NORTHWEST INSTITUTE *for* ADVANCED COMPUTING

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

UNIVERSITY *of* WASHINGTON

# AMATH 483/583
# High Performance Scientific Computing

# Lecture 5:
# CPUs, hierarchical memory, matrices

Andrew Lumsdaine

Northwest Institute for Advanced Computing

Pacific Northwest National Laboratory

University of Washington

Seattle, WA

# Overview

- Classes, Vectors, const, overloading
- Tour of computer architecture
- Class Matrix
- Matrix matrix product

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

# SC'19 Student Cluster Competition Call-Out!

- Teams work with advisor and vendor to design and build a cutting-edge, commercially available cluster constrained by the 3000-watt power limit
- Cluster run a variety of HPC workflows, ranging from being limited by CPU performance to being memory bandwidth limited to I/O intensive
- Teams are comprised of six undergrad or high-school students plus advisor

https://sc19.supercomputing.org/program/studentssc/student-cluster-competition/

Team Meetings
Mondays 5:30PM-8:00PM

NORTHWEST INSTITUTE for ADVANCED COMPUTING

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

UNIVERSITY of WASHINGTON

# C++ Core Guidelines related to classes

- [C.1: Organize related data into structures (structs or classes)](#)
- [C.3: Represent the distinction between an interface and an implementation using a class](#)
- [C.4: Make a function a member only if it needs direct access to the representation of a class](#)
- [C.10: Prefer concrete types over class hierarchies](#)
- [C.11: Make concrete types regular](#)

NORTHWEST INSTITUTE *for* ADVANCED COMPUTING

http://isocpp.github.io/CppCoreGuidelines/CppCoreGuidelines

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

# Anatomy of a C++ class

Declares interface

Hides definition

```cpp
class Vector {
public:
  Vector(size_t M) : num_rows_(M), storage_(num_rows_) {}

  double& operator()(size_t i) { return storage_[i]; }

  size_t num_rows() const { return num_rows_; }

private:
  size_t              num_rows_;
  std::vector<double> storage_;
};
```

Public accessors

Private data

Maintain invariants

NORTHWEST INSTITUTE for ADVANCED COMPUTING

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

UNIVERSITY of WASHINGTON

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

# Anatomy of a C++ class

Declares interface

Hides

Encapsulation

```cpp
class Vector {
public:
  Vector(size_t M) : num_rows_(M), storage_(num_rows_) {}

  double& operator()(size_t i) { return storage_[i]; }

  size_t num_rows() const { return num_rows_; }

private:
  size_t              num_rows_;
  std::vector<double> storage_;
};
```

Public accessors

Private data

Maintain invariants

NORTHWEST INSTITUTE for ADVANCED COMPUTING

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

UNIVERSITY of WASHINGTON

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

# The Vector class so Far

- Encapsulates vector data
- Member data for dimensions (rows) and for storing elements
- Member function to get number of rows
- Separate interface and implementation via public / private

- Three more things:
  - How to bring a Vector into being ("constructors")
  - Function for getting vector data
  - Function for setting vector data

Revisit operator()

Also called function call operator

Can create *function objects*

NORTHWEST INSTITUTE for ADVANCED COMPUTING

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

UNIVERSITY of
WASHINGTON

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

# Constructors

- The C++ compiler "knows" about built-in types
- When a variable of a built-in type is declared, the compiler just needs to allocate space for it
- C++ classes are user-defined
- Compiler can do its best (default constructor), but usually we need to do more to create a well-defined object

- For example, a well-defined vector should be given its (positive) dimension **when it is created**. (And the data initialized.)

NORTHWEST INSTITUTE *for* ADVANCED COMPUTING

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

UNIVERSITY *of*
WASHINGTON

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

# Constructors

Built-in type, compiler allocates known amount of space

Default constructor is invoked when variable is declared with no arguments

```cpp
int x = 42;
```

Compiler creates x with *default constructor*

```cpp
Vector x;
```

In this case, the constructor that takes an integer argument

```cpp
Vector x(27);
```

Compiler creates x by making a call to a specific constructor

```cpp
std::cout << "x is " << x.num_rows() << " in length." << std::cout;
```

Create a Vector x with 27 elements

Because that is how we defined the constructor

NORTHWEST INSTITUTE for ADVANCED COMPUTING

NATIONAL LABORATORY

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

UNIVERSITY of WASHINGTON

# Declaring Constructors

```cpp
#include <vector>

class Vector {
public:
  Vector();
  Vector(size_t M);

  size_t num_rows() const { r

private:
  size_t                num_rows_;
  std::vector<double> storage_;
};
```

A constructor is defined using the name of the class

And then the arguments

Can be **overloaded** (different functions distinguished by argument types)

Where have we already seen overloading?

