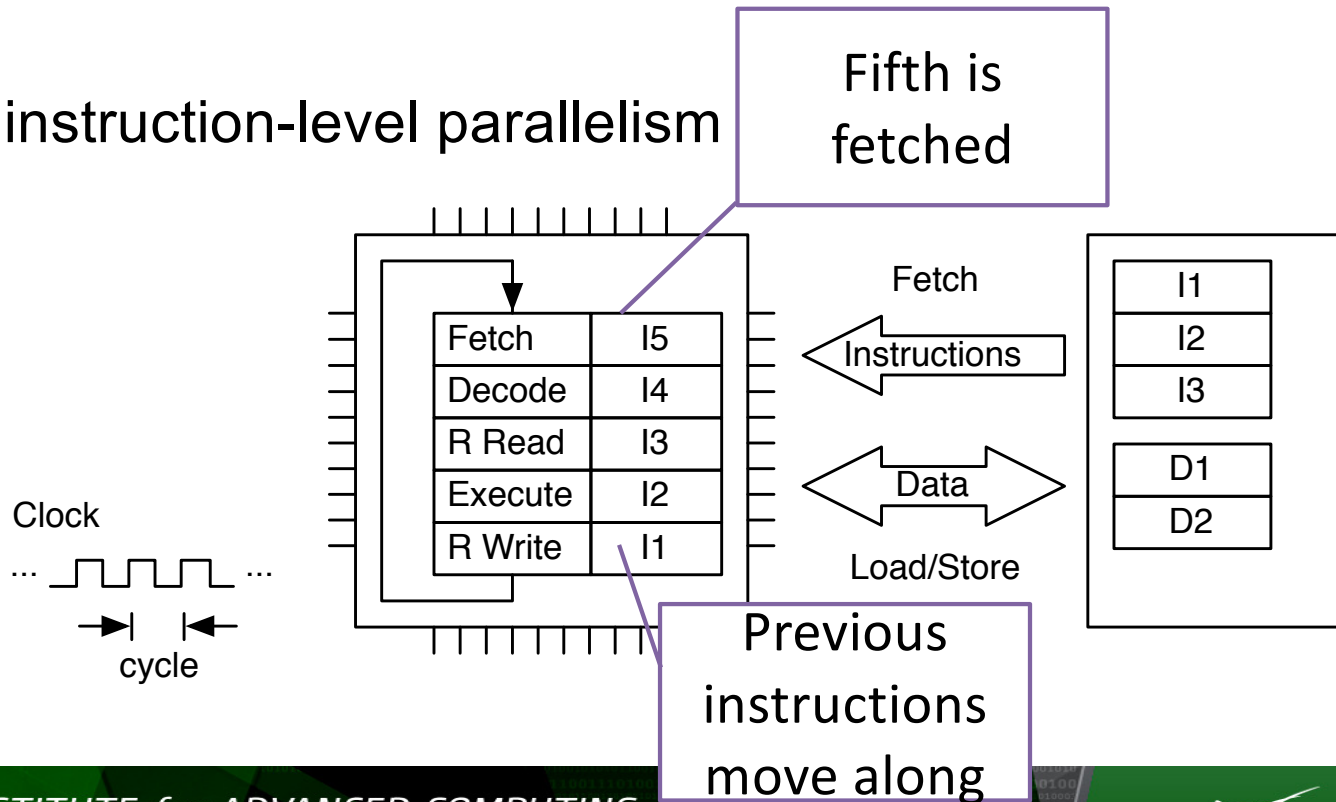University of Washington by Andrew Lumsdaine

# Processor Core Instruction Handling

- By pipelining, multiple instructions can be executed at each clock cycle
- Form of instruction-level parallelism



Third is fetched

Second moves to decode

First moves to Read

| Fetch | I3 |
| Decode | I2 |
| R Read | I1 |
| Execute | |
| R Write | |

Clock

... cycle

Fetch

Instructions

Data

Load/Store

| I1 |
| I2 |
| I3 |

| D1 |
| D2 |

NORTHWEST INSTITUTE for ADVANCED COMPUTING

Pacific Northwest
NATIONAL LABORATORY
Proudly Operated by Battelle
for the U.S. Department of Energy
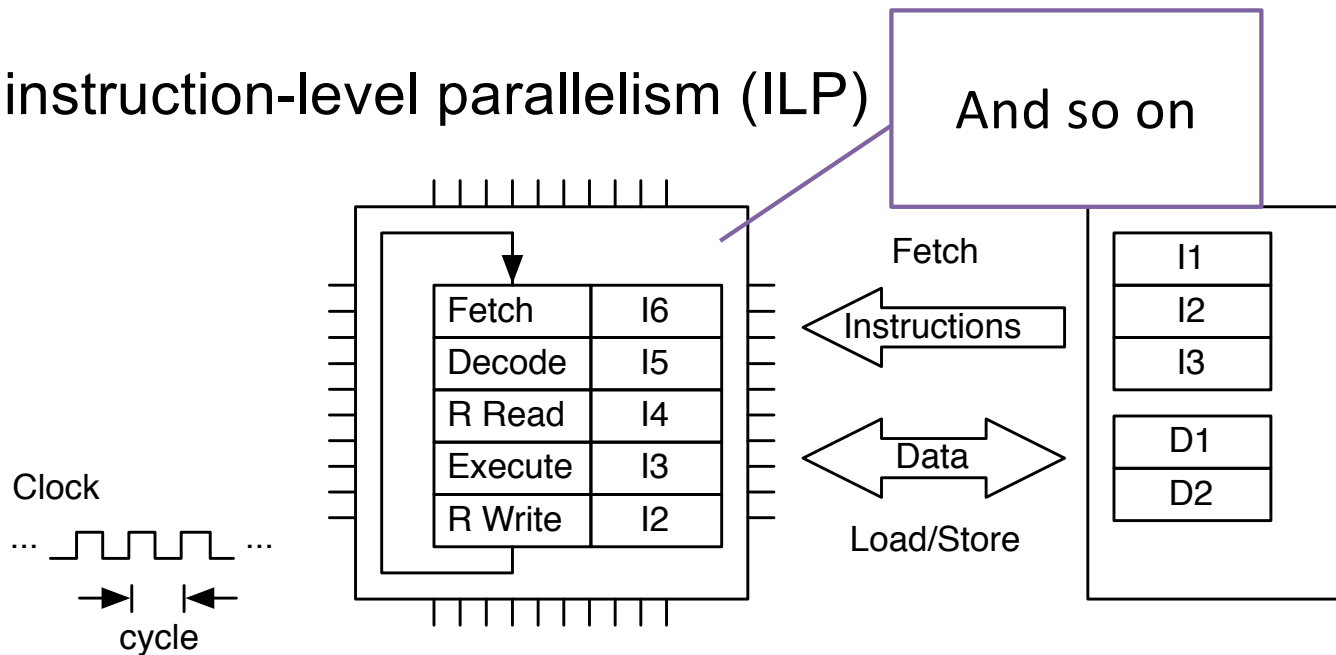
UNIVERSITY of WASHINGTON

# Processor Core Instruction Handling

- By pipelining, multiple instructions can be executed at each clock cycle
- Form of instruction-level parallelism

Fourth is fetched

Previous instructions move along

| Fetch | I4 |
|---|---|
| Decode | I3 |
| R Read | I2 |
| Execute | I1 |
| R Write | |

Clock

... cycle

Fetch

Instructions

Data

Load/Store

| I1 |
|---|
| I2 |
| I3 |

| D1 |
|---|
| D2 |

NORTHWEST INSTITUTE for ADVANCED COMPUTING

Pacific Northwest
NATIONAL LABORATORY
Proudly Operated by Battelle
for the U.S. Department of Energy
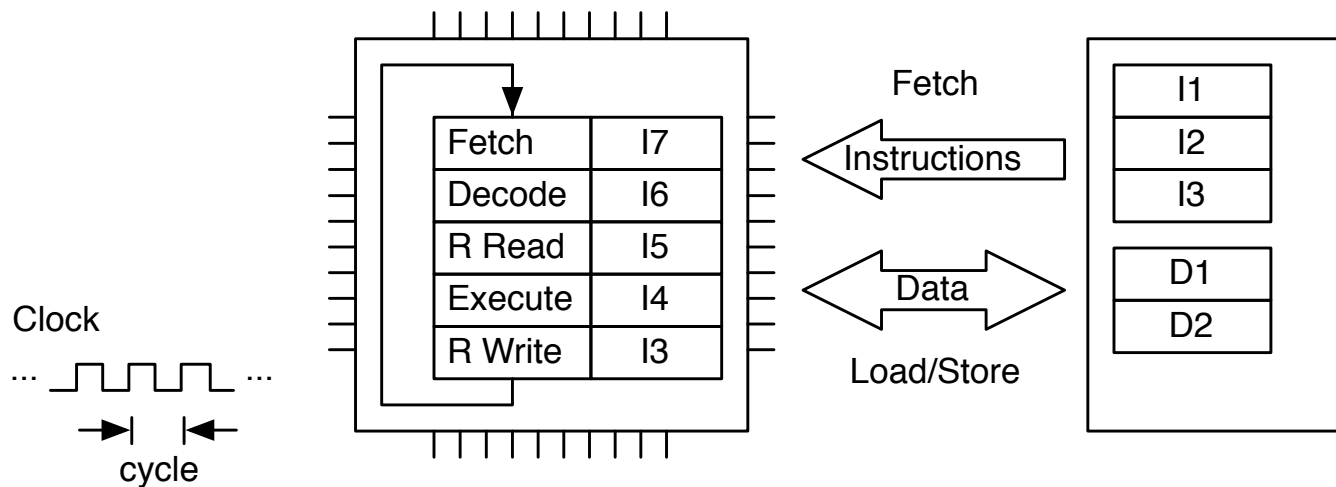
UNIVERSITY of
WASHINGTON

# Processor Core Instruction Handling

- By pipelining, multiple instructions can be executed at each clock cycle
- Form of instruction-level parallelism

Fifth is fetched

Previous instructions move along

Clock

... ⊓⊔⊓⊔⊓ ...

cycle

| Fetch | I5 |
| Decode | I4 |
| R Read | I3 |
| Execute | I2 |
| R Write | I1 |

Fetch
Instructions

Data

Load/Store

| I1 |
| I2 |
| I3 |

| D1 |
| D2 |

NORTHWEST INSTITUTE for ADVANCED COMPUTING

Pacific Northwest
NATIONAL LABORATORY
Proudly Operated by Battelle
for the U.S. Department of Energy
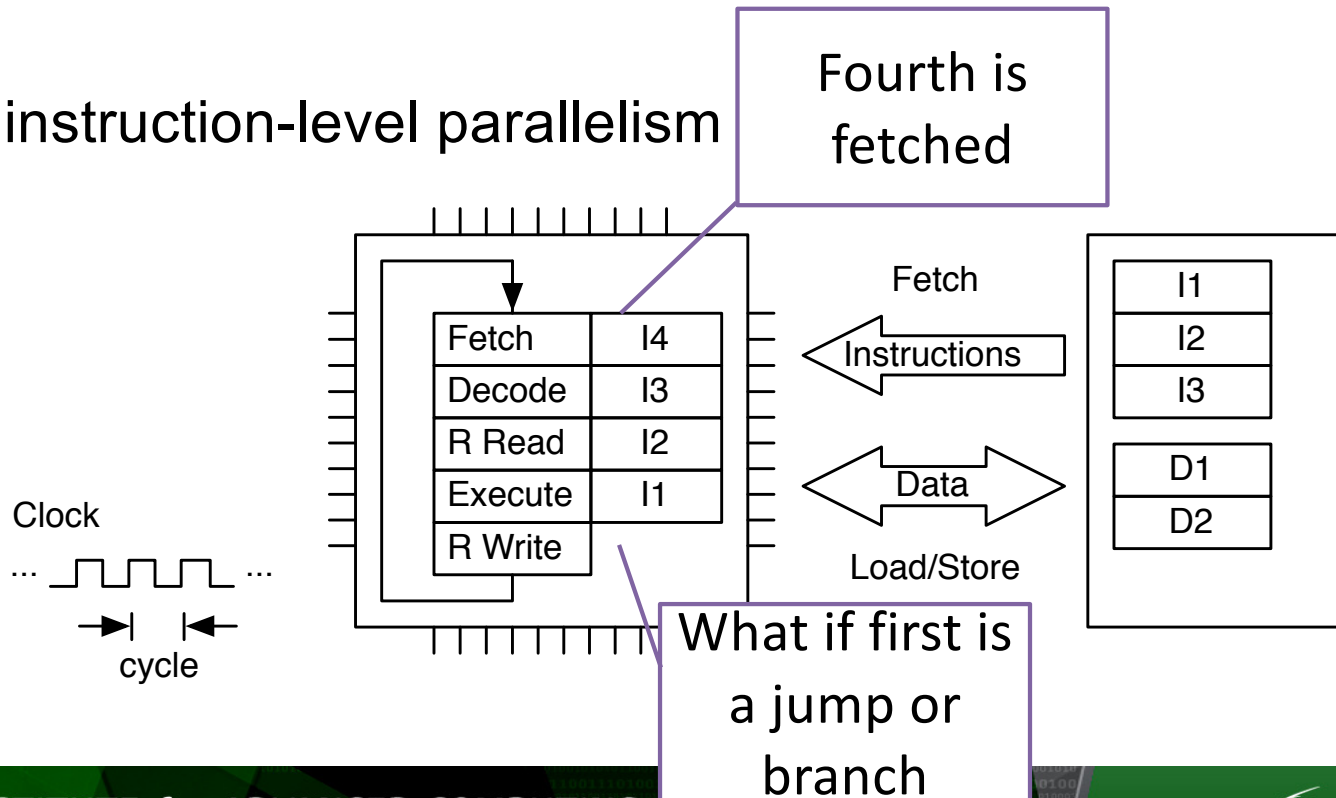
UNIVERSITY of
WASHINGTON

# Processor Core Instruction Handling

- By pipelining, multiple instructions can be executed at each clock cycle
- Form of instruction-level parallelism (ILP)

And so on

NORTHWEST INSTITUTE for ADVANCED COMPUTING

Pacific Northwest
NATIONAL LABORATORY
Proudly Operated by Battelle
for the U.S. Department of Energy

UNIVERSITY of
WASHINGTON

# Processor Core Instruction Handling

- By pipelining, multiple instructions can be executed at each clock cycle
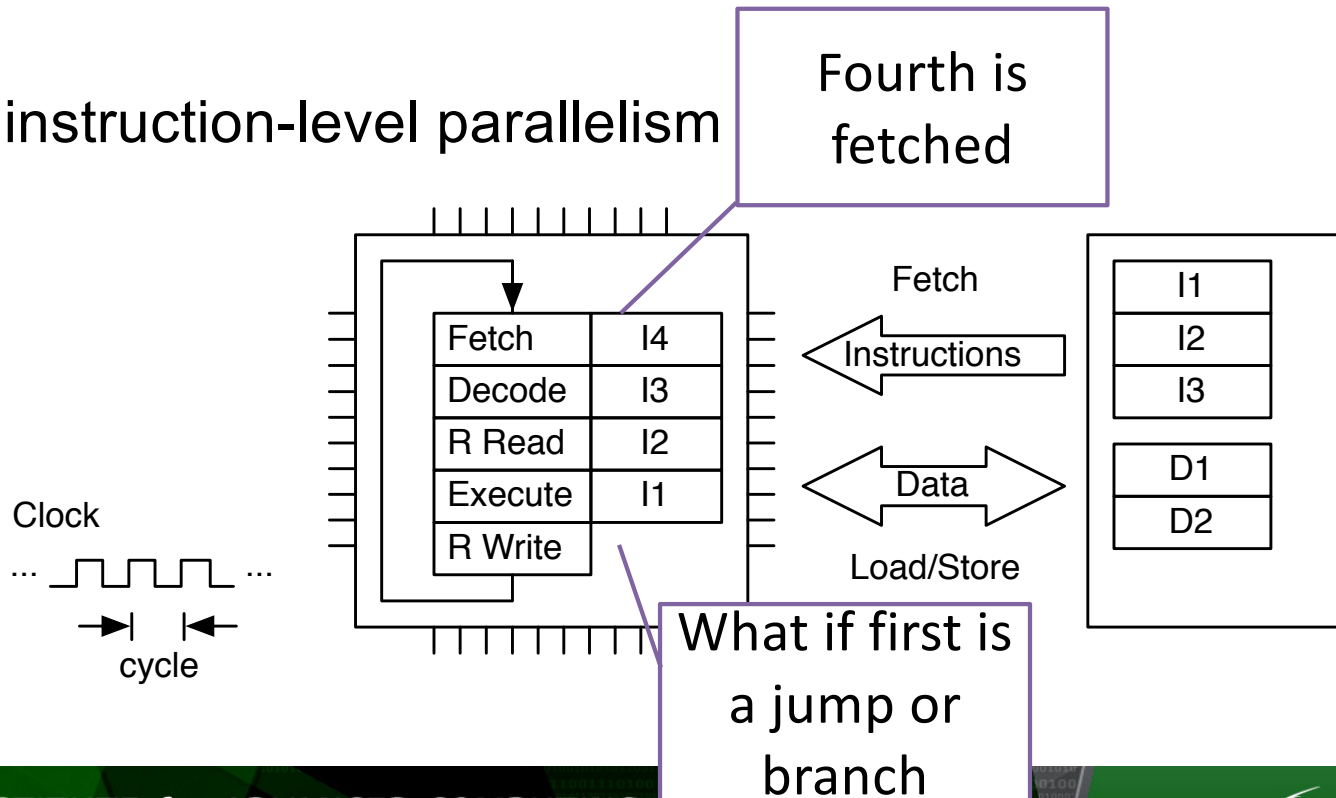- Form of instruction-level parallelism (ILP)

# Pipeline Stall

- By pipelining, multiple instructions can be executed at each clock cycle
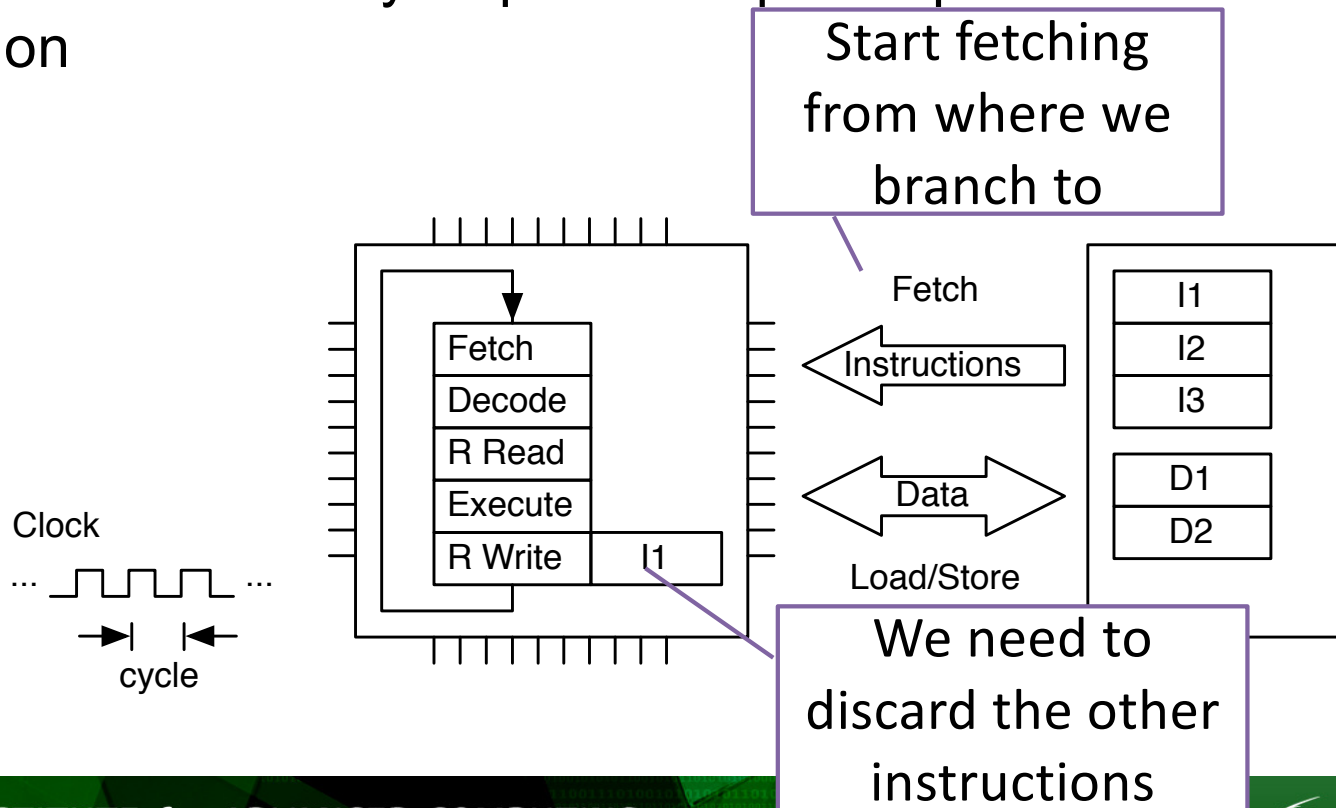- Form of instruction-level parallelism

Fourth is fetched

| Fetch | I4 |
| Decode | I3 |
| R Read | I2 |
| Execute | I1 |
| R Write | |

Clock

... cycle ...

Fetch

Instructions

Data

Load/Store

| I1 |
| I2 |
| I3 |

| D1 |
| D2 |

What if first is a jump or branch

NORTHWEST INSTITUTE for ADVANCED COMPUTING

Pacific Northwest
NATIONAL LABORATORY
Proudly Operated by Battelle
for the U.S. Department of Energy

UNIVERSITY of
WASHINGTON

# Pipeline Stall

- By pipelining, multiple instructions can be executed at each clock cycle
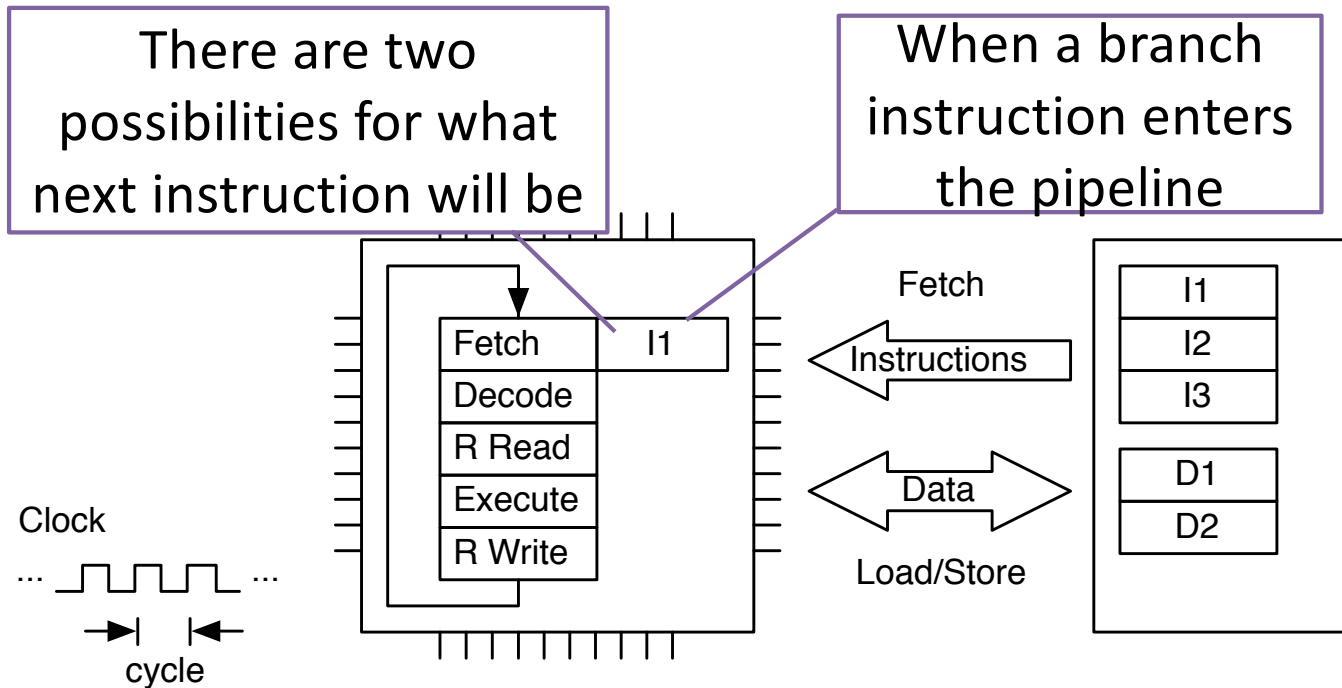- Form of instruction-level parallelism

Fourth is fetched

| Fetch | I4 |
|---|---|
| Decode | I3 |
| R Read | I2 |
| Execute | I1 |
| R Write | |

Fetch
Instructions

Data

Load/Store

| I1 |
|---|
| I2 |
| I3 |

| D1 |
|---|
| D2 |

Clock

... ⊓⊔⊓⊔⊓⊔ ...

cycle

What if first is a jump or branch

NORTHWEST INSTITUTE *for* ADVANCED COMPUTING

Pacific Northwest
NATIONAL LABORATORY

*Proudly Operated by Battelle
for the U.S. Department of Energy*

W
UNIVERSITY of
WASHINGTON

# Pipeline Stall

- A single instruction may require multiple steps from fetch to completion

Start fetching from where we branch to

Fetch

Instructions

Data

Load/Store

| Fetch |
| Decode |
| R Read |
| Execute |
| R Write | I1 |

Clock

... ⊓⊔⊓⊔⊓ ...

cycle

| I1 |
| I2 |
| I3 |

| D1 |
| D2 |

We need to discard the other instructions

NORTHWEST INSTITUTE for ADVANCED COMPUTING

Pacific Northwest
NATIONAL LABORATORY
Proudly Operated by Battelle
for the U.S. Department of Energy

UNIVERSITY of WASHINGTON

# Branch Prediction

- Load the instructions we think will be branched to

There are two possibilities for what next instruction will be

When a branch instruction enters the pipeline

Fetch

| I1 |
| I2 |
| I3 |

| Fetch | I1 |
| Decode | |
| R Read | |
| Execute | |
| R Write | |

Instructions

Data

Load/Store

| D1 |
| D2 |

Clock

... cycle

# Branch Prediction

- Load the instructions we think will be branched to

This is the instruction the CPU predicts will branch to

First instruction moves to decode

Fetch

| Fetch | I2 |
|-------|----|
| Decode | I1 |
| R Read | |
| Execute | |
| R Write | |

Clock

… cycle

Instructions

Data

| I1 |
|----|
| I2 |
| I3 |

| D1 |
|----|
| D2 |

NORTHWEST INSTITUTE for ADVANCED COMPUTING

26

Pacific Northwest
NATIONAL LABORATORY
Proudly Operated by Battelle
for the U.S. Department of Energy

UNIVERSITY of
WASHINGTON

# Branch Prediction

- Load the instructions we think will be branched to
- And their successors

Third is fetched

Second moves to decode

First moves to Read

| Fetch | I3 |
|-------|----|
| Decode | I2 |
| R Read | I1 |
| Execute | |
| R Write | |

Fetch

Instructions

Data

Load/Store

| I1 |
|----|
| I2 |
| I3 |

| D1 |
|----|
| D2 |

Clock

... cycle

NORTHWEST INSTITUTE for ADVANCED COMPUTING

Pacific Northwest
NATIONAL LABORATORY
Proudly Operated by Battelle
for the U.S. Department of Energy

UNIVERSITY of
WASHINGTON

# Instruction Pipelining

- When instruction is executed we were either right
  - Continue the pipeline
- Or wrong
  - Flush the pipeline

Fourth is fetched

Fetch

Instructions

| Fetch | I4 |
|-------|-----|
| Decode | I3 |
| R Read | I2 |
| Execute | I1 |
| R Write | |

| I1 |
|----|
| I2 |
| I3 |

Data

| D1 |
|----|
| D2 |

Load/Store

Clock

... cycle

Previous instructions move along

# Pipeline Stall from Mis-Predict

- A single instruction may require multiple steps from fetch to completion

Start fetching from where we branch to

We need to discard the other instructions

Fetch

Instructions

Data

Load/Store

| I1 |
| I2 |
| I3 |

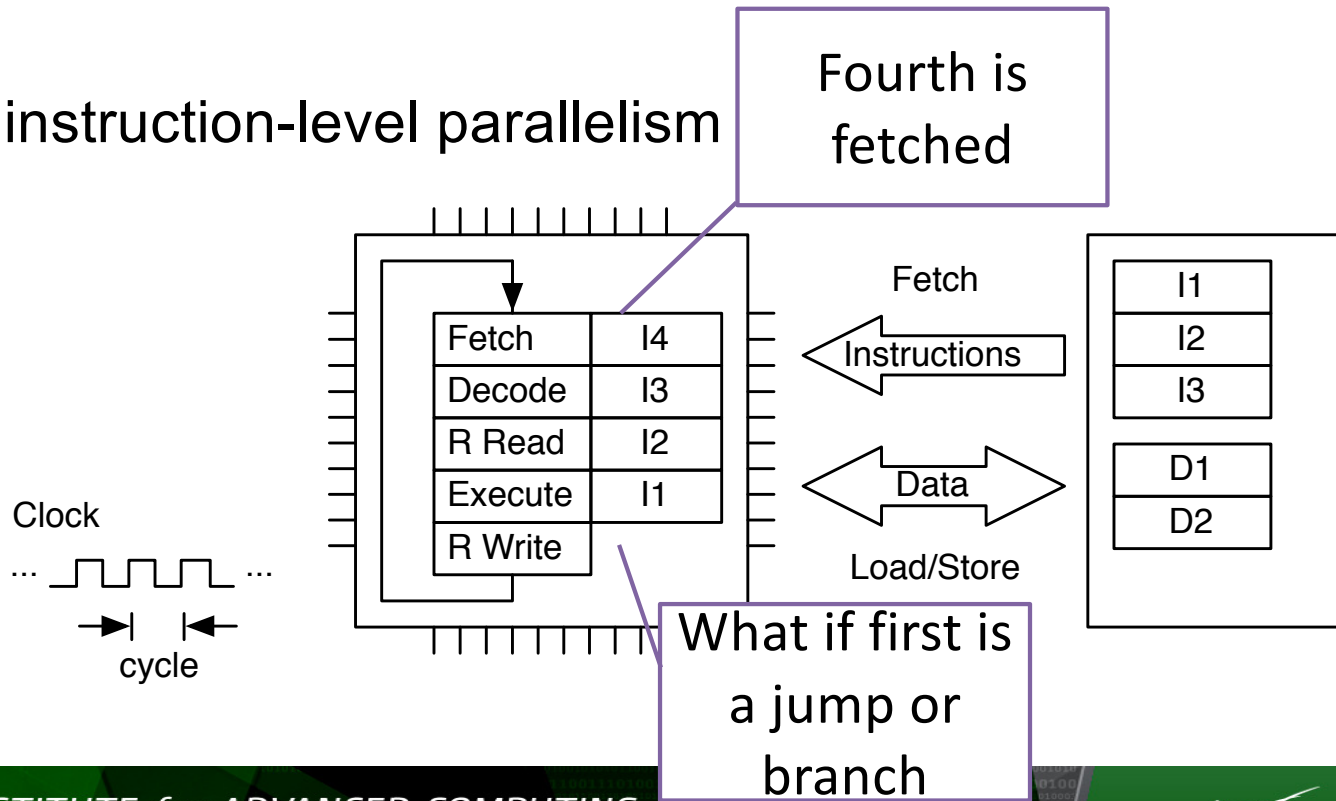| D1 |
| D2 |

Fetch
Decode
R Read
Execute
R Write    I1

Clock

... ⊓⊔⊓⊔⊓ ...

cycle

# Performance-Oriented Architecture Features

- Execution Pipeline
  - Stages of functionality to process issued instructions
  - Hazards are conflicts with continued execution
  - Forwarding supports closely associated operations exhibiting precedence constraints
- Out of Order Execution
  - Uses reservation stations
  - Hides some core latencies and provide fine grain asynchronous operation supporting concurrency
- Branch Prediction
  - Permits computation to proceed at a conditional branch point prior to resolving predicate value
  - Overlaps follow-on computation with predicate resolution
  - Requires roll-back or equivalent to correct false guesses
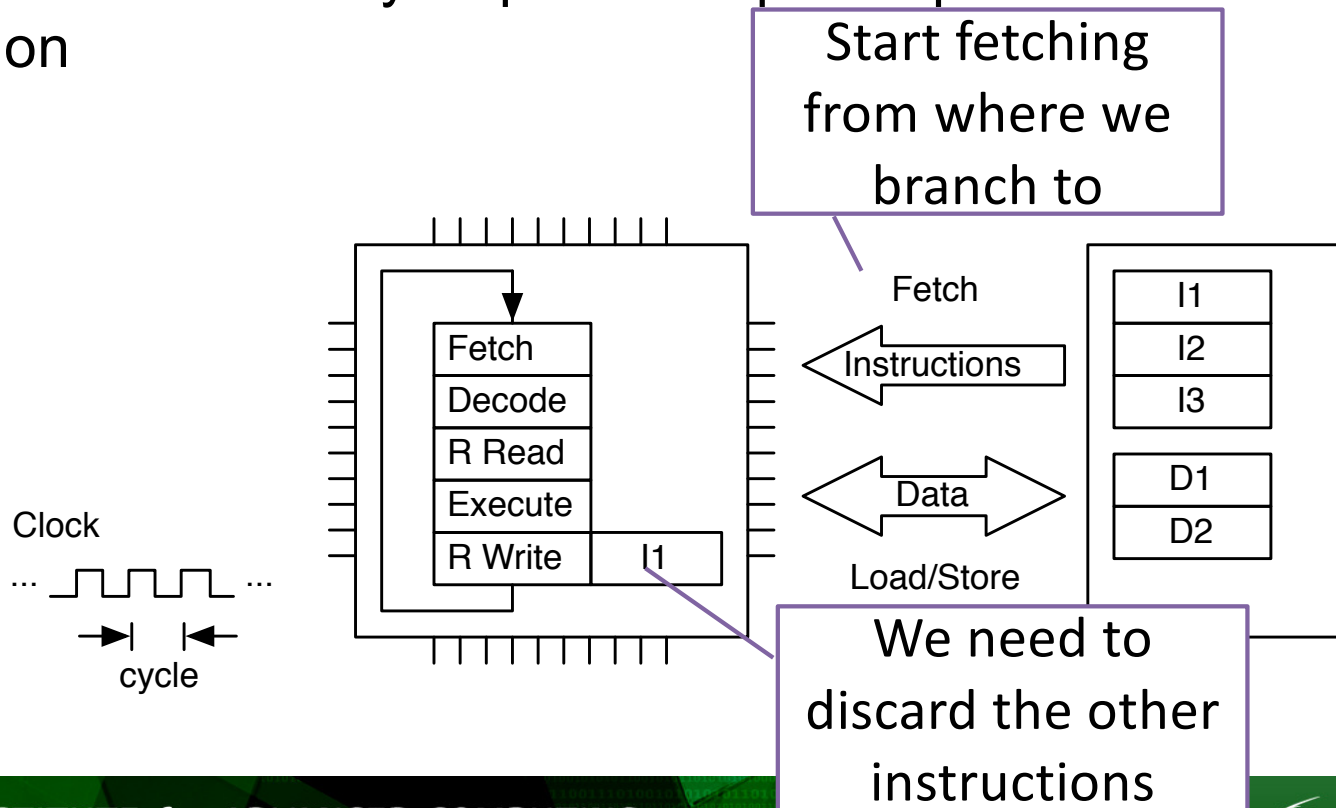  - Sometimes follows both paths, and several deep

NORTHWEST INSTITUTE for ADVANCED COMPUTING

Pacific Northwest
NATIONAL LABORATORY
Proudly Operated by Battelle
for the U.S. Department of Energy

UNIVERSITY of WASHINGTON

30

AMATH 483/583 High-Performance Scientific Computing Spring 2019
University of Washington by Andrew Lumsdaine

# Pipeline Stall

- By pipelining, multiple instructions can be executed at each clock cycle
- Form of instruction-level parallelism

Fourth is fetched

| Fetch | I4 |
|---|---|
| Decode | I3 |
| R Read | I2 |
| Execute | I1 |
| R Write | |

Clock

... ⊓⊔⊓⊔ ...

cycle

Fetch

Instructions

Data

Load/Store

| I1 |
|---|
| I2 |
| I3 |

| D1 |
|---|
| D2 |

What if first is a jump or branch

NORTHWEST INSTITUTE for ADVANCED COMPUTING

Pacific Northwest
NATIONAL LABORATORY
Proudly Operated by Battelle
for the U.S. Department of Energy

UNIVERSITY of
WASHINGTON

# Pipeline Stall

- A single instruction may require multiple steps from fetch to completion



Start fetching from where we branch to

We need to discard the other instructions

Fetch

Decode

R Read

Execute

R Write    I1

Clock

... cycle

Fetch

Instructions

Data

Load/Store

I1
I2
I3

D1
D2

NORTHWEST INSTITUTE for ADVANCED COMPUTING

Pacific Northwest
NATIONAL LABORATORY
Proudly Operated by Battelle
for the U.S. Department of Energy

UNIVERSITY of WASHINGTON

# Compiling functions

```cpp
#include <iostream>
#include <cmath>

double sqrt583(double z) {
  double x = 1.0;

  for (size_t i = 0; i < 32; ++i) {
    double dx = - (x*x-z) / (2.0*x) ;
    x += dx;
    if (abs(dx) < 1.e-9) break;
  }

  return x;
}


int main () {

  std::cout << sqrt583(2.0) << std::endl;


  return 0;
}
```

```
$ c++ main.cpp
$ ./a.out
1.4142
```

Compile main.cpp

```
$ c++ main.cpp
```

Translate it into a language the cpu can run

```
$ ./a.out
```

The executable (program that the cpu can run)

But what is this really?

Pacific Northwest
NATIONAL LABORATORY
Proudly Operated by Battelle
for the U.S. Department of Energy

UNIVERSITY of
WASHINGTON

# Compiled language

```cpp
#include <iostream>
#include <cmath>

double sqrt583(double z) {
  double x = 1.0;

  for (size_t i = 0; i < 32; ++i) {
    double dx = - (x*x-z) / (2.0*x) ;
    x += dx;
    if (abs(dx) < 1.e-9) break;
  }

  return x;
}


int main () {

  std::cout << sqrt583(2.0) << std::endl;


  return 0;
}
```
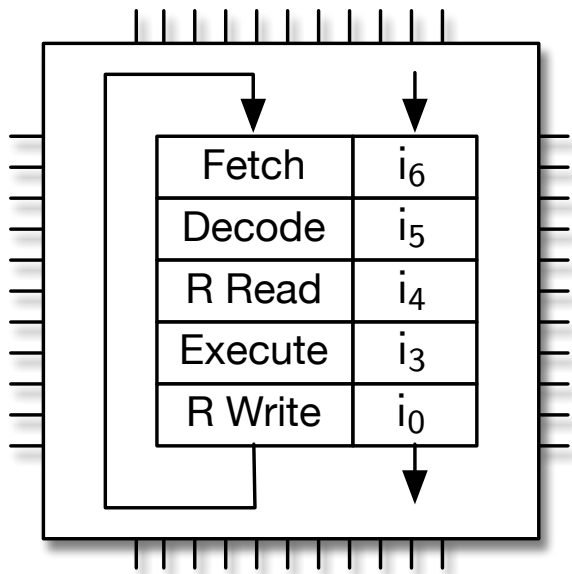
$ c++ main.cpp

main

```asm
subq     $64, %rsp
movsd    LCPI1_0(%rip), %xmm0
movl     $0, -36(%rbp)
movsd    %xmm0, -48(%rbp)
movsd    -48(%rbp), %xmm0
callq    __Z7sqrt583d
movq     %rax, -24(%rbp)
movq     %rdi, -32(%rbp)
movq     -24(%rbp), %rdi
subq     *-32(%rbp)
...
```

sqrt583

```asm
movsd    LCPI0_0(%rip), %xmm1
movsd    %xmm0, -16(%rbp)
movsd    %xmm1, -24(%rbp)
movq     $0, -32(%rbp)
cmpq     $32, -32(%rbp)
jae      LBB0_6
movsd    LCPI0_1(%rip), %xmm0
movsd    LCPI0_3(%rip), %xmm1
movabsq  $-9223372036854, %rax
movsd    -24(%rbp), %xmm2
...
```

Pacific Northwest
NATIONAL LABORATORY
Proudly Operated by Battelle
for the U.S. Department of Energy

UNIVERSITY of
WASHINGTON

# Fetch Decode Execute

CPU instructions are stored in memory

"main" entry point

"main" function



Instructions

"sqrt583" entry point

Data

"sqrt583" function

```
main    subq     $64, %rsp
        movsd    LCPI1_0(%rip), %xmm0
        movl     $0, −36(%rbp)
        movsd    %xmm0, −48(%rbp)
        movsd    −48(%rbp), %xmm0
        callq    __Z7sqrt583d
        movq     %rax, −24(%rbp)
        movq     %rdi, −32(%rbp)
        movq     −24(%rbp), %rdi
        subq     *−32(%rbp)
        . . .

sqrt583 movsd    LCPI0_0(%rip), %xmm1
        movsd    %xmm0, −16(%rbp)
        movsd    %xmm1, −24(%rbp)
        movq     $0, −32(%rbp)
        cmpq     $32, −32(%rbp)
        jae      LBB0_6
        movsd    LCPI0_1(%rip), %xmm0
        movsd    LCPI0_3(%rip), %xmm1
        movabsq  $−9223372036854, %rax
        movsd    −24(%rbp), %xmm2
        . . .
```

# Function Call



Call sqrt583

"main" function
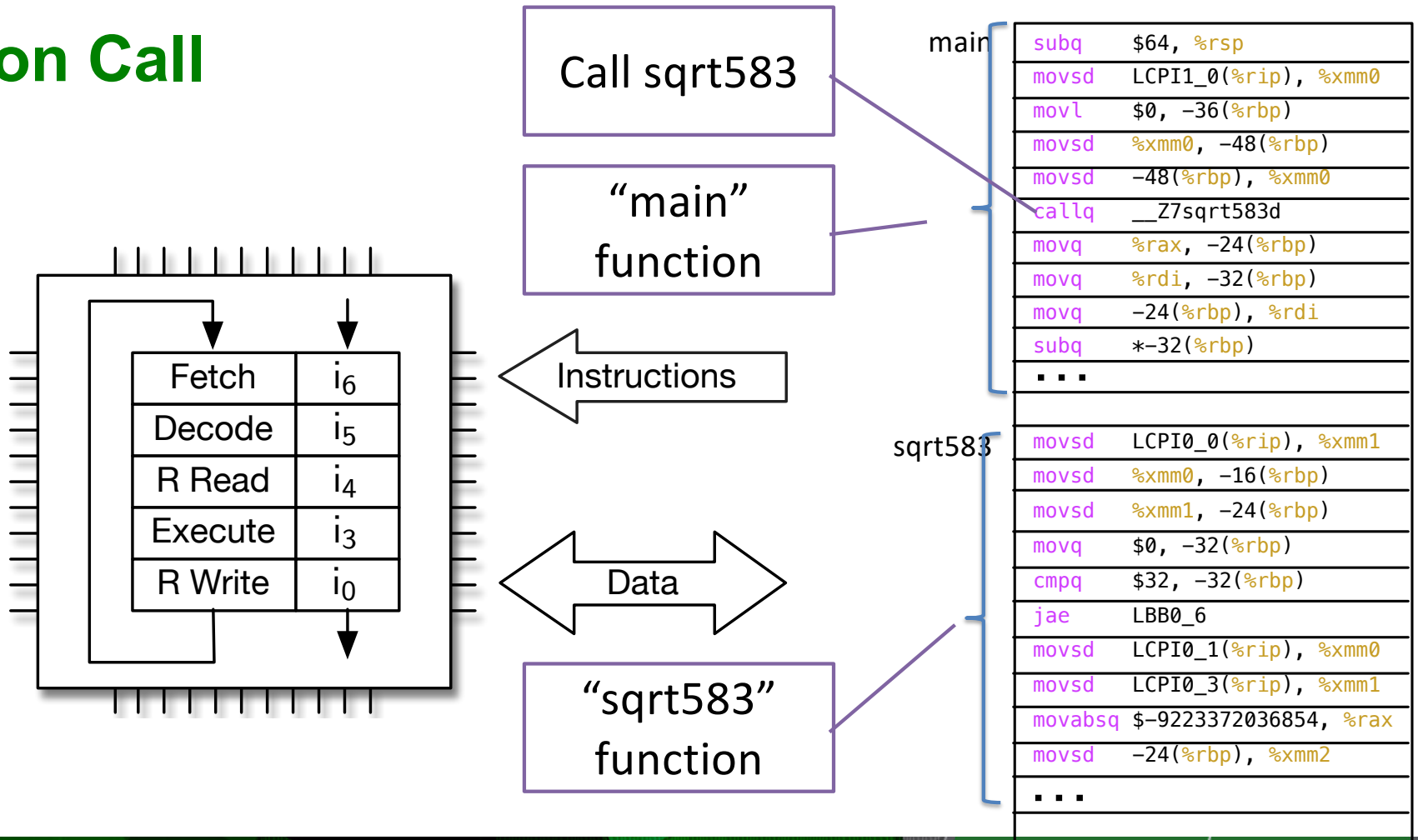
Instructions

Data

"sqrt583" function
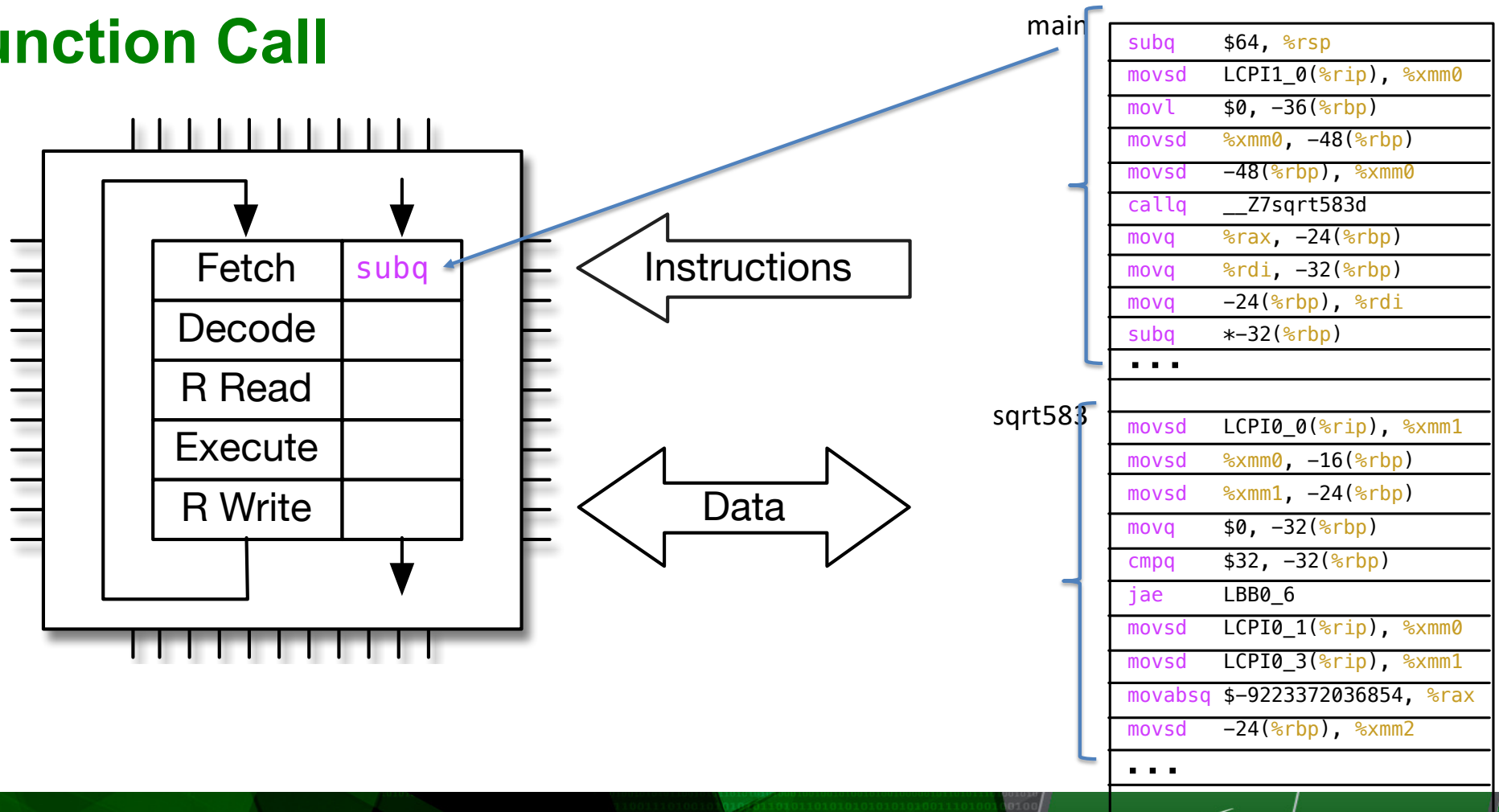
main

```
subq    $64, %rsp
movsd   LCPI1_0(%rip), %xmm0
movl    $0, -36(%rbp)
movsd   %xmm0, -48(%rbp)
movsd   -48(%rbp), %xmm0
callq   __Z7sqrt583d
movq    %rax, -24(%rbp)
movq    %rdi, -32(%rbp)
movq    -24(%rbp), %rdi
subq    *-32(%rbp)
. . .
```

sqrt583

```
movsd    LCPI0_0(%rip), %xmm1
movsd    %xmm0, -16(%rbp)
movsd    %xmm1, -24(%rbp)
movq     $0, -32(%rbp)
cmpq     $32, -32(%rbp)
jae      LBB0_6
movsd    LCPI0_1(%rip), %xmm0
movsd    LCPI0_3(%rip), %xmm1
movabsq  $-9223372036854, %rax
movsd    -24(%rbp), %xmm2
. . .
```
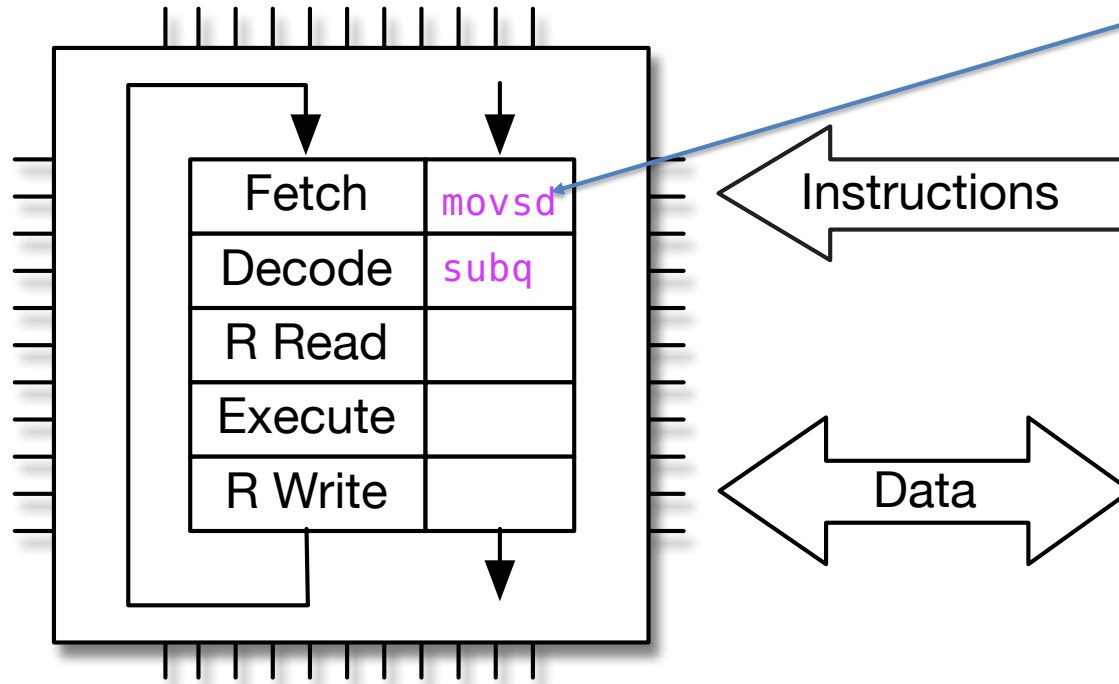
Fetch | $i_6$
Decode | $i_5$
R Read | $i_4$
Execute | $i_3$
R Write | $i_0$

NORTHWEST INSTITUTE for ADVANCED COMPUTING

# Function Call



Instructions

Data

| main | | |
|---|---|---|
| subq | $64, %rsp | |
| movsd | LCPI1_0(%rip), %xmm0 | |
| movl | $0, −36(%rbp) | |
| movsd | %xmm0, −48(%rbp) | |
| movsd | −48(%rbp), %xmm0 | |
| callq | __Z7sqrt583d | |
| movq | %rax, −24(%rbp) | |
| movq | %rdi, −32(%rbp) | |
| movq | −24(%rbp), %rdi | |
| subq | *−32(%rbp) | |
| . . . | | |

| sqrt583 | | |
|---|---|---|
| movsd | LCPI0_0(%rip), %xmm1 | |
| movsd | %xmm0, −16(%rbp) | |
| movsd | %xmm1, −24(%rbp) | |
| movq | $0, −32(%rbp) | |
| cmpq | $32, −32(%rbp) | |
| jae | LBB0_6 | |
| movsd | LCPI0_1(%rip), %xmm0 | |
| movsd | LCPI0_3(%rip), %xmm1 | |
| movabsq | $−9223372036854, %rax | |
| movsd | −24(%rbp), %xmm2 | |
| . . . | | |

# Function Call



Processor diagram showing Fetch (movsd), Decode (subq), R Read, Execute, R Write pipeline stages, with Instructions and Data arrows.
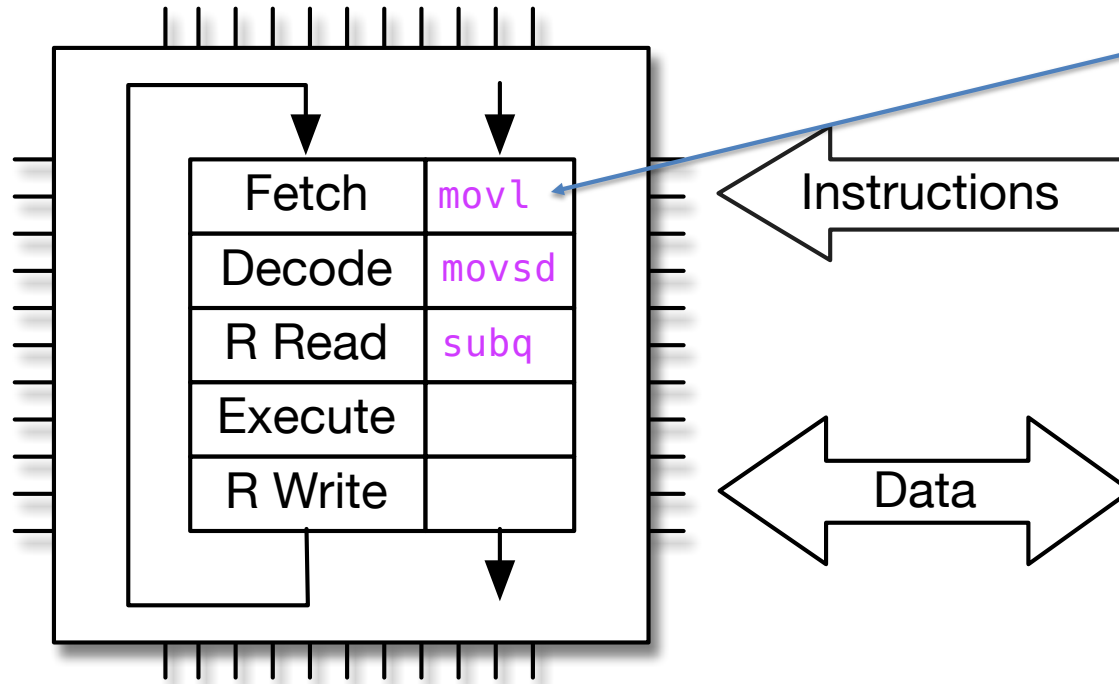
main
```
subq     $64, %rsp
movsd    LCPI1_0(%rip), %xmm0
movl     $0, -36(%rbp)
movsd    %xmm0, -48(%rbp)
movsd    -48(%rbp), %xmm0
callq    __Z7sqrt583d
movq     %rax, -24(%rbp)
movq     %rdi, -32(%rbp)
movq     -24(%rbp), %rdi
subq     *-32(%rbp)
. . .
```

sqrt583
```
movsd    LCPI0_0(%rip), %xmm1
movsd    %xmm0, -16(%rbp)
movsd    %xmm1, -24(%rbp)
movq     $0, -32(%rbp)
cmpq     $32, -32(%rbp)
jae      LBB0_6
movsd    LCPI0_1(%rip), %xmm0
movsd    LCPI0_3(%rip), %xmm1
movabsq  $-9223372036854, %rax
movsd    -24(%rbp), %xmm2
. . .
```

NORTHWEST INSTITUTE for ADVANCED COMPUTING

38

AMATH 483/583 High-Performance Scientific Computing Spring 2019
University of Washington by Andrew Lumsdaine

Pacific Northwest
NATIONAL LABORATORY
Proudly Operated by Battelle
for the U.S. Department of Energy

UNIVERSITY of
WASHINGTON

# Function Call
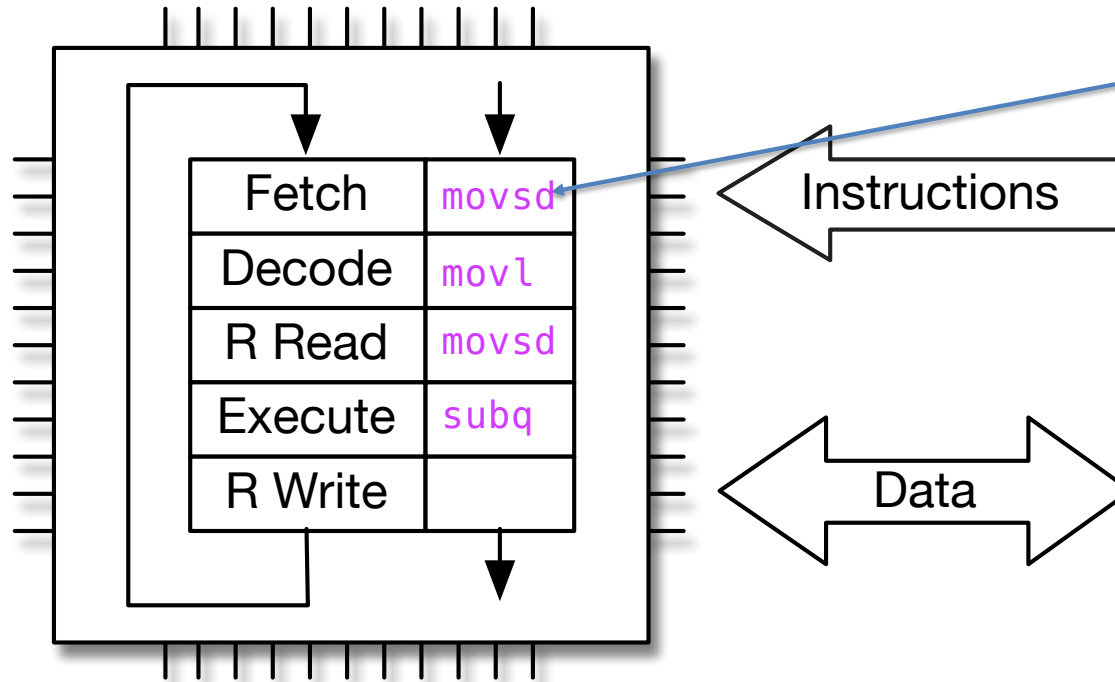


```
main
    subq      $64, %rsp
    movsd     LCPI1_0(%rip), %xmm0
    movl      $0, -36(%rbp)
    movsd     %xmm0, -48(%rbp)
    movsd     -48(%rbp), %xmm0
    callq     __Z7sqrt583d
    movq      %rax, -24(%rbp)
    movq      %rdi, -32(%rbp)
    movq      -24(%rbp), %rdi
    subq      *-32(%rbp)
    . . .

sqrt583
    movsd     LCPI0_0(%rip), %xmm1
    movsd     %xmm0, -16(%rbp)
    movsd     %xmm1, -24(%rbp)
    movq      $0, -32(%rbp)
    cmpq      $32, -32(%rbp)
    jae       LBB0_6
    movsd     LCPI0_1(%rip), %xmm0
    movsd     LCPI0_3(%rip), %xmm1
    movabsq   $-9223372036854, %rax
    movsd     -24(%rbp), %xmm2
    . . .
```

NORTHWEST INSTITUTE for ADVANCED COMPUTING

39

AMATH 483/583 High-Performance Scientific Computing Spring 2019
University of Washington by Andrew Lumsdaine

Pacific Northwest
NATIONAL LABORATORY
Proudly Operated by Battelle
for the U.S. Department of Energy

UNIVERSITY of
WASHINGTON

# Function Call



| | | |
|---|---|---|
| subq | $64, %rsp | main |
| movsd | LCPI1_0(%rip), %xmm0 | |
| movl | $0, -36(%rbp) | |
| movsd | %xmm0, -48(%rbp) | |
| movsd | -48(%rbp), %xmm0 | |
| callq | __Z7sqrt583d | |
| movq | %rax, -24(%rbp) | |
| movq | %rdi, -32(%rbp) | |
| movq | -24(%rbp), %rdi | |
| subq | *-32(%rbp) | |
| . . . | | |

| | | |
|---|---|---|
| movsd | LCPI0_0(%rip), %xmm1 | sqrt583 |
| movsd | %xmm0, -16(%rbp) | |
| movsd | %xmm1, -24(%rbp) | |
| movq | $0, -32(%rbp) | |
| cmpq | $32, -32(%rbp) | |
| jae | LBB0_6 | |
| movsd | LCPI0_1(%rip), %xmm0 | |
| movsd | LCPI0_3(%rip), %xmm1 | |
| movabsq | $-9223372036854, %rax | |
| movsd | -24(%rbp), %xmm2 | |
| . . . | | |

Fetch — movsd
Decode — movl
R Read — movsd
Execute — subq
R Write

Instructions

Data

NORTHWEST INSTITUTE for ADVANCED COMPUTING

40

AMATH 483/583 High-Performance Scientific Computing Spring 2019
University of Washington by Andrew Lumsdaine

Pacific Northwest
NATIONAL LABORATORY
Proudly Operated by Battelle
for the U.S. Department of Energy

UNIVERSITY of
WASHINGTON

# Function Call



```
main    subq    $64, %rsp
        movsd   LCPI1_0(%rip), %xmm0
        movl    $0, -36(%rbp)
        movsd   %xmm0, -48(%rbp)
        movsd   -48(%rbp), %xmm0
        callq   __Z7sqrt583d
        movq    %rax, -24(%rbp)
        movq    %rdi, -32(%rbp)
        movq    -24(%rbp), %rdi
        subq    *-32(%rbp)
        . . .

sqrt583 movsd   LCPI0_0(%rip), %xmm1
        movsd   %xmm0, -16(%rbp)
        movsd   %xmm1, -24(%rbp)
        movq    $0, -32(%rbp)
        cmpq    $32, -32(%rbp)
        jae     LBB0_6
        movsd   LCPI0_1(%rip), %xmm0
        movsd   LCPI0_3(%rip), %xmm1
        movabsq $-9223372036854, %rax
        movsd   -24(%rbp), %xmm2
        . . .
```
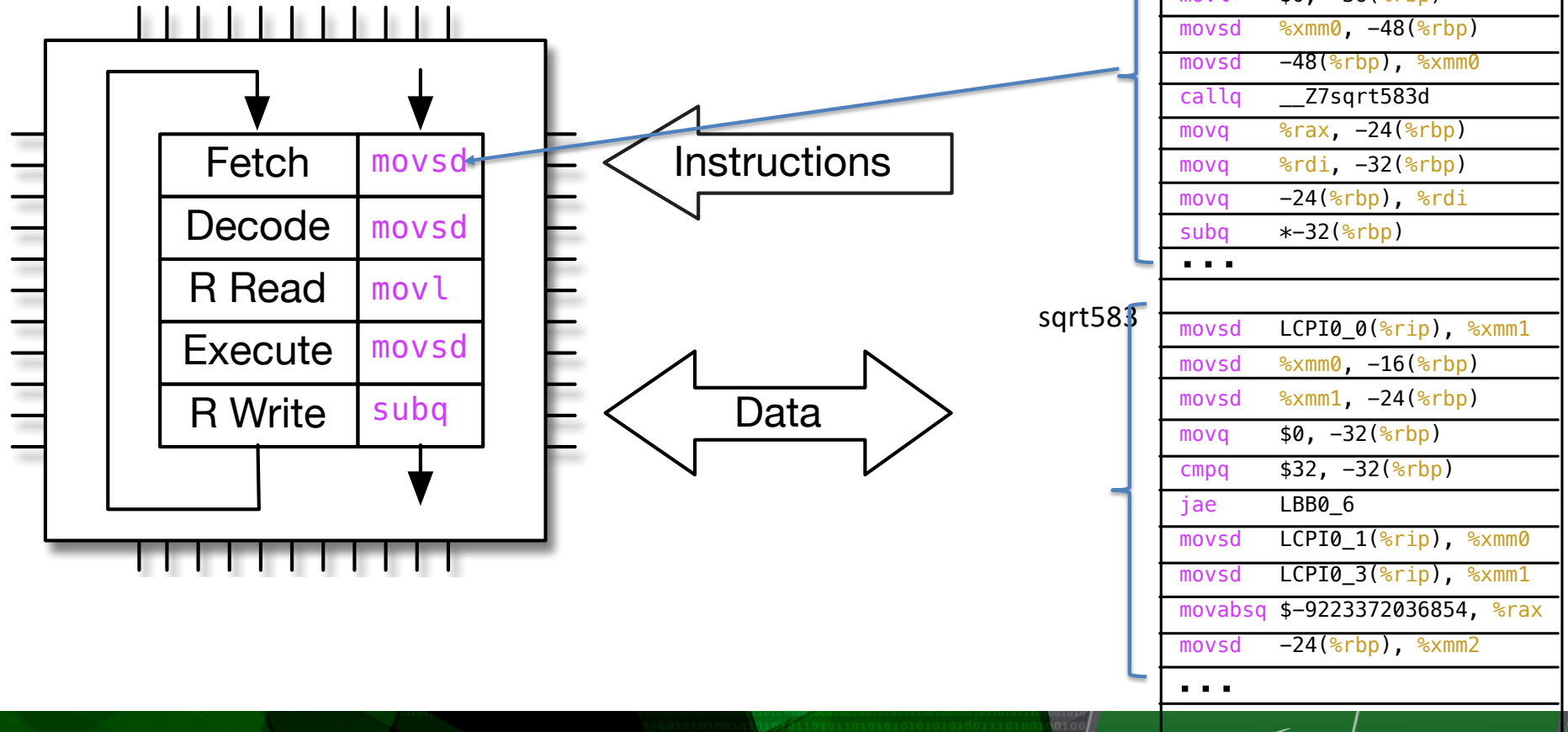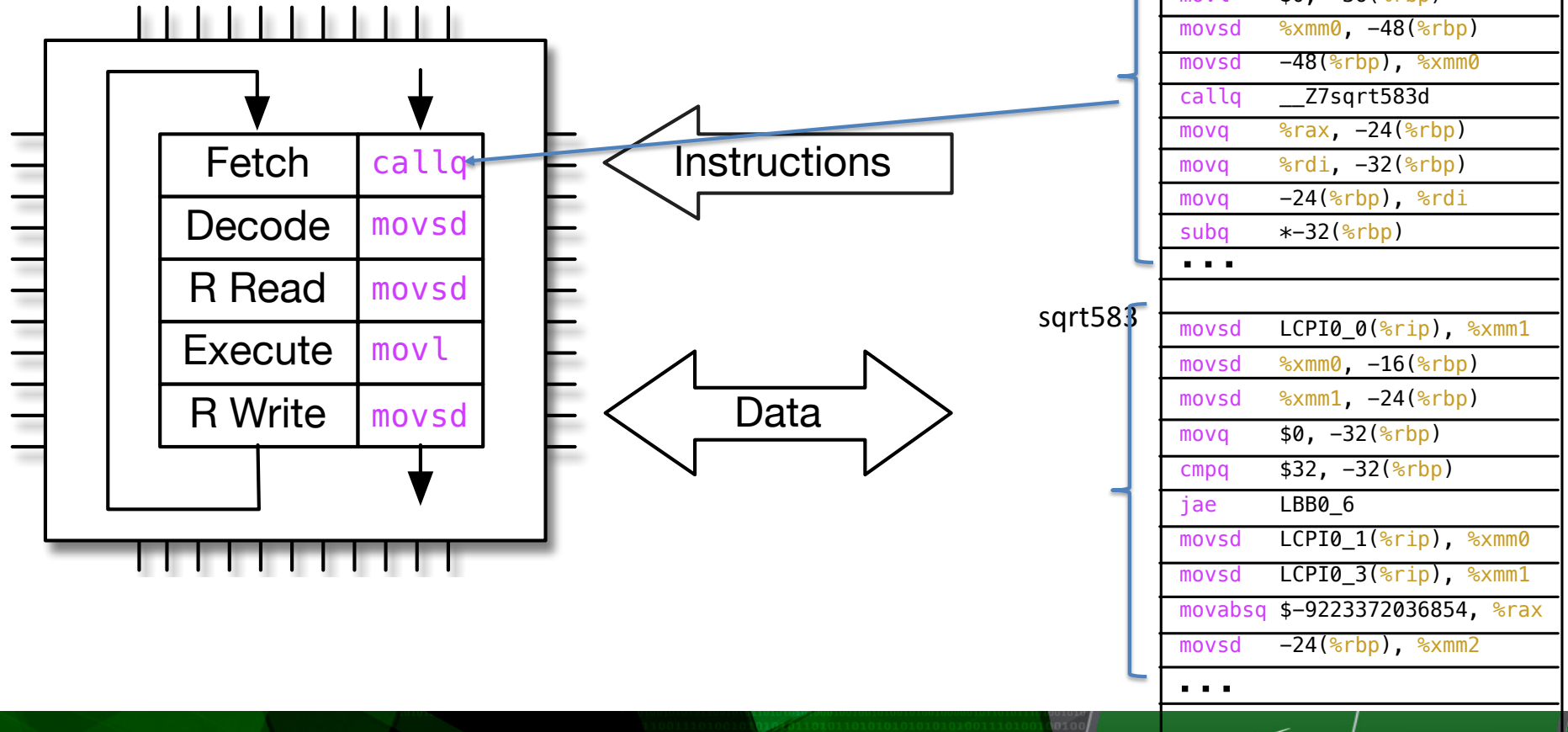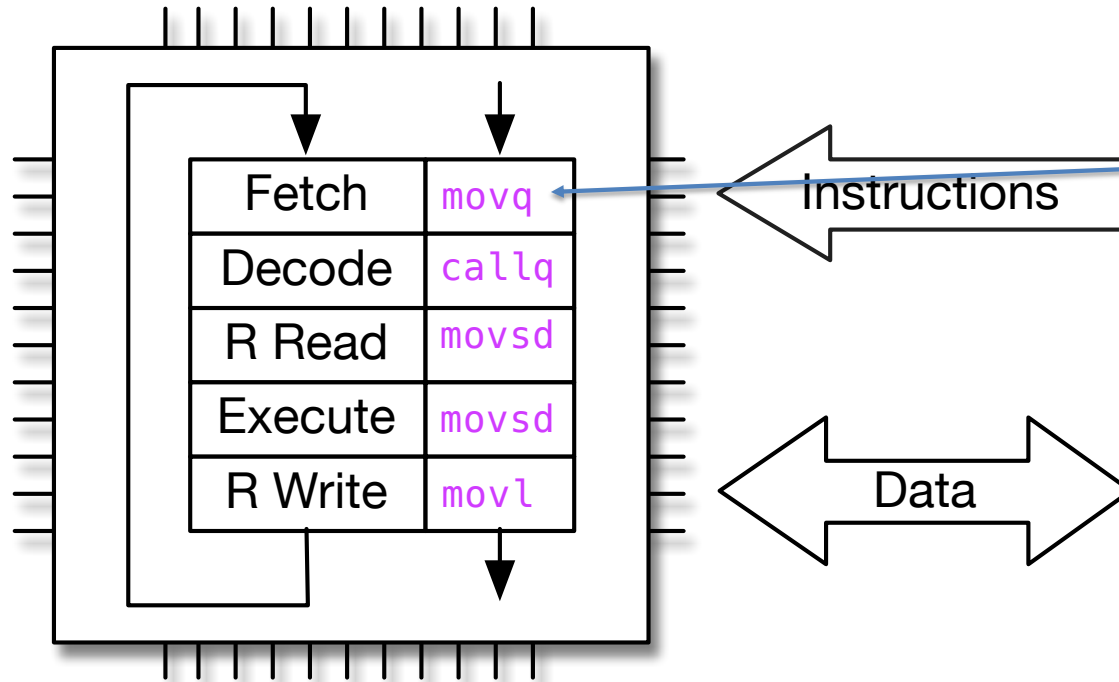
# Function Call



main
| | |
|---|---|
| subq | $64, %rsp |
| movsd | LCPI1_0(%rip), %xmm0 |
| movl | $0, −36(%rbp) |
| movsd | %xmm0, −48(%rbp) |
| movsd | −48(%rbp), %xmm0 |
| callq | __Z7sqrt583d |
| movq | %rax, −24(%rbp) |
| movq | %rdi, −32(%rbp) |
| movq | −24(%rbp), %rdi |
| subq | *−32(%rbp) |
| . . . | |

sqrt583
| | |
|---|---|
| movsd | LCPI0_0(%rip), %xmm1 |
| movsd | %xmm0, −16(%rbp) |
| movsd | %xmm1, −24(%rbp) |
| movq | $0, −32(%rbp) |
| cmpq | $32, −32(%rbp) |
| jae | LBB0_6 |
| movsd | LCPI0_1(%rip), %xmm0 |
| movsd | LCPI0_3(%rip), %xmm1 |
| movabsq | $−9223372036854, %rax |
| movsd | −24(%rbp), %xmm2 |
| . . . | |

Fetch — callq
Decode — movsd
R Read — movsd
Execute — movl
R Write — movsd

Instructions

Data

NORTHWEST INSTITUTE for ADVANCED COMPUTING

42

AMATH 483/583 High-Performance Scientific Computing Spring 2019
University of Washington by Andrew Lumsdaine

Pacific Northwest
NATIONAL LABORATORY
Proudly Operated by Battelle
for the U.S. Department of Energy

UNIVERSITY of
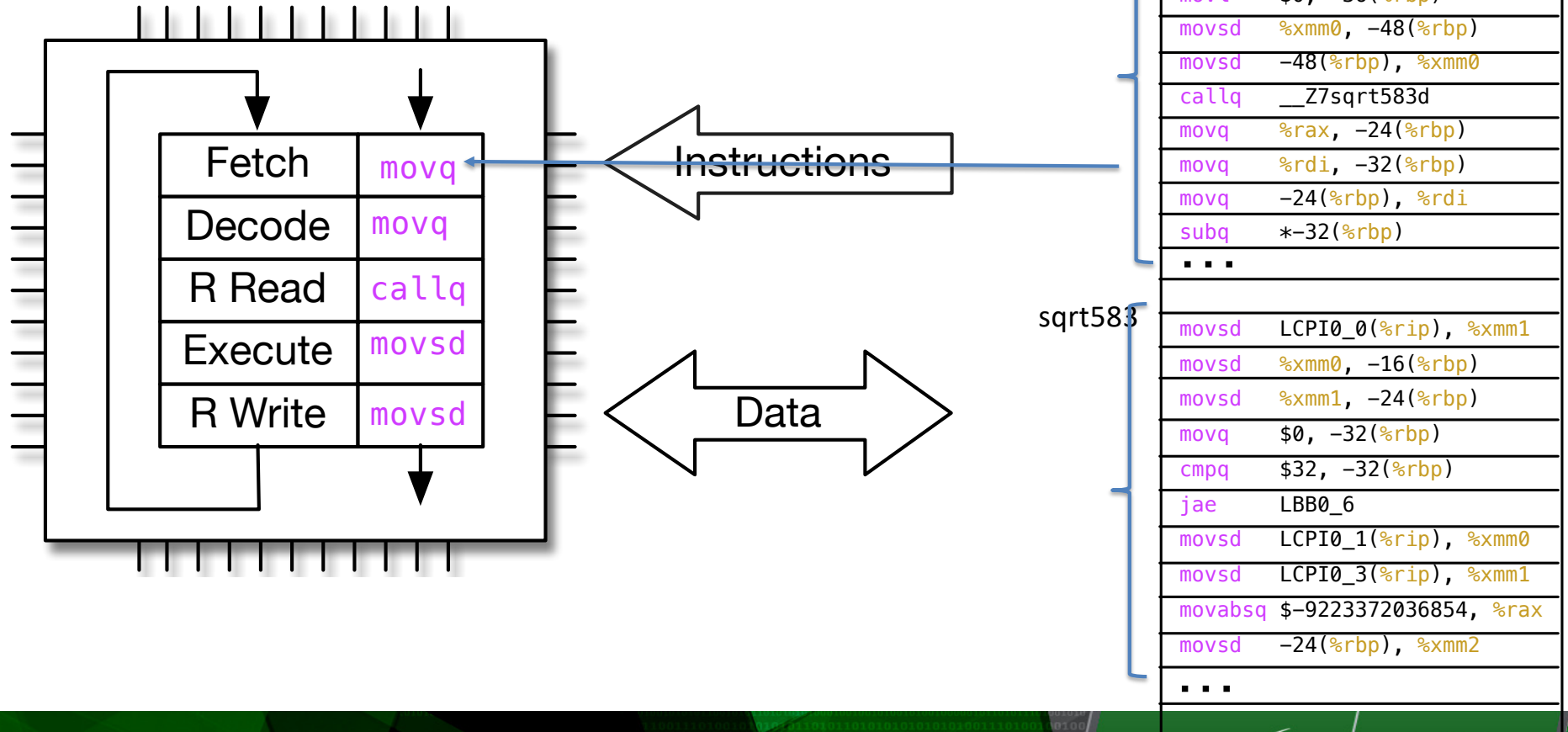WASHINGTON

# Function Call

# Function Call



main
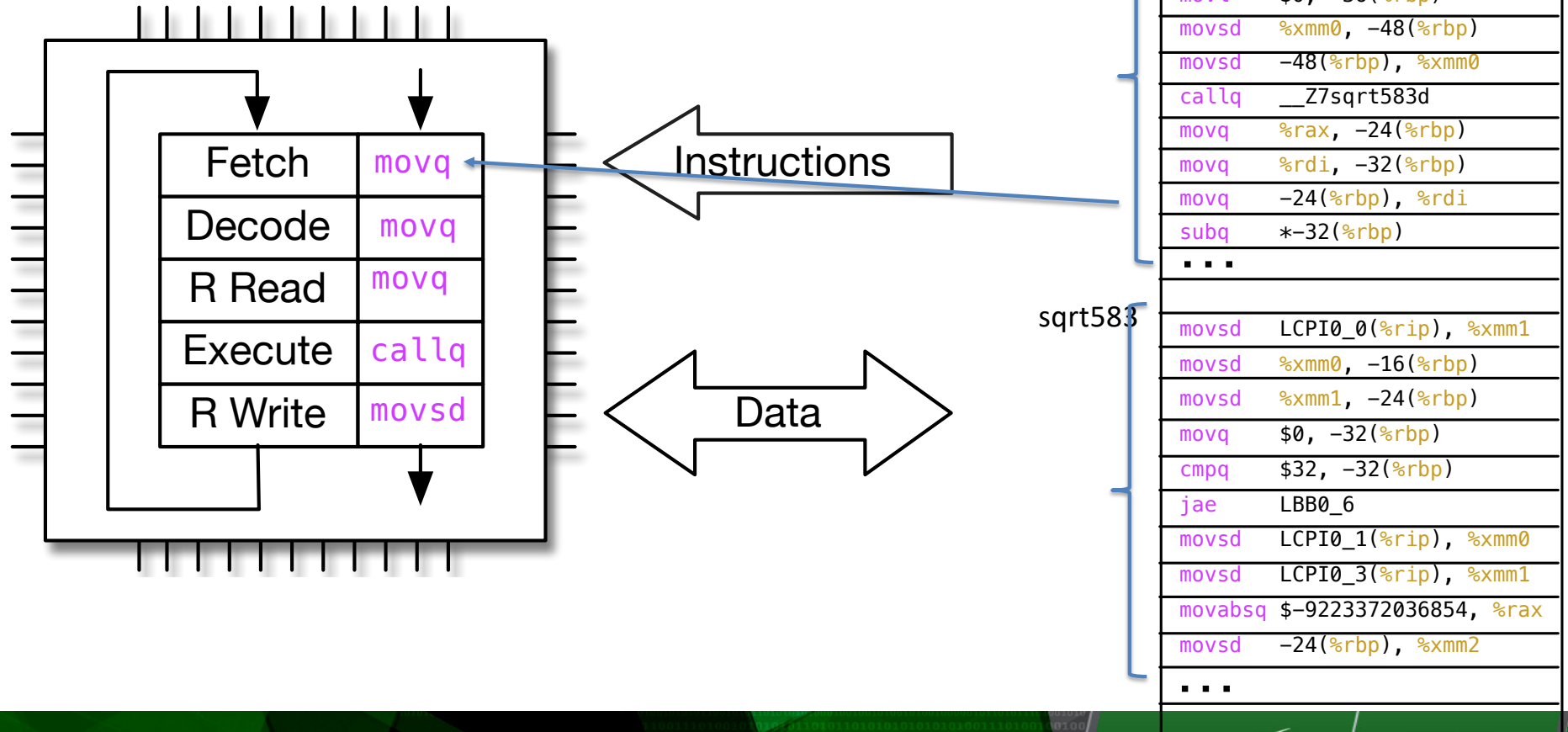```
subq     $64, %rsp
movsd    LCPI1_0(%rip), %xmm0
movl     $0, -36(%rbp)
movsd    %xmm0, -48(%rbp)
movsd    -48(%rbp), %xmm0
callq    __Z7sqrt583d
movq     %rax, -24(%rbp)
movq     %rdi, -32(%rbp)
movq     -24(%rbp), %rdi
subq     *-32(%rbp)
. . .
```

sqrt583
```
movsd    LCPI0_0(%rip), %xmm1
movsd    %xmm0, -16(%rbp)
movsd    %xmm1, -24(%rbp)
movq     $0, -32(%rbp)
cmpq     $32, -32(%rbp)
jae      LBB0_6
movsd    LCPI0_1(%rip), %xmm0
movsd    LCPI0_3(%rip), %xmm1
movabsq  $-9223372036854, %rax
movsd    -24(%rbp), %xmm2
. . .
```
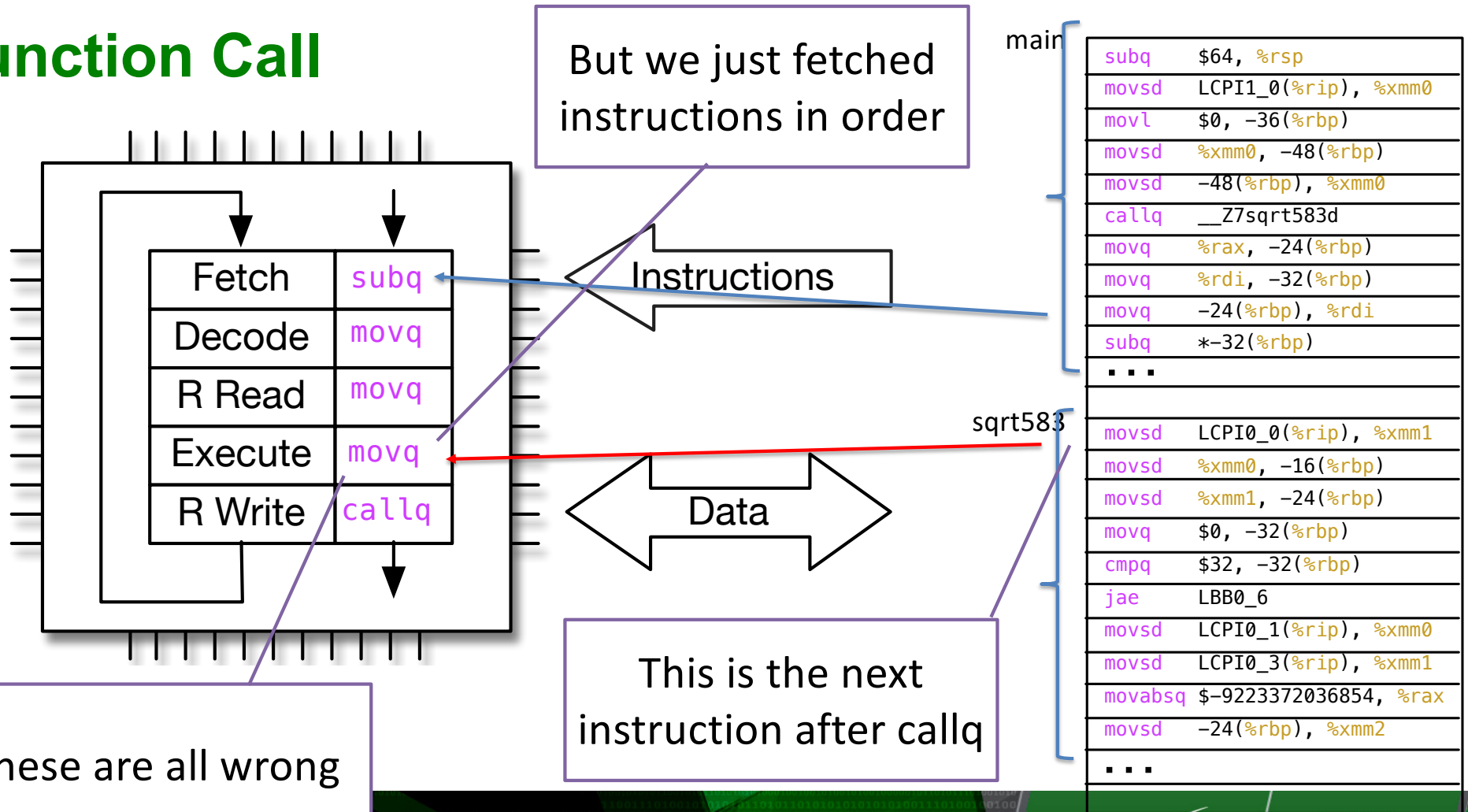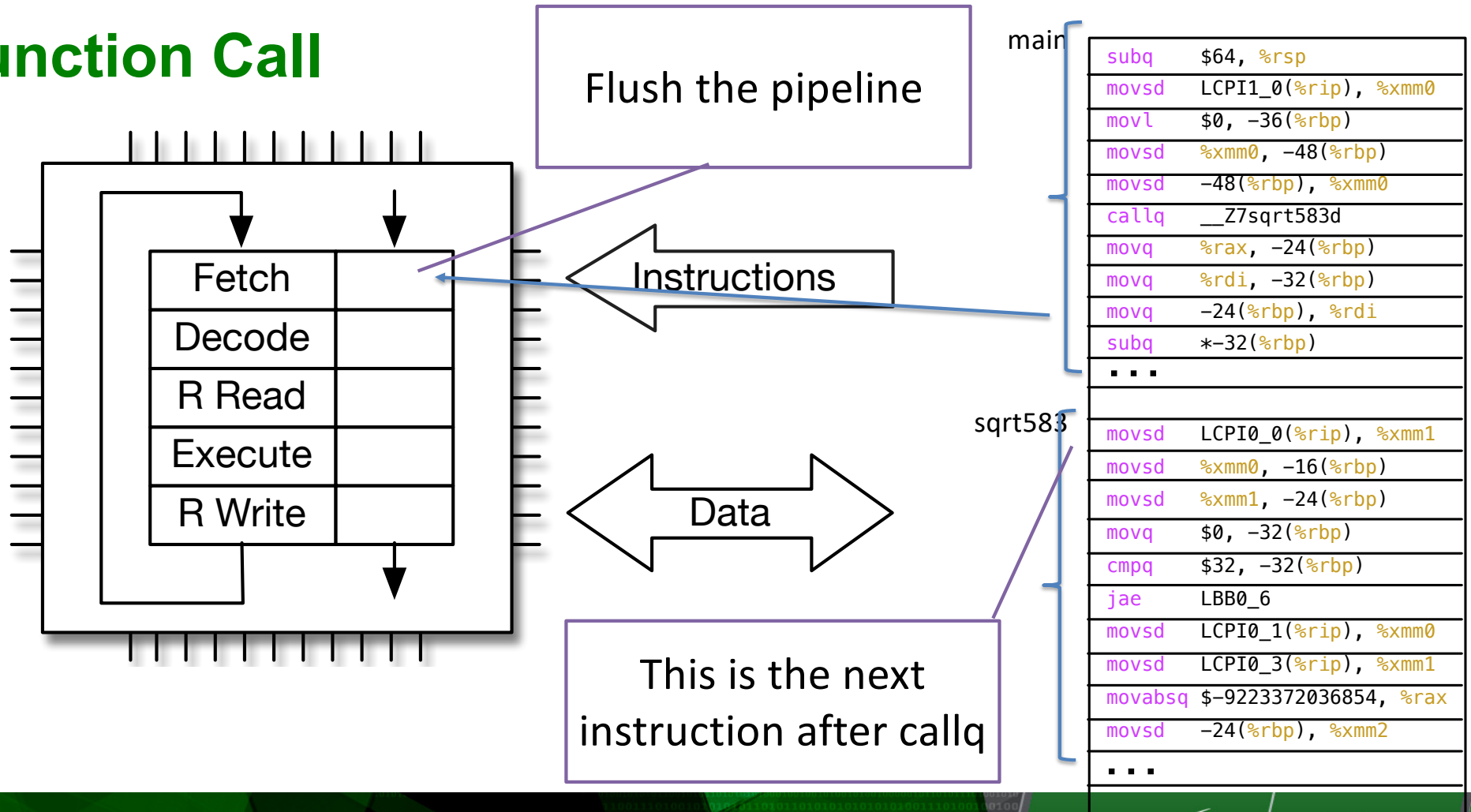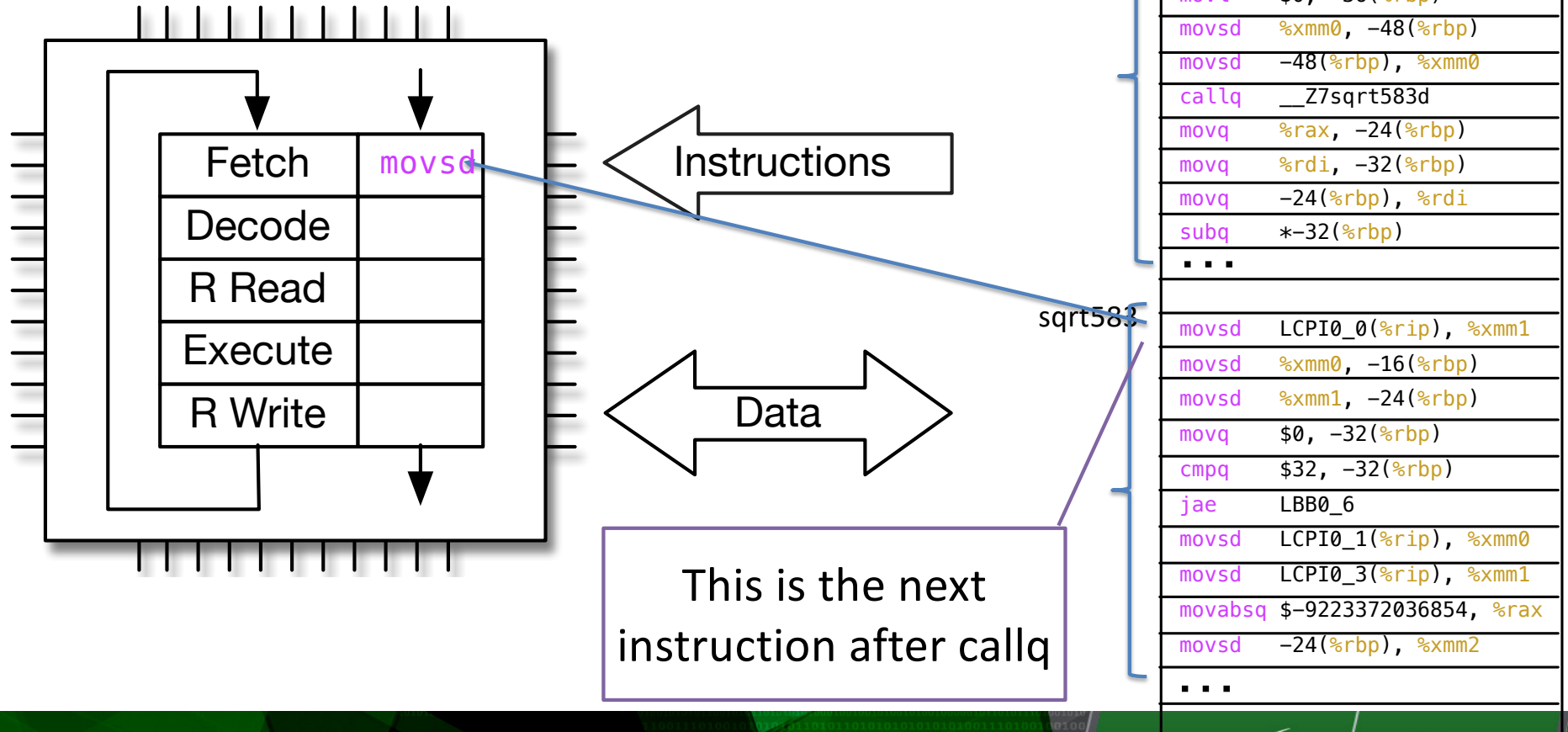
Instructions

Data

Fetch — movq
Decode — movq
R Read — callq
Execute — movsd
R Write — movsd

NORTHWEST INSTITUTE for ADVANCED COMPUTING

44

AMATH 483/583 High-Performance Scientific Computing Spring 2019
University of Washington by Andrew Lumsdaine

Pacific Northwest
NATIONAL LABORATORY
Proudly Operated by Battelle
for the U.S. Department of Energy

UNIVERSITY of
WASHINGTON

# Function Call



main

| | |
|---|---|
| subq | $64, %rsp |
| movsd | LCPI1_0(%rip), %xmm0 |
| movl | $0, −36(%rbp) |
| movsd | %xmm0, −48(%rbp) |
| movsd | −48(%rbp), %xmm0 |
| callq | __Z7sqrt583d |
| movq | %rax, −24(%rbp) |
| movq | %rdi, −32(%rbp) |
| movq | −24(%rbp), %rdi |
| subq | *−32(%rbp) |
| . . . | |

sqrt583

| | |
|---|---|
| movsd | LCPI0_0(%rip), %xmm1 |
| movsd | %xmm0, −16(%rbp) |
| movsd | %xmm1, −24(%rbp) |
| movq | $0, −32(%rbp) |
| cmpq | $32, −32(%rbp) |
| jae | LBB0_6 |
| movsd | LCPI0_1(%rip), %xmm0 |
| movsd | LCPI0_3(%rip), %xmm1 |
| movabsq | $−9223372036854, %rax |
| movsd | −24(%rbp), %xmm2 |
| . . . | |

# Function Call

But we just fetched instructions in order

| | |
|---|---|
| Fetch | subq |
| Decode | movq |
| R Read | movq |
| Execute | movq |
| R Write | callq |

Instructions

Data

These are all wrong

This is the next instruction after callq

main

```
subq      $64, %rsp
movsd     LCPI1_0(%rip), %xmm0
movl      $0, -36(%rbp)
movsd     %xmm0, -48(%rbp)
movsd     -48(%rbp), %xmm0
callq     __Z7sqrt583d
movq      %rax, -24(%rbp)
movq      %rdi, -32(%rbp)
movq      -24(%rbp), %rdi
subq      *-32(%rbp)
. . .
```

sqrt583

```
movsd     LCPI0_0(%rip), %xmm1
movsd     %xmm0, -16(%rbp)
movsd     %xmm1, -24(%rbp)
movq      $0, -32(%rbp)
cmpq      $32, -32(%rbp)
jae       LBB0_6
movsd     LCPI0_1(%rip), %xmm0
movsd     LCPI0_3(%rip), %xmm1
movabsq   $-9223372036854, %rax
movsd     -24(%rbp), %xmm2
. . .
```

ADVANCED COMPUTING

Pacific Northwest
NATIONAL LABORATORY

Proudly Operated by Battelle
for the U.S. Department of Energy

UNIVERSITY of
WASHINGTON

# Function Call



Flush the pipeline

Instructions

Data

This is the next instruction after callq

main

```
subq    $64, %rsp
movsd   LCPI1_0(%rip), %xmm0
movl    $0, -36(%rbp)
movsd   %xmm0, -48(%rbp)
movsd   -48(%rbp), %xmm0
callq   __Z7sqrt583d
movq    %rax, -24(%rbp)
movq    %rdi, -32(%rbp)
movq    -24(%rbp), %rdi
subq    *-32(%rbp)
. . .
```

sqrt583

```
movsd   LCPI0_0(%rip), %xmm1
movsd   %xmm0, -16(%rbp)
movsd   %xmm1, -24(%rbp)
movq    $0, -32(%rbp)
cmpq    $32, -32(%rbp)
jae     LBB0_6
movsd   LCPI0_1(%rip), %xmm0
movsd   LCPI0_3(%rip), %xmm1
movabsq $-9223372036854, %rax
movsd   -24(%rbp), %xmm2
. . .
```

# Function Call



main
```
subq     $64, %rsp
movsd    LCPI1_0(%rip), %xmm0
movl     $0, -36(%rbp)
movsd    %xmm0, -48(%rbp)
movsd    -48(%rbp), %xmm0
callq    __Z7sqrt583d
movq     %rax, -24(%rbp)
movq     %rdi, -32(%rbp)
movq     -24(%rbp), %rdi
subq     *-32(%rbp)
. . .
```

sqrt583
```
movsd    LCPI0_0(%rip), %xmm1
movsd    %xmm0, -16(%rbp)
movsd    %xmm1, -24(%rbp)
movq     $0, -32(%rbp)
cmpq     $32, -32(%rbp)
jae      LBB0_6
movsd    LCPI0_1(%rip), %xmm0
movsd    LCPI0_3(%rip), %xmm1
movabsq  $-9223372036854, %rax
movsd    -24(%rbp), %xmm2
. . .
```

Instructions

Data

This is the next instruction after callq

# Function Call



main

| | |
|---|---|
| subq | $64, %rsp |
| movsd | LCPI1_0(%rip), %xmm0 |
| movl | $0, -36(%rbp) |
| movsd | %xmm0, -48(%rbp) |
| movsd | -48(%rbp), %xmm0 |
| callq | __Z7sqrt583d |
| movq | %rax, -24(%rbp) |
| movq | %rdi, -32(%rbp) |
| movq | -24(%rbp), %rdi |
| subq | *-32(%rbp) |
| ... | |

sqrt583

| | |
|---|---|
| movsd | LCPI0_0(%rip), %xmm1 |
| movsd | %xmm0, -16(%rbp) |
| movsd | %xmm1, -24(%rbp) |
| movq | $0, -32(%rbp) |
| cmpq | $32, -32(%rbp) |
| jae | LBB0_6 |
| movsd | LCPI0_1(%rip), %xmm0 |
| movsd | LCPI0_3(%rip), %xmm1 |
| movabsq | $-9223372036854, %rax |
| movsd | -24(%rbp), %xmm2 |
| ... | |

Instructions

Data

This is the next instruction after callq

NORTHWEST INSTITUTE for ADVANCED COMPUTING

49

AMATH 483/583 High-Performance Scientific Computing Spring 2019
University of Washington by Andrew Lumsdaine

Pacific Northwest
NATIONAL LABORATORY
Proudly Operated by Battelle
for the U.S. Department of Energy

UNIVERSITY of WASHINGTON

# Function Call



main
| | |
|---|---|
| subq | $64, %rsp |
| movsd | LCPI1_0(%rip), %xmm0 |
| movl | $0, -36(%rbp) |
| movsd | %xmm0, -48(%rbp) |
| movsd | -48(%rbp), %xmm0 |
| callq | __Z7sqrt583d |
| movq | %rax, -24(%rbp) |
| movq | %rdi, -32(%rbp) |
| movq | -24(%rbp), %rdi |
| subq | *-32(%rbp) |
| . . . | |

sqrt583
| | |
|---|---|
| movsd | LCPI0_0(%rip), %xmm1 |
| movsd | %xmm0, -16(%rbp) |
| movsd | %xmm1, -24(%rbp) |
| movq | $0, -32(%rbp) |
| cmpq | $32, -32(%rbp) |
| jae | LBB0_6 |
| movsd | LCPI0_1(%rip), %xmm0 |
| movsd | LCPI0_3(%rip), %xmm1 |
| movabsq | $-9223372036854, %rax |
| movsd | -24(%rbp), %xmm2 |
| . . . | |

Instructions

Data

This is the next instruction after callq

NORTHWEST INSTITUTE for ADVANCED COMPUTING

50

AMATH 483/583 High-Performance Scientific Computing Spring 2019
University of Washington by Andrew Lumsdaine

Pacific Northwest
NATIONAL LABORATORY
Proudly Operated by Battelle
for the U.S. Department of Energy

UNIVERSITY of
WASHINGTON

# Function Call



main
```
subq      $64, %rsp
movsd     LCPI1_0(%rip), %xmm0
movl      $0, -36(%rbp)
movsd     %xmm0, -48(%rbp)
movsd     -48(%rbp), %xmm0
callq     __Z7sqrt583d
movq      %rax, -24(%rbp)
movq      %rdi, -32(%rbp)
movq      -24(%rbp), %rdi
subq      *-32(%rbp)
. . .
```

sqrt583
```
movsd     LCPI0_0(%rip), %xmm1
movsd     %xmm0, -16(%rbp)
movsd     %xmm1, -24(%rbp)
movq      $0, -32(%rbp)
cmpq      $32, -32(%rbp)
jae       LBB0_6
movsd     LCPI0_1(%rip), %xmm0
movsd     LCPI0_3(%rip), %xmm1
movabsq   $-9223372036854, %rax
movsd     -24(%rbp), %xmm2
. . .
```

Instructions

Data

This is the next instruction after callq

**NORTHWEST INSTITUTE for ADVANCED COMPUTING**

AMATH 483/583 High-Performance Scientific Computing Spring 2019
University of Washington by Andrew Lumsdaine

Pacific Northwest
NATIONAL LABORATORY
Proudly Operated by Battelle
for the U.S. Department of Energy

UNIVERSITY of WASHINGTON

# Function Call



main
```
subq     $64, %rsp
movsd    LCPI1_0(%rip), %xmm0
movl     $0, -36(%rbp)
movsd    %xmm0, -48(%rbp)
movsd    -48(%rbp), %xmm0
callq    __Z7sqrt583d
movq     %rax, -24(%rbp)
movq     %rdi, -32(%rbp)
movq     -24(%rbp), %rdi
subq     *-32(%rbp)
...
```

sqrt583
```
movsd    LCPI0_0(%rip), %xmm1
movsd    %xmm0, -16(%rbp)
movsd    %xmm1, -24(%rbp)
movq     $0, -32(%rbp)
cmpq     $32, -32(%rbp)
jae      LBB0_6
movsd    LCPI0_1(%rip), %xmm0
movsd    LCPI0_3(%rip), %xmm1
movabsq  $-9223372036854, %rax
movsd    -24(%rbp), %xmm2
...
```

Fetch — cmpq
Decode — movq
R Read — movsd
Execute — movsd
R Write — movsd

Instructions

Data

This is the next instruction after callq

# Function Call



```
main
        subq    $64, %rsp
        movsd   LCPI1_0(%rip), %xmm0
        movl    $0, -36(%rbp)
        movsd   %xmm0, -48(%rbp)
        movsd   -48(%rbp), %xmm0
        callq   __Z7sqrt583d
        movq    %rax, -24(%rbp)
        movq    %rdi, -32(%rbp)
        movq    -24(%rbp), %rdi
        subq    *-32(%rbp)
        . . .

sqrt583
        movsd   LCPI0_0(%rip), %xmm1
        movsd   %xmm0, -16(%rbp)
        movsd   %xmm1, -24(%rbp)
        movq    $0, -32(%rbp)
        cmpq    $32, -32(%rbp)
        jae     LBB0_6
        movsd   LCPI0_1(%rip), %xmm0
        movsd   LCPI0_3(%rip), %xmm1
        movabsq $-9223372036854, %rax
        movsd   -24(%rbp), %xmm2
        . . .
```

This is the next instruction after callq

# Pipeline flush: Bad

Flush the pipeline

Instructions

Data

This is the next instruction after callq

| Fetch |
| Decode |
| R Read |
| Execute |
| R Write |

main
```
subq      $64, %rsp
movsd     LCPI1_0(%rip), %xmm0
movl      $0, -36(%rbp)
movsd     %xmm0, -48(%rbp)
movsd     -48(%rbp), %xmm0
callq     __Z7sqrt583d
movq      %rax, -24(%rbp)
movq      %rdi, -32(%rbp)
movq      -24(%rbp), %rdi
subq      *-32(%rbp)
. . .
```

sqrt583
```
movsd     LCPI0_0(%rip), %xmm1
movsd     %xmm0, -16(%rbp)
movsd     %xmm1, -24(%rbp)
movq      $0, -32(%rbp)
cmpq      $32, -32(%rbp)
jae       LBB0_6
movsd     LCPI0_1(%rip), %xmm0
movsd     LCPI0_3(%rip), %xmm1
movabsq   $-9223372036854, %rax
movsd     -24(%rbp), %xmm2
. . .
```

NORTHWEST INSTITUTE for ADVANCED COMPUTING

54

AMATH 483/583 High-Performance Scientific Computing Spring 2019
University of Washington by Andrew Lumsdaine

Pacific Northwest
NATIONAL LABORATORY
Proudly Operated by Battelle
for the U.S. Department of Energy

UNIVERSITY of WASHINGTON

# Memory Access

The next one may be cheaper

- What are typical costs for accessing memory?
- What is typical clock cycle time?
- How many clock cycles to fetch an instruction? `200`
- How many clock cycles to execute load / store instruction? `400`
- CPI for load / store?

`600`

Clock

... ⊓_⊓_⊓_ ...

→| |← cycle

`0.5 ns`

CPU

Fetch
← Instructions

Data ⟷

Load/Store

`100 ns`

Memory

NORTHWEST INSTITUTE for ADVANCED COMPUTING

Pacific Northwest
NATIONAL LABORATORY
Proudly Operated by Battelle
for the U.S. Department of Energy

UNIVERSITY of
WASHINGTON

# Memory Access Costs

- Access to main memory has huge impact on performance

# Memory Access Costs

- Access to main memory has huge impact on performance
- *Latency*: How long does the first access to data take
- *Bandwidth*: How much data can we continuously fetch

# Memory Access Costs

- Access to main memory has huge impact on performance (600X)
- Processor would be idle almost all the time

# Cache

Small memory near FDE unit

Very very fast (and expensive)

Fetch
Decode
R Read
Execute
R Write

Fetch

Instructions

Data

Load/Store

0.5 ns

Cache

Fetch

Instructions

Data

Load/Store

100 ns

I1
I2
I3

D1
D2

Clock

... ⊓⊔⊓ ...

cycle

# Hierarchical Memory

Registers
(immediately fast)

Level 1 Cache
(very very fast)

Separate L1 for
instuctions/data

| F | r0 |
|---|---|
| D | r1 |
| R | r2 |
| E | r3 |
| | r4 |
| W | r5 |

Fetch

Instructions

L1
(I)

L2

Fetch

Instructions

| I1 |
|---|
| I2 |
| I3 |

| D1 |
|---|
| D2 |

Data

Load/Store

L1
(D)

Data

Load/Store

Clock

... ⊓⊔⊓⊔⊓ ...

cycle

0.5 ns

5 ns

100 ns

Level 2 Cache
(pretty fast)

# Hierarchical Memory



FDE works with data in registers

There is also an MMU and TLB

Data goes from L2 to L1

Data goes from L1 to registers

Data goes from main memory to L2

F D R E W

r0 r1 r2 r3 r4 r5

Fetch
Instructions

Data
Load/Store

L1 (I)

L1 (D)

L2

Fetch
Instructions

Data
Load/Store

I1
I2
I3

D1
D2

Clock

… cycle …

0.5 ns

5 ns

100 ns

AMATH 483/583 High-Performance Scientific Computing Spring 2019
University of Washington by Andrew Lumsdaine

Pacific NATIONAL LABORATORY
Proudly Operated by Battelle
for the U.S. Department of Energy

UNIVERSITY of WASHINGTON

# Cache and Multicore



Separate L1 and L2 for each core

Cores work on separate register sets and instrs

Cores work on separate register sets and instrs

Shared L3

Main memory is shared

NORTHWEST INSTITUTE for ADVANCED COMPUTING

62

AMATH 483/583 High-Performance Scientific Computing Spring 2019
University of Washington by Andrew Lumsdaine

Pacific Northwest
NATIONAL LABORATORY
Proudly Operated by Battelle
for the U.S. Department of Energy

UNIVERSITY of WASHINGTON

# Performance

# Locality → Performance

Keep as much data here as possible

Keep as much data here as possible

Keep as much data here as possible

| F | r0 |
|---|----|
| D | r1 |
| R | r2 |
| E | r3 |
| W | r4 |
|   | r5 |

Fetch
Instructions

Data
Load/Store

L1 (I)

L1 (D)

L2

Fetch
Instructions

Data
Load/Store

| I1 |
|----|
| I2 |
| I3 |

| D1 |
|----|
| D2 |

Clock

... ⊓_⊓_⊓_ ...

cycle

NORTHWEST INSTITUTE for ADVANCED COMPUTING

64

AMATH 483/583 High-Performance Scientific Computing Spring 2019
University of Washington by Andrew Lumsdaine

Pacific Northwest
NATIONAL LABORATORY
Proudly Operated by Battelle
for the U.S. Department of Energy

UNIVERSITY of
WASHINGTON

# Locality → Strategy

If we need an operand here

We first look for it here (L1)

If it is in L2 (*hit*), copy to L1

Can data be missing from main memory?



Fetch
Instructions

L1 (I)

L2

Fetch
Instructions

I1
I2
I3

D1
D2

Data
Load/Store

L1 (D)

Data
Load/Store

F
D
R
E
W

r0
r1
r2
r3
r4
r5

Clock

··· ⊓⊔⊓⊔⊓⊔ ···

cycle

If it is there (*hit*), load

If it is not there (*miss*), look in L2

If it is not in L2, get from main memory

# Locality → Strategy



When we need the next operand

We want it to be here

Or here

On a miss, copy the data we want *and its neighbors*

Clock

··· ⊓⊔⊓⊔⊓⊔ ···

|→| |←|

cycle

# Locality → Strategy

The next operand may be "near" the last

It could be "near" in time or space

Near in time (*temporal locality*): the next operand is a previous operand

Near in space (*spatial locality*): the next operand is in a nearby memory location to a previous operand

F D R E W

r0 r1 r2 r3 r4 r5

Fetch
Instructions

Data
Load/Store

L1 (I)

L1 (D)

L2

Fetch
Instructions

Data

I1 I2 I3

D1 D2

Clock
… ⌐_⌐_⌐_ …
→| |←
cycle

NORTHWEST INSTITUTE for ADVANCED COMPUTING

67

AMATH 483/583 High-Performance Scientific Computing Spring 2019
University of Washington by Andrew Lumsdaine

Pacific Northwest
NATIONAL LABORATORY
Proudly Operated by Battelle
for the U.S. Department of Energy

UNIVERSITY of
WASHINGTON

# Locality → Performance

- Caches are much smaller than main memory. How do we decide what data to keep in cache to effect higher performance (more accesses)?
- **Temporal Locality**: if a program accesses a memory location, there is a much higher than random probability that the same location will be accessed again
  - Cache replacement policies attempt to keep cached elements in the cache for as long as possible
- **Spatial Locality**: if a program accesses a memory location, there is a much higher than random probability that nearby locations will also be accessed (soon)
  - Cache policies read contiguous chunks of data – a referenced element and its neighbors – not just single elements

NORTHWEST INSTITUTE for ADVANCED COMPUTING

68

AMATH 483/583 High-Performance Scientific Computing Spring 2019
University of Washington by Andrew Lumsdaine

Pacific Northwest
NATIONAL LABORATORY
Proudly Operated by Battelle
for the U.S. Department of Energy

UNIVERSITY of WASHINGTON

# Matrix Vector Product

- Recall for ANN $x^{i+1} = S(W^i x^i)$

- In general $y \leftarrow A \times x$

$$x^{i+1} \leftarrow W^i \times x^i$$

num_cols()

$$y_i = \sum_{i=0}^{N-1} A_{ij} x_j, \quad i = 0, \dots, M$$

summation

Two nested loops

num_rows()

```
for (size_t i = 0; i < A.num_rows(); ++i) {
  for (size_t j = 0; j < A.num_cols(); ++j) {
    y(i) += A(i, j) * x(j);
  }
}
```

How many flops?

How many times is this done?

How much data?

# Matrix-matrix product

$$C_{ij} = \sum_{k=0}^{K-1} A_{ik}B_{kj}$$

Three nested loops

```
for (size_t i = 0; i < C.num_rows(); ++i) {
  for (size_t j = 0; j < C.num_cols(); ++j) {
    for (size_t k = 0; k < A.num_cols(); ++k) {
      C(i, j) += A(i, k) * B(k, j);
    }
  }
}
```

How many flops?

How many times is this done?

How much data?

NORTHWEST INSTITUTE for ADVANCED COMPUTING

70

Pacific Northwest
NATIONAL LABORATORY

UNIVERSITY of
WASHINGTON

AMATH 483/583 High-Performance Scientific Computing Spring 2019
University of Washington by Andrew Lumsdaine

Proudly Operated by Battelle
for the U.S. Department of Energy

# Timing and Benchmarking

- Humans have pathological need to see who is better at everything
- But ordering requires a single number corresponding to "goodness"
- Which is impossible of course
- So we take one task and turn that into the definition of goodness (cf IQ)
  - (What is IQ? It's the thing that the IQ test measures.) — My personal rant
- In HPC, we take performance on a particular computational task to rank the worlds computers with the 500 best scores on this task
  - Linear system solution – matrix matrix product at the core
  - Performance = FLOPS = (Total computations) / (Time to compute)
  - Linpack $\rightarrow$ $2N^3$/ (Time to compute)

NORTHWEST INSTITUTE for ADVANCED COMPUTING

71

AMATH 483/583 High-Performance Scientific Computing Spring 2019
University of Washington by Andrew Lumsdaine

Pacific Northwest
NATIONAL LABORATORY
Proudly Operated by Battelle
for the U.S. Department of Energy

UNIVERSITY of
WASHINGTON

# Timing a Program

- The time program in Linux (Unix) will measure time resources a process uses

```
$ time ls -lR /usr > /dev/null

real    0m0.464s
user    0m0.080s
sys     0m0.380s
```

Elapsed Wall Clock time

Time Spent running user code

Time Spent running system code

This is what we'll be using

But finer grained control

# C++ Timer

And this will be provided to you

All you need to worry about

```cpp
class Timer {
private:
  typedef std::chrono::time_point<std::chrono::system_clock> time_t;

public:
  Timer() : startTime(), stopTime() {}

  time_t start()    { return (startTime = std::chrono::system_clock::now()); }
  time_t stop()     { return (stopTime  = std::chrono::system_clock::now()); }
  double elapsed()  { return
      std::chrono::duration_cast<std::chrono::milliseconds>(stopTime-startTime).count(); }

private:
  time_t startTime, stopTime;
};
```

# Measuring Matrix Matrix Product

```cpp
#include <iostream>
#include "Matrix.hpp"
#include "Timer.hpp"
using namespace std;

int main() {
  cout << "N\tElapsed" << endl;

  for (int N = 8; N < 1024; N *= 2) {
    Matrix A(N, N), B(N, N), C(N, N), D(N, N);

    Timer T; T.start();
    A = B*C;
    T.stop();

    cout << N << "\t" << T.elapsed() << endl;
  }

  return 0;
}
```

Declare Timer T

Start Timer T

Stop Timer T

And???

Print Elapsed Time

Insufficient resolution

```
$ ./a.out
N        Elapsed
8        0
16       0
32       0
64       0
128      2
256      28
512      315
```

74

# What All Are We Timing

```cpp
Matrix operator*(const Matrix& A, const Matrix& B) {
    Matrix C(A.num_rows(), B.num_cols());
    zeroize(C);
    for (size_t i = 0; i < A.num_rows
        for (size_t j = 0; j < B.num_cols(); ++j ) {
            for (size_t k = 0; k < A.num_cols(); ++k )
                C(i, j) += A(i, k) * B(k, j);
        }
    }
    return C;
}
```

Allocating a Matrix

Zeroing it out

The actual matrix product

Never allocate new memory in performance critical sections of code

NORTHWEST INSTITUTE for ADVANCED COMPUTING

75

AMATH 483/583 High-Performance Scientific Computing Spring 2019
University of Washington by Andrew Lumsdaine

Pacific Northwest
NATIONAL LABORATORY
Proudly Operated by Battelle
for the U.S. Department of Energy

UNIVERSITY of WASHINGTON

# Just For Benchmarking

```cpp
Matrix operator*(const Matrix& A, const Matrix&B) {
  Matrix C(A.num_rows(), B.num_cols());
  zeroizeMatrix(C);
  multiply(A, B, C);
  return C;
}

void multiply(const Matrix& A, const Matrix&B, Matrix&C) {
  for (size_t i = 0; i < A.num_rows(); ++i) {
    for (size_t j = 0; j < B.num_cols(); ++j) {
      for (size_t k = 0; k < A.num_cols(); ++k) {
        C(i,j) += A(i,k) * B(k,j);
      }
    }
  }
}
```

C++ Core Guideline Violation

F.20: For "out" output values, prefer return values to output parameters
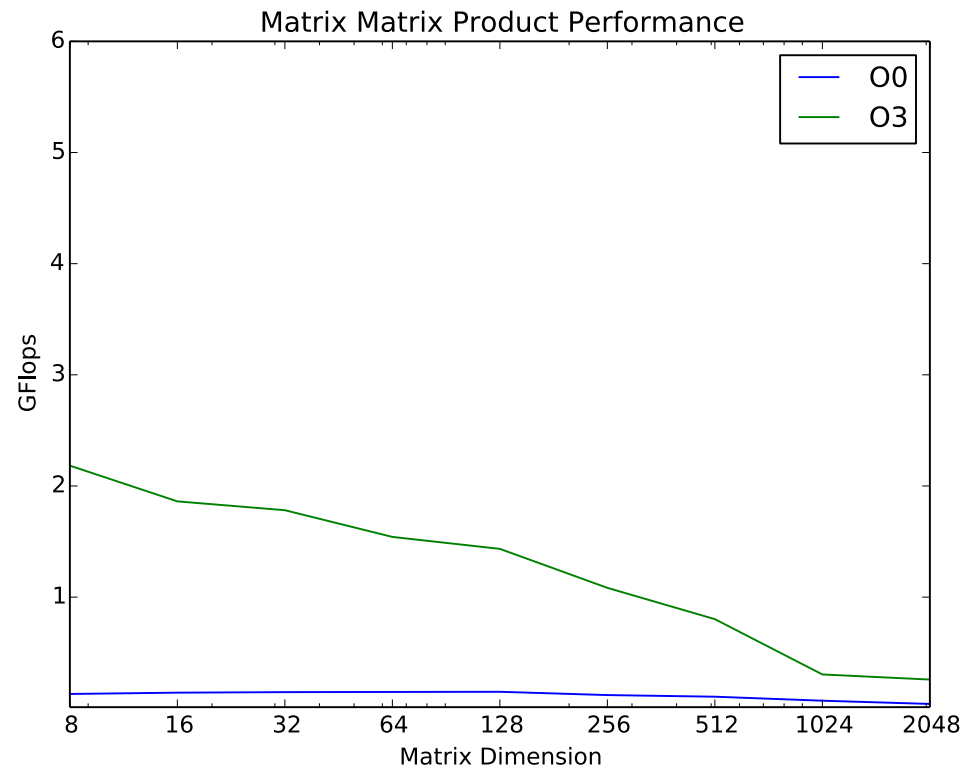
UNIVERSITY of WASHINGTON

# Benchmarking

```cpp
double benchmark(int M, int N, int K, long numruns) {
  Matrix A(M, K), B(K, N), C(M, N);

  Timer T;
  T.start();
  for (int i = 0; i < numruns; ++i) {
    multiply(A, B, C);
  }
  T.stop();

  return T.elapsed();
}
```

Run the core loop many times to get sufficient resolution for small(er) sizes

NORTHWEST INSTITUTE for ADVANCED COMPUTING

Pacific Northwest
NATIONAL LABORATORY

Proudly Operated by Battelle
for the U.S. Department of Energy

UNIVERSITY of
WASHINGTON

# Bonus Question (Advanced Topic)

```cpp
double benchmark(int M, int N, int K, long n
  Matrix A(M, K), B(K, N), C(M, N);

  Timer T;
  T.start();
  for (int i = 0; i < numruns; ++i) {
    multiply(A, B, C);
  }
  T.stop();

  return T.elapsed();
}
```

If we have different multiply routines (and we will), how many of these do we write?

By how much do they differ?

How can we parameterize that?

NORTHWEST INSTITUTE for ADVANCED COMPUTING

Pacific Northwest
NATIONAL LABORATORY
Proudly Operated by Battelle
for the U.S. Department of Energy

UNIVERSITY of WASHINGTON

# Bonus Question (Advanced Topic)

```
double benchmark(int M, int N, int K, long nu
                 <something> f) {
  Matrix A(M, K), B(K, N), C(M, N);

  Timer T;
  T.start();
  for (int i = 0; i < numruns; ++i) {
    f(A, B, C);
  }
  T.stop();

  return T.elapsed();
}
```

We want to pass in something

That we call like a function

Double bonus: It just needs an operator()()

Let's not get carried away

AMATH 483/583 High-Performance Scientific Computing Spring 2019
University of Washington by Andrew Lumsdaine

Pacific Northwest
NATIONAL LABORATORY
Proudly Operated by Battelle
for the U.S. Department of Energy

UNIVERSITY of
WASHINGTON

# Functions as Data

Is a function

And takes two const Matrix& and a Matrix& for args

```
#include <functional>

double benchmark(int M, int N, int K, long numruns,

  function<void (const Matrix&, const Matrix&, Matrix&)>f) {

Matrix A(M, K), B(K, N), C(M, N);

Timer T;
T.start();
for (int i = 0; i < numruns; ++i) {
  f(A, B, C);
}
T.stop();

return T.elapsed();
}
```

Parameter f

That returns void

Like multiply()

```
void multiply(const Matrix& A, const Matrix&B, Matrix&C);
```

AMATH 483/583 High-Performance Scientific Computing Spring 2019
University of Washington by Andrew Lumsdaine

# Functions as Data (Advanced)

Functions returning void

And taking two const Matrix& and a Matrix& for args

```cpp
void multiply(const Matrix& A, const Matrix &B, Matrix& C);
void multiply_2(const Matrix& A, const Matrix &B, Matrix& C);
void yet_another(const Matrix& A, const Matrix &B, Matrix& C);

//...

double t1 = benchmark(100, 100, 100, multiply);
double t2 = benchmark(100, 100, 100, multiply_2);
double t2 = benchmark(100, 100, 100, yet_another);
```

Pass them into another function

# Let's Start Benchmarking

```cpp
double benchmark(int M, int N, int K, long numruns) {
  Matrix A(M, K), B(K, N), C(M, N);

  Timer T;
  T.start();
  for (int i = 0; i < numruns; ++i) {
    multiply(A, B, C);
  }
  T.stop();

  return T.elapsed();
}
```

```
bench: bench.o Matrix.o
c++ -std=c++11 bench.o Matrix.o -o bench

bench.o: bench.cpp Matrix.hpp
c++ -std=c++11 -c bench.cpp -o bench.o


Matrix.o: Matrix.cpp Matrix.hpp
c++ -std=c++11 -c Matrix.cpp -o Matrix.o
```

NORTHWEST INSTITUTE for ADVANCED COMPUTING

82

AMATH 483/583 High-Performance Scientific Computing Spring 2019
University of Washington by Andrew Lumsdaine

Pacific Northwest
NATIONAL LABORATORY
Proudly Operated by Battelle
for the U.S. Department of Energy

UNIVERSITY of
WASHINGTON

# Base Performance Results



Matrix Matrix Product Performance

# Let's Make One Small Change

```cpp
double benchmark(int M, int N, int K, long numruns) {
  Matrix A(M, K), B(K, N), C(M, N);

  Timer T;
  T.start();
  for (int i = 0; i < numruns; ++i) {
    multiply(A, B, C);
  }
  T.stop();

  return T.elapsed();
}
```

Tell the compiler to use optimization level 3

```
bench: bench.o Matrix.o
c++ -O3 -std=c++11 bench.o Matrix.o -o bench

bench.o: bench.cpp Matrix.hpp
c++ -O3 -std=c++11 -c bench.cpp -o bench.o

Matrix.o: Matrix.cpp Matrix.hpp
c++ -O3 -std=c++11 -c Matrix.cpp -o Matrix.o
```

NORTHWEST INSTITUTE for ADVANCED COMPUTING

84

AMATH 483/583 High-Performance Scientific Computing Spring 2019
University of Washington by Andrew Lumsdaine

Pacific Northwest
NATIONAL LABORATORY
Proudly Operated by Battelle
for the U.S. Department of Energy

UNIVERSITY of
WASHINGTON

# Base Performance Results



Matrix Matrix Product Performance

# The Three Most Important Requirements for HPC

- Locality
- Locality
- Locality

NORTHWEST INSTITUTE for ADVANCED COMPUTING

86

AMATH 483/583 High-Performance Scientific Computing Spring 2019
University of Washington by Andrew Lumsdaine

Pacific Northwest
NATIONAL LABORATORY

Proudly Operated by Battelle
for the U.S. Department of Energy

UNIVERSITY of
WASHINGTON

# Locality -> Performance

- Caches are much smaller than main memory. How do we decide what data to keep in cache to effect higher performance (more accesses)?
- **Temporal Locality:** if a program accesses a memory location, there is a much higher than random probability that the same location will be accessed again
  - Cache replacement policies attempt to keep cached elements in the cache for as long as possible
- **Spatial Locality:** if a program accesses a memory location, there is a much higher than random probability that nearby locations will also be accessed (soon)
  - Cache policies read contiguous chunks of data – a referenced element and its neighbors – not just single elements

# Improving Locality

- Consider each step of inner loop

```
for (int i = 0; i < M; ++i)
    for (int j = 0; j < N; ++j)
        for (int k = 0; k < K; ++k)
            C(i,j) += A(i,k) * B(k,j);
```

- Load `C(i,j)` into register
- Load `A(i,k)` into register
- Load `B(k,j)` into register
- Multiply
- Add
- Store `C(i,j)`

What can be reused?

- Four memory operations and two floating point operations per iteration
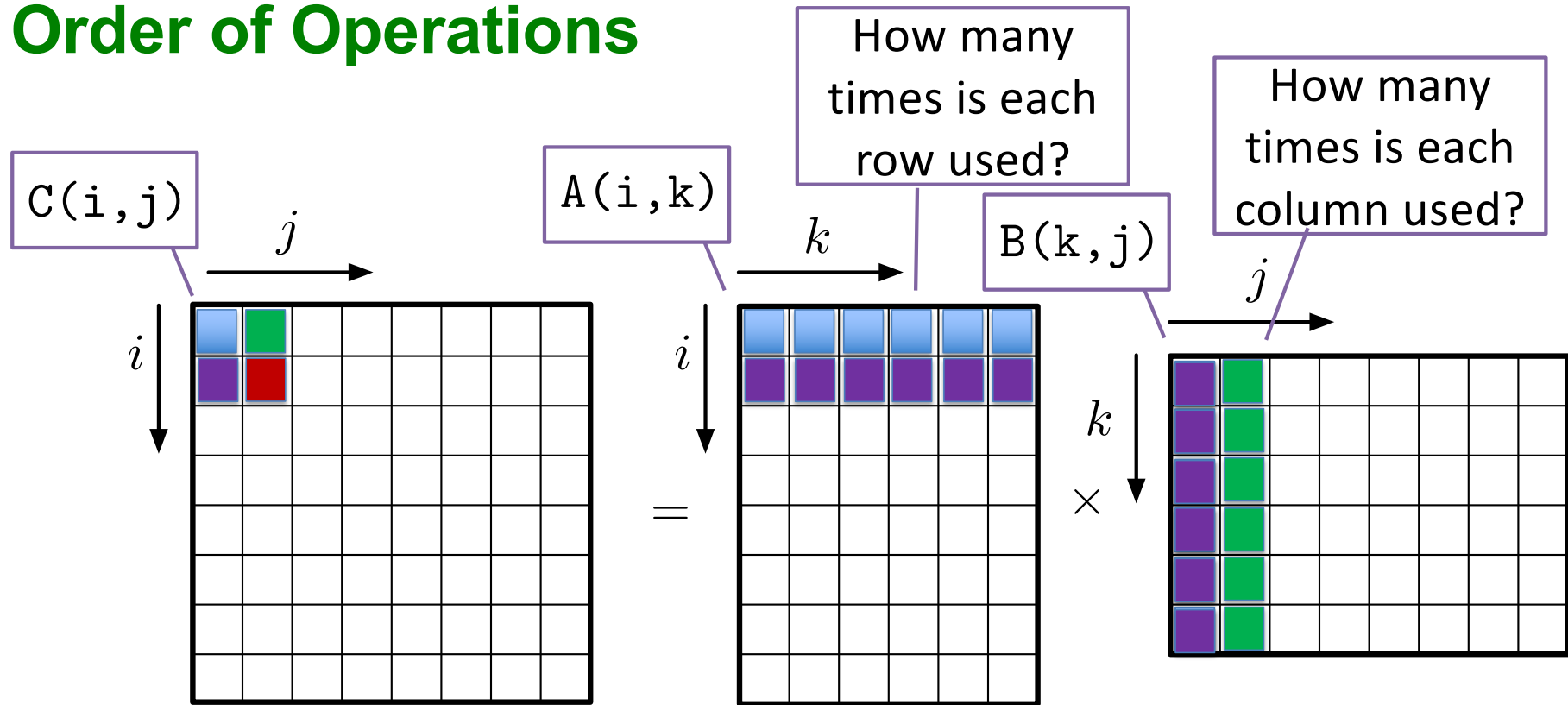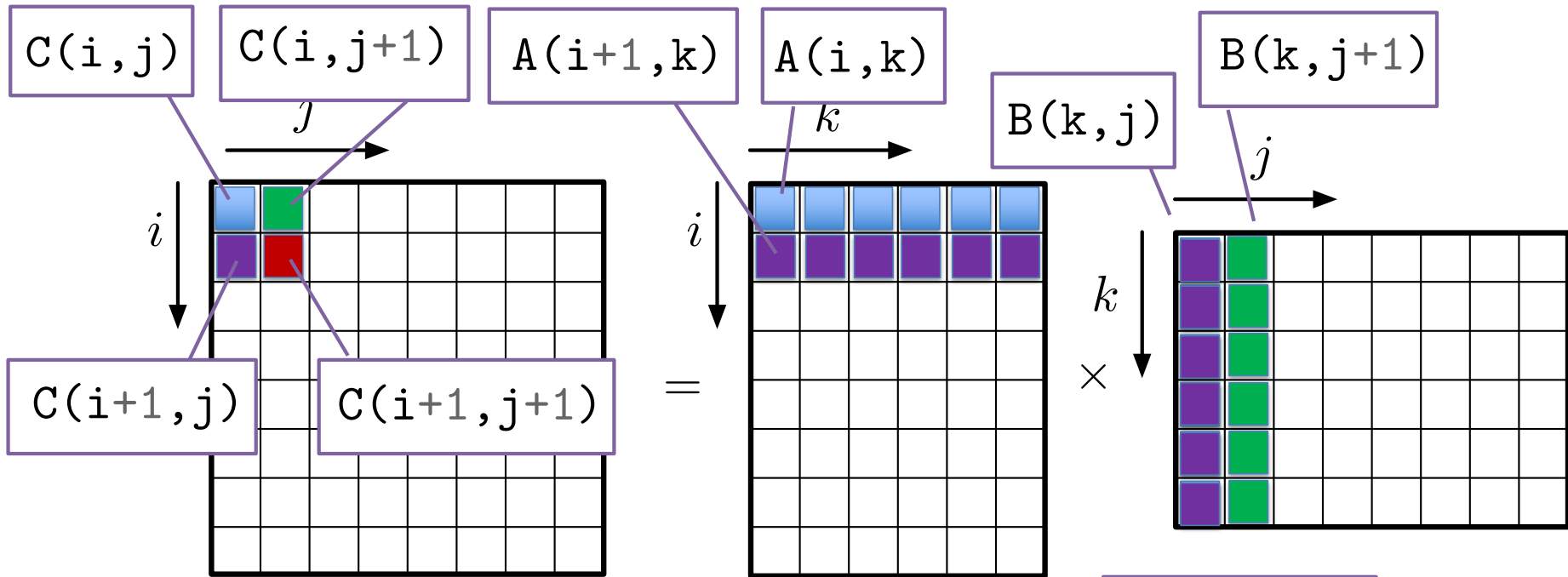- 1/3 flop per cycle (if each operation is one cycle)

# Improving Locality

```cpp
void multiply(const Matrix& A, const Matrix&B, Matrix&C) {
  for (size_t i = 0; i < A.num_rows(); ++i) {
    for (size_t j = 0; j < B.num_cols(); ++j) {
      for (size_ k = 0; k < A.num_cols(); ++k) {
        C(i,j) += A(i,k) * B(k,j);
      }
    }
  }
}
```

What can be reused?

- Load `C(i,j)` into register
- Load `A(i,k)` into register
- Load `B(k,j)` into register
- Multiply
- Add
- Store `C(i,j)`

- Four memory operations and two floating point operations per iteration
- 2/6 = 1/3 flop per cycle (if each operation is one cycle)

NORTHWEST INSTITUTE for ADVANCED COMPUTING

Pacific Northwest
NATIONAL LABORATORY
Proudly Operated by Battelle
for the U.S. Department of Energy

UNIVERSITY of
WASHINGTON

# Hoisting

Hoist C(i,j)

```
void multiply(const Matrix& A, const Matrix&B, Matrix&C) {
  for (size_t i = 0; i < A.num_rows(); ++i) {
    for (size_t j = 0; j < B.num_cols(); ++j) {
      double t = C(i,j);
      for (size_t k = 0; k < A.num_cols(); ++k) {
        t += A(i,k) * B(k,j);
      }
      C(i,j) = t;
    }
  }
}
```

- Load `A(i,k)`
- Load `B(k,j)`
- Multiply
- Add

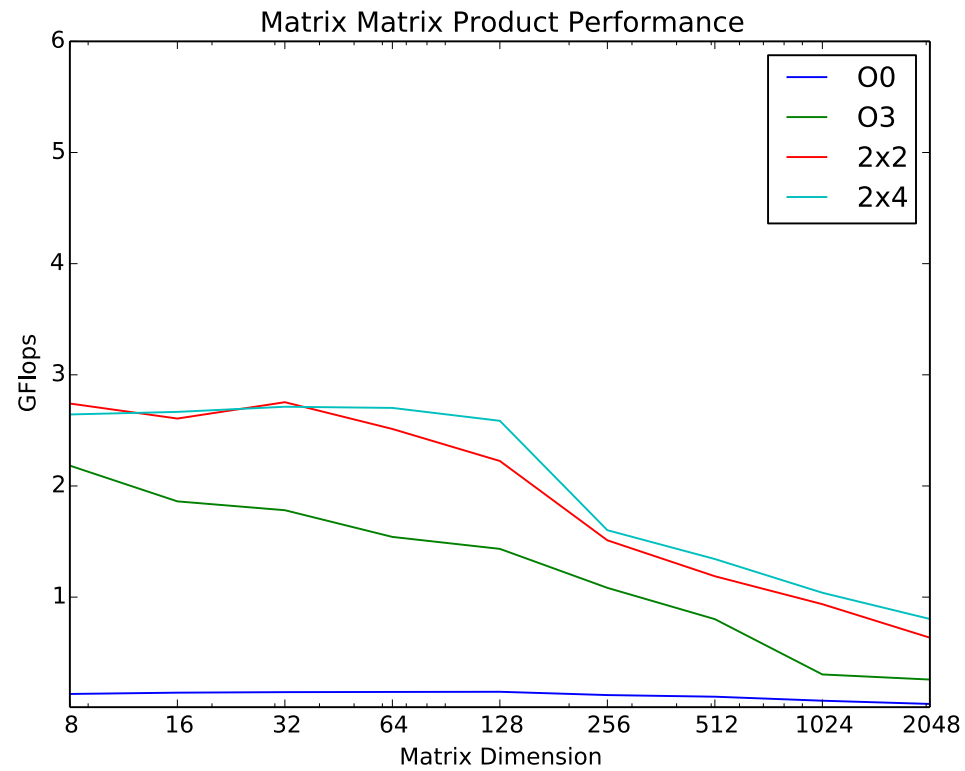- Two memory operations and two floating point operations per iteration
- 2/4 = 1/2 flop per cycle (if each operation is one cycle)

NORTHWEST INSTITUTE for ADVANCED COMPUTING

90

AMATH 483/583 High-Performance Scientific Computing Spring 2019
University of Washington by Andrew Lumsdaine

Pacific Northwest
NATIONAL LABORATORY
Proudly Operated by Battelle
for the U.S. Department of Energy

UNIVERSITY of
WASHINGTON

# Order of Operations



C(i,j)

A(i,k)

B(k,j)

NORTHWEST INSTITUTE for ADVANCED COMPUTING

Pacific Northwest
NATIONAL LABORATORY

Proudly Operated by Battelle
for the U.S. Department of Energy

UNIVERSITY of
WASHINGTON

# Order of Operations

C(i,j)

A(i,k)

B(k,j)

$j$

$i$

$j$

$k$

$i$

$k$

$=$

$\times$

NORTHWEST INSTITUTE for ADVANCED COMPUTING

92

AMATH 483/583 High-Performance Scientific Computing Spring 2019
University of Washington by Andrew Lumsdaine

Pacific Northwest
NATIONAL LABORATORY

Proudly Operated by Battelle
for the U.S. Department of Energy

UNIVERSITY of
WASHINGTON

# Order of Operations

C(i,j)

$j$

$i$

A(i,k)

How many times is each row used?

$i$

$k$

B(k,j)

How many times is each column used?

$j$

$k$

$=$   $\times$

NORTHWEST INSTITUTE for ADVANCED COMPUTING

93

AMATH 483/583 High-Performance Scientific Computing Spring 2019
University of Washington by Andrew Lumsdaine

Pacific Northwest
NATIONAL LABORATORY

Proudly Operated by Battelle
for the U.S. Department of Energy

UNIVERSITY of
WASHINGTON

# Reuse: How Many Times Are Data Reused?

C(i,j)  C(i,j+1)  A(i+1,k)  A(i,k)  B(k,j)  B(k,j+1)

C(i+1,j)  C(i+1,j+1)

=  ×

Each is used twice

# Improving Locality: Unroll and Ja...

```cpp
void tiledMultiply2x2(const Matrix& A, const Matrix&B
  for (size_t i = 0; i < A.num_rows(); i += 2) {
    for (size_t j = 0; j < B.num_cols(); j += 2) {
      for (size_t k = 0; k < A.num_cols(); ++k) {
        C(i  , j  ) += A(i  , k) * B(k, j  );
        C(i  , j+1) += A(i  , k) * B(k, j+1);
        C(i+1, j  ) += A(i+1, k) * B(k, j  );
        C(i+1, j+1) += A(i+1, k) * B(k, j+1);
      }
    }
  }
}
```

B(k,j) is used twice

B(k,j+1) is used twice

A(i,k) is used twice

A(i+1,k) is used twice

Can also hoist (independent of k)

- Four memory operations and eight floating point operations per iteration
- 8/12 = 2/3 flop per cycle (if each operation is one cycle) – 2X the base case

NORTHWEST INSTITUTE *for* ADVANCED COMPUTING

95

AMATH 483/583 High-Performance Scientific Computing Spring 2019
University of Washington by Andrew Lumsdaine

Pacific Northwest
NATIONAL LABORATORY
Proudly Operated by Battelle
for the U.S. Department of Energy

UNIVERSITY of WASHINGTON

# Example: Register Locality

# 2 by 4



Matrix Matrix Product Performance

AMATH 483/583 High-Performance Scientific Computing Spring 2019
University of Washington by Andrew Lumsdaine

NORTHWEST INSTITUTE for ADVANCED COMPUTING

# 4 by 2



Matrix Matrix Product Performance

NORTHWEST INSTITUTE for ADVANCED COMPUTING

# 4 by 4



Matrix Matrix Product Performance

NORTHWEST INSTITUTE for ADVANCED COMPUTING

# Tiling and Hoisting

```cpp
void hoistedTiledMultiply2x2(const Matrix& A, const Matrix&B, Matrix&C) {
  for (size_t i = 0; i < A.num_rows(); i += 2) {
    for (size_t j = 0; j < B.num_cols(); j += 2) {
      double t00 = C(i,  j);      double t01 = C(i,  j+1);
      double t10 = C(i+1,j);      double t11 = C(i+1,j+1);
      for (size_t k = 0; k < A.num_cols(); ++k) {
          t00 += A(i  , k) * B(k, j  );
          t01 += A(i  , k) * B(k, j+1);
          t10 += A(i+1, k) * B(k, j  );
          t11 += A(i+1, k) * B(k, j+1);
      }
      C(i,  j) = t00;  C(i,  j+1) = t01;
      C(i+1,j) = t10;  C(i+1,j+1) = t11;
    }
  }
}
```

Hoist 2x2 tile

AMATH 483/583 High-Performance Scientific Computing Spring 2019
University of Washington by Andrew Lumsdaine

# Tiling and Hoisting



Matrix Matrix Product Performance

NORTHWEST INSTITUTE for ADVANCED COMPUTING

# Improving Locality: Cache

- Large matrix problems won't fit completely into cache
- Use blocked algorithm – work with blocks that will fit into cache

$$C_{IJ} = \sum_K A_{IK} B_{KJ}$$

| $C_{00}$ | $C_{01}$ | $C_{02}$ | $C_{03}$ |
|---|---|---|---|
| $C_{10}$ | $C_{11}$ | $C_{12}$ | $C_{13}$ |
| $C_{20}$ | $C_{21}$ | $C_{22}$ | $C_{23}$ |
| $C_{30}$ | $C_{31}$ | $C_{32}$ | $C_{33}$ |

=

| $A_{00}$ | $A_{01}$ | $A_{02}$ | $A_{03}$ |
|---|---|---|---|
| $A_{10}$ | $A_{11}$ | $A_{12}$ | $A_{13}$ |
| $A_{20}$ | $A_{21}$ | $A_{22}$ | $A_{23}$ |
| $A_{30}$ | $A_{31}$ | $A_{32}$ | $A_{33}$ |

×

| $B_{00}$ | $B_{01}$ | $B_{02}$ | $B_{03}$ |
|---|---|---|---|
| $B_{10}$ | $B_{11}$ | $B_{12}$ | $B_{13}$ |
| $B_{20}$ | $B_{21}$ | $B_{22}$ | $B_{23}$ |
| $B_{30}$ | $B_{31}$ | $B_{32}$ | $B_{33}$ |

- Each product term fits completely into cache and runs at high-performance
- Cache misses amortized           work with           data

$$C_{21} = A_{20}B_{01} + A_{21}B_{11} + A_{22}B_{21} + A_{23}B_{31}$$

$$O(N^3) \qquad O(N^2)$$

NORTHWEST INSTITUTE for ADVANCED COMPUTING

102

Pacific Northwest
NATIONAL LABORATORY
Proudly Operated by Battelle
for the U.S. Department of Energy

UNIVERSITY of
WASHINGTON

# Blocking and Tiling

```cpp
void blockedTiledMultiply2x2(const Matrix& A, const Matrix&B, Matrix&C) {
  const int blocksize = std::min(A.num_rows(), 32);

  for (size_t ii = 0; ii < A.num_rows(); ii += blocksize) {
    for (size_t jj = 0; jj < B.num_cols(); jj += blocksize) {
      for (size_t kk = 0; kk < A.num_cols(); kk += blocksize) {

        for (size_t i = ii; i < ii+blocksize; i += 2) {
          for (size_t j = jj; j < jj+blocksize; j += 2) {
            for (size_t k = kk; k < kk+blocksize; ++k) {
              C(i  , j  ) += A(i  , k) * B(k, j  );
              C(i  , j+1) += A(i  , k) * B(k, j+1);
              C(i+1, j  ) += A(i+1, k) * B(k, j  );
              C(i+1, j+1) += A(i+1, k) * B(k, j+1);
            }
          }
        }
      }
    }
  }
}
```

Outer loops work
across blocks
(for each block)

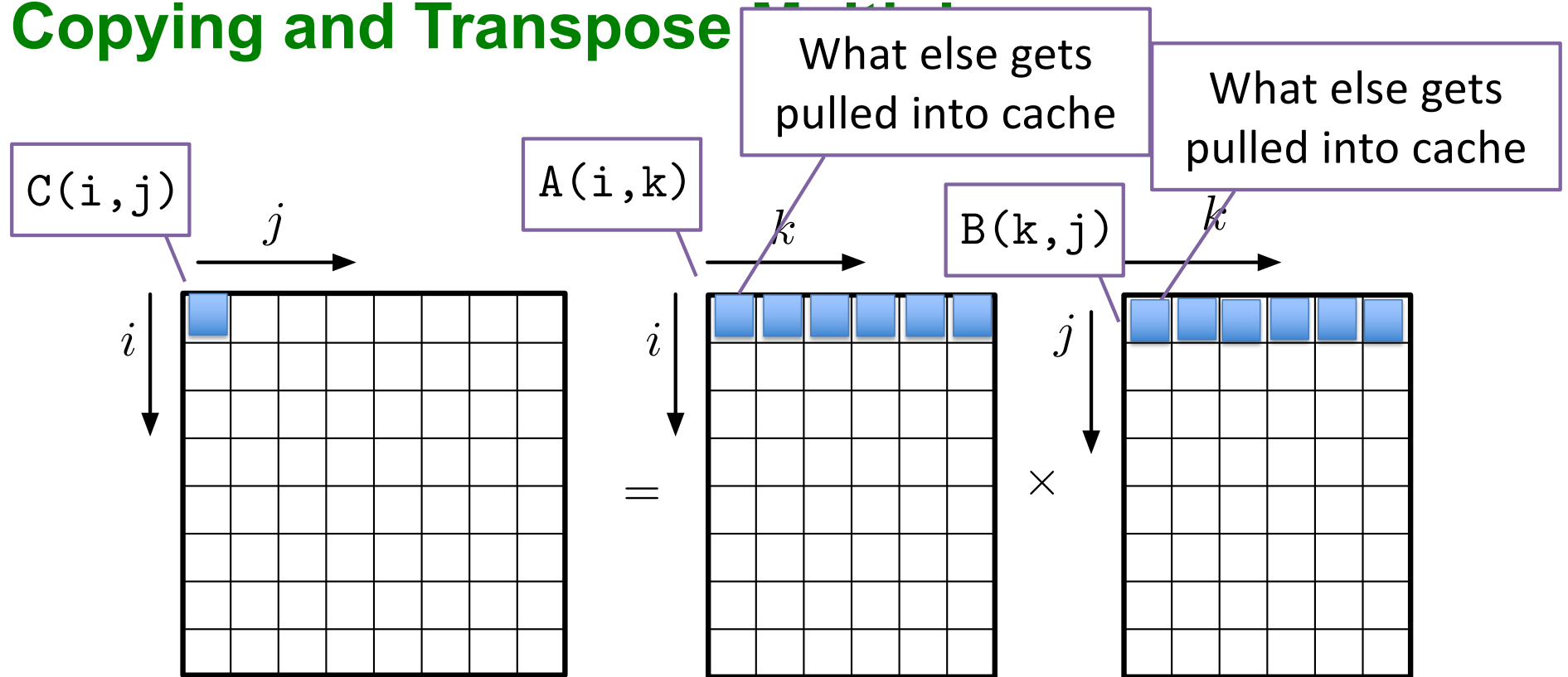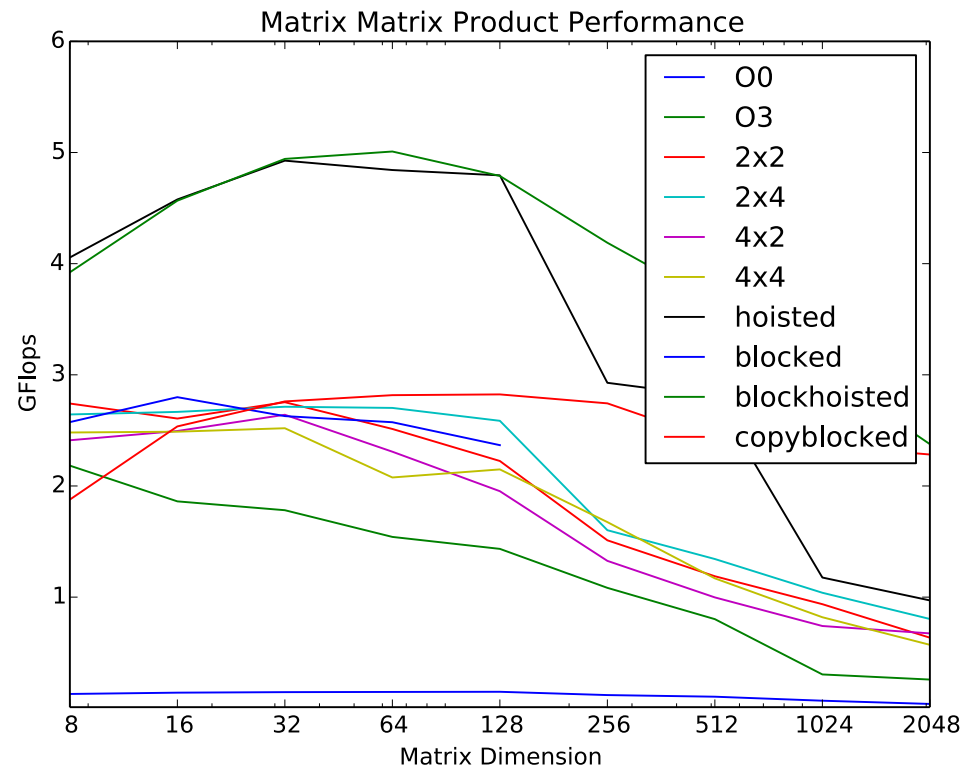Inner loops
work on blocks

# Blocking and Tiling and Hoisting



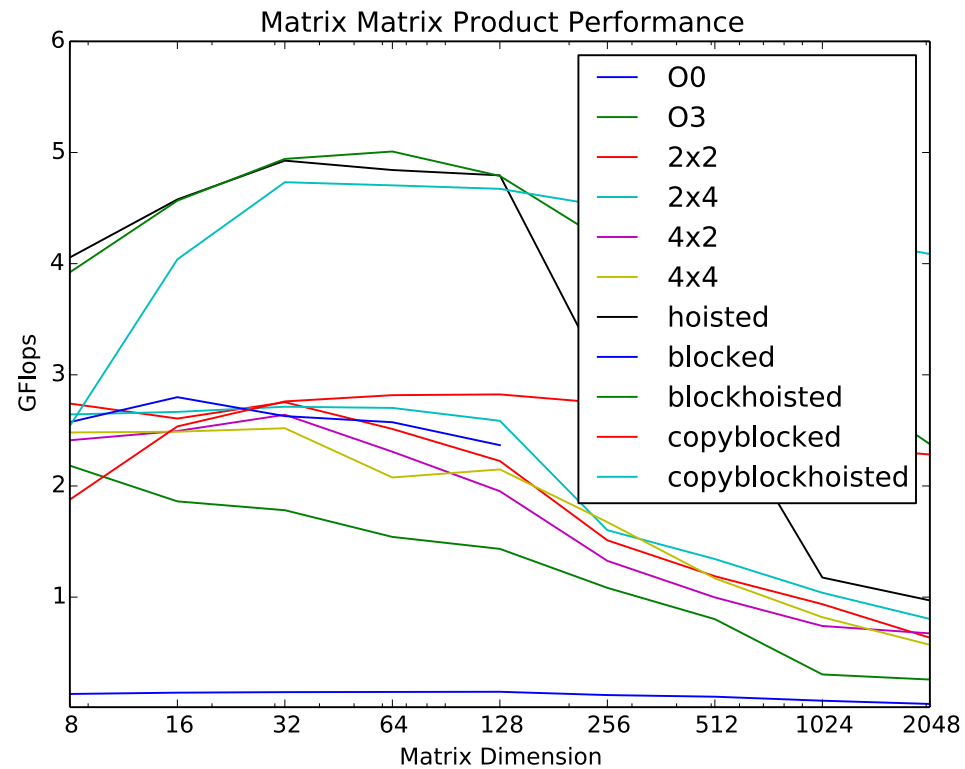Matrix Matrix Product Performance

# Copying

C(i,j)

A(i,k)

What else gets pulled into cache

B(k,j)

What else gets pulled into cache

$j$

$i$

$=$

$k$

$i$

$\times$

$k$

$j$

# Copying and Transpose Multiply

C(i,j)

A(i,k)

What else gets pulled into cache

B(k,j)

What else gets pulled into cache

NORTHWEST INSTITUTE for ADVANCED COMPUTING

106

AMATH 483/583 High-Performance Scientific Computing Spring 2019
University of Washington by Andrew Lumsdaine

Pacific Northwest
NATIONAL LABORATORY

Proudly Operated by Battelle
for the U.S. Department of Energy

UNIVERSITY of
WASHINGTON

# Copying and Blocking and Tiling



Matrix Matrix Product Performance

AMATH 483/583 High-Performance Scientific Computing Spring 2019
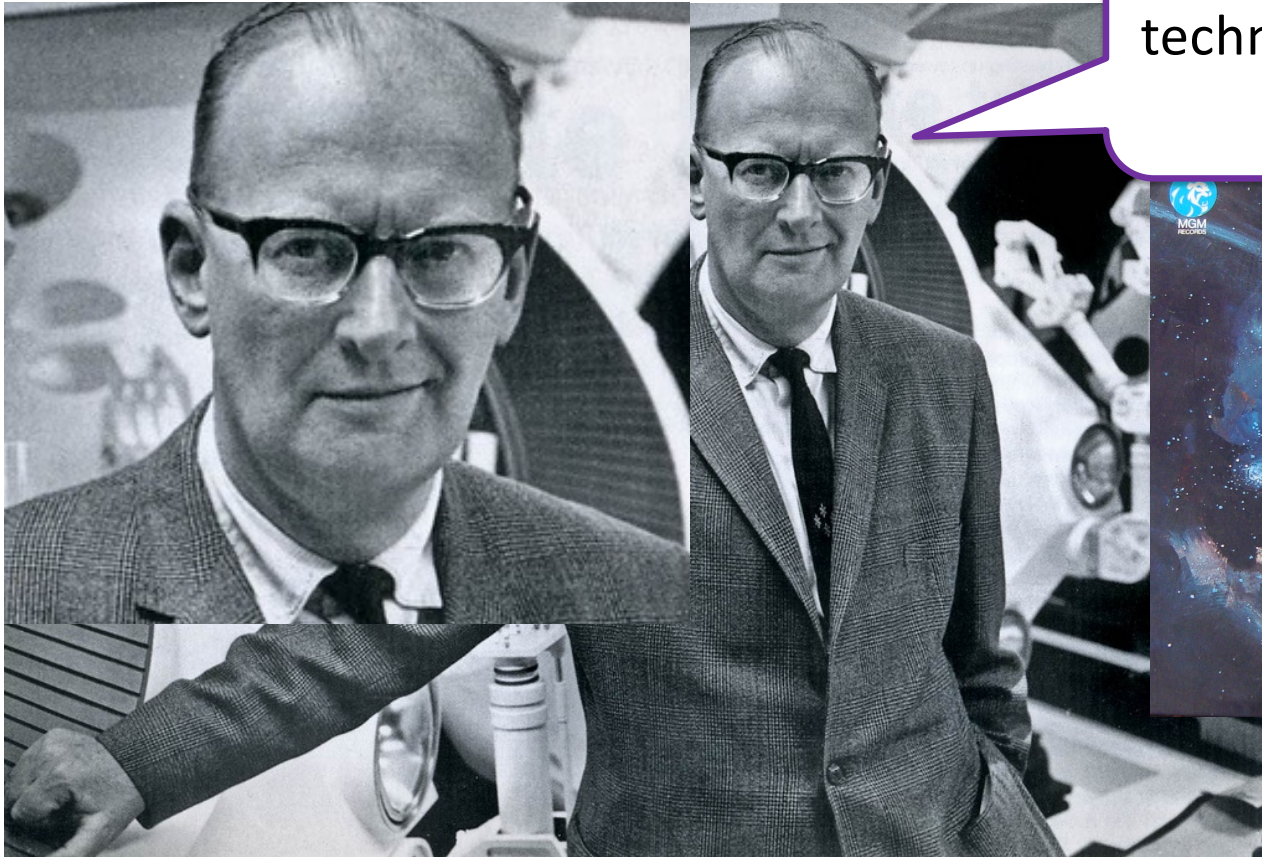University of Washington by Andrew Lumsdaine

# Blocking and Tiling and Hoisting and Copying

# Recap

- Locality: Write software so hardware can leverage it
- Register locality (tiling / unroll and jam)
- Hoisting
- Blocking
- Copying / transpose multiply
- Always use –O3 for release (not for debug)

# Name This Famous Person

Any sufficiently advanced technology is indistinguishable from magic

NORTHWEST INSTITUTE *for* ADVANCED COMPUTING

110

AMATH 483/583 High-Performance Scientific Computing Spring 2019
University of Washington by Andrew Lumsdaine

Pacific Northwest
NATIONAL LABORATORY
*Proudly Operated by Battelle*
*for the U.S. Department of Energy*

UNIVERSITY of WASHINGTON

# This Nearly Famous Person Says

Optimizing compilers are sufficiently advanced technology

And so are modern microprocessors

But especially optimizing compilers for modern microprocessors

Magic: the power of apparently influencing the course of events by using mysterious or supernatural forces

NORTHWEST INSTITUTE for ADVANCED COMPUTING

Pacific Northwest
NATIONAL LABORATORY
Proudly Operated by Battelle
for the U.S. Department of Energy

UNIVERSITY of WASHINGTON

111

AMATH 483/583 High-Performance Scientific Computing Spring 2019
University of Washington by Andrew Lumsdaine

# Tuning

- Starting with base code
- Various compiler optimizations help
- Tiling (which size)
- Blocking (what size)
- What size works best for Tiling and Blocking *together?*

- What loop ordering?  Matrix matrix product has six different orderings?  What block ordering?

- What about when we add AVX, and threads, etc?

> How do we find the optimal combination?

> Magic: the power of apparently influencing the course of events by using mysterious or supernatural forces

> The answer will be different for different CPUs

NORTHWEST INSTITUTE for ADVANCED COMPUTING

Pacific Northwest
NATIONAL LABORATORY
Proudly Operated by Battelle
for the U.S. Department of Energy

UNIVERSITY of
WASHINGTON

112

AMATH 483/583 High-Performance Scientific Computing Spring 2019
University of Washington by Andrew Lumsdaine

# Finding the Sweet Spot

- Exhaustive parameter space search
  - Tiling, Blocking, Compiler flags, AVX inst, loop ordering
- Original project at UC Berkeley phiPAC (Bilmes et al)
- Further developed by Whaley and Dongarra → Automatically Tuned Linear Algebra Subprograms (ATLAS)
  - Recently honored with "test of time" award

And wrote a program to generate different multiply functions

This started as a final course project

The competition was to write fastest matrix-matrix product
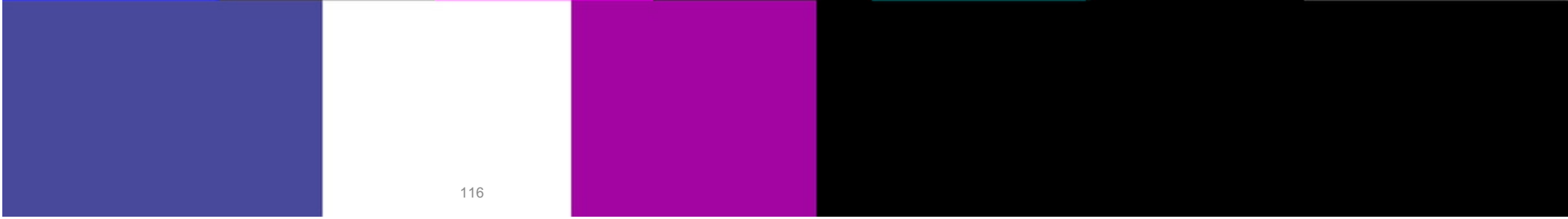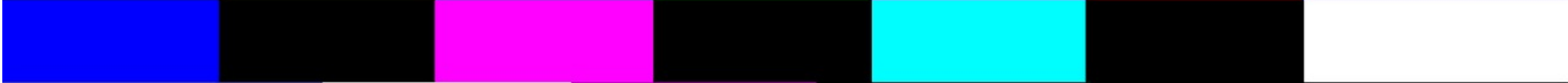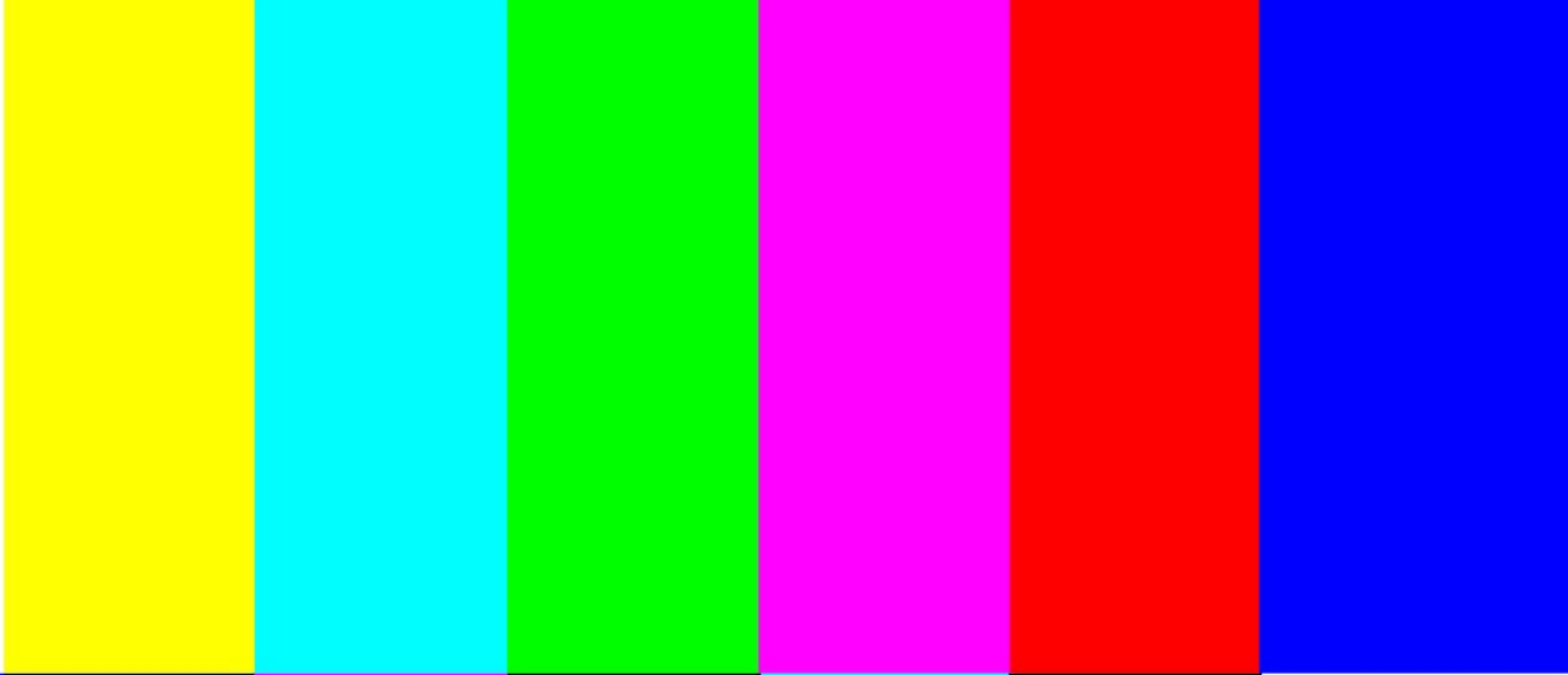
Students were the good kind of lazy

NORTHWEST INSTITUTE for ADVANCED COMPUTING

Pacific Northwest
NATIONAL LABORATORY
Proudly Operated by Battelle
for the U.S. Department of Energy

UNIVERSITY of WASHINGTON

113

AMATH 483/583 High-Performance Scientific Computing Spring 2019
University of Washington by Andrew Lumsdaine

# Thank you!

NORTHWEST INSTITUTE *for* ADVANCED COMPUTING

Pacific Northwest
NATIONAL LABORATORY

*Proudly Operated by* **Battelle**
*for the U.S. Department of Energy*

UNIVERSITY *of*
WASHINGTON

NORTHWEST INSTITUTE *for* ADVANCED COMPUTING

115

AMATH 483/583 High-Performance Scientific Computing Spring 2019
University of Washington by Andrew Lumsdaine

Pacific Northwest
NATIONAL LABORATORY
Proudly Operated by Battelle
for the U.S. Department of Energy

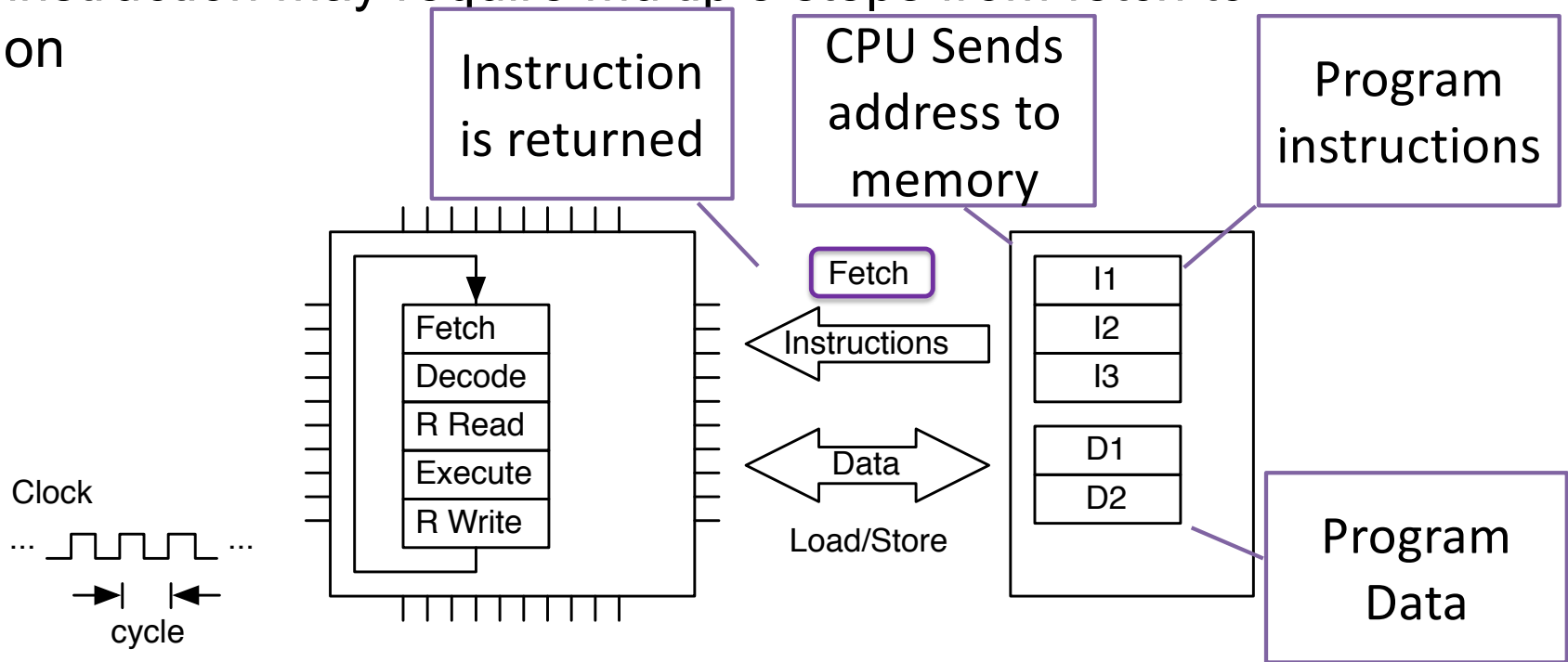UNIVERSITY of WASHINGTON

# Microprocessors

- Basic operation: read and execute program instructions stored in memory

- Fundamental performance / efficiency metric: cycles per instruction (CPI) also FLC



**Instructions can only be run in CPU**

**Transitions move data through CPU**

Clock

... ⊓⊔⊓ ...

cycle

CPU

**Data can only be operated on in CPU**

Fetch

Instructions

Data

Load/Store

Memory

**Program instructions and data**

# Processor Core Instruction Handling

- A single instruction may require multiple steps from fetch to completion



Instruction is returned

CPU Sends address to memory

Program instructions

Fetch

I1
I2
I3

Instructions

D1
D2

Data

Load/Store

Clock

... cycle ...

Fetch
Decode
R Read
Execute
R Write

Program Data

NORTHWEST INSTITUTE for ADVANCED COMPUTING

Pacific Northwest
NATIONAL LABORATORY
Proudly Operated by Battelle
for the U.S. Department of Energy

UNIVERSITY of WASHINGTON

# Processor Core Instruction Handling

- By pipelining, multiple instructions can be executed at each clock cycle
- Form of instruction-level parallelism