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

# Defining Constructors

```cpp
#include <vector>

class Vector {
public:
  Vector();
  Vector(size_t M);

  size_t num_rows() const  { return num_rows; }

private:
  size_t            num_rows_;
  std::vector<double> storage_;
};
```

```cpp
#include "Vector.hpp"

Vector::Vector(size_t M) {
  num_rows_ = M;
  storage_ = std::vector<double>(num_rows);
}


Vector::Vector() {
  num_rows_ = 1;
  storage_ = std::vector<double>(num_rows_);
}
```

# Defining Constructors

Vector.hpp

```cpp
#include <vector>

class Vector {
public:
  Vector() {
    num_rows_ = 1;
    storage_  = std::vector<double>(num_rows);
  }
  Vector(size_t M) {
    num_rows_ = M;
    storage_  = std::vector<double>(num_rows);
  }

  size_t num_rows() const { return num_rows; }

private:
  size_t              num_rows_;
  std::vector<double> storage_;
};
```

# Initialization

- We have said that variables should always be initialized
- Different syntaxes

```cpp
int a = 42;

int b = int(42);

int c(42);

int d = { 42 };

std::vector<double> x = std::vector<double>(27);

std::vector<double> y(27);
```

c(42)

y(27)

# Defining Constructors

Note order of initialization

Initialization syntax
Introduce with :
Construct data members

```cpp
#include <vector>

class Vector {
public:
  Vector(size_t M) : num_rows_(M), storage_(num_rows_) {}

  size_t num_rows() const { return num_rows_;

private:
  size_t                num_rows_;
  std::vector<double> storage_;
};
```

Omit default constructor (why?)

Note order of declaration

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

# Defining Constructors

Vector.hpp

```cpp
#include <vector>

class Vector {
public:
  Vector(size_t M) : num_rows_(M), storage_(num_rows_) {}

  size_t num_rows() const { return num_rows_; }

private:
  size_t              num_rows_;
  std::vector<double> storage_;
};
```

Initialization

Primordial

Object doesn't yet exist

Object exists

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

# What Should operator() return?

```cpp
class Vector
public:
    double& operator()(size_t i);

private:
    size_t               num_rows_;
    std::vector<double>  storage_;
};
```

Return a **reference** to internal member data

```cpp
Vector x(5);

double foo = x(3);
x(2) = 0.0;
```

Can assign to internal data through the reference

Can read from internal data through the reference

Vector x(5);

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

# All Together

```cpp
#include <vector>

class Vector {
public:
  Vector(size_t M) : num_rows_(M), storage_(num_rows_) {}

  double& operator()(size_t i) { return storage_[i]; }

  size_t num_rows() const { return num_rows_; }

private:
  size_t            num_rows_;
  std::vector<double> storage_;
};
```

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

# Reprise operator+()

```cpp
#include <vector>

class Vector {
public:
  Vector operator+(const Vector& y);

private:
  size_t                num_rows_;
  std::vector<double> storage_;
};
```

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

# Reprise operator+()

Does this need to be a member?

Data for z

```cpp
#include <vector>

class Vector {
public:
  Vector operator+(const Vector& y) {
    Vector z(num_rows_);
    for (size_t i = 0; i < num_rows_; ++i) {
      z.storage_[i] = storage_[i] + y.storage[i];
    }
  }

private:
  size_t              num_rows_;
  std::vector<double> storage_;
};
```

Data for "x"

Data for y

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

# All Together

Vector.hpp

```cpp
#include <vector>

class Vector {
public:
  Vector(size_t M) : num_rows_(M), storage_(num_rows_) {}

  double& operator()(size_t i) { return storage_[i]; }

  size_t num_rows() const { return num_rows_; }

private:
  size_t             num_rows_;
  std::vector<double> storage_;
};
```

Can access via operator()

Don't need access to internals

Amath583.cpp

```cpp
#include "Vector.hpp"

Vector operator+(const Vector& x, const Vector& y) {
  Vector z(x.num_rows());
  for (size_t i = 0; i < z.num_rows(); ++i) {
    z(i) = x(i) + y(i);
  }
}
```

Return a Vector

Take args by const reference

Nicely symmetric

NORTHWEST INSTITUTE for AD

# All Together

**Vector.hpp**

```cpp
#include <vector>

class Vector {
public:
  Vector(size_t M) : num_rows_(M), storage_(num_rows_) {}

  double& operator()(size_t i) { return storage_[i]; }

  size_t num_rows() const { return num_rows_; }

private:
  size_t            num_rows_;
  std::vector<double> storage_;
};
```

**Amath583.hpp**

```cpp
#include "Vector.hpp"

Vector operator+(const Vector& x, const Vector& y);
```

**Amath583.cpp**

```cpp
#include "Vector.hpp"
#include "amath583.hpp"

Vector operator+(const Vector& x, const Vector& y) {
  Vector z(x.num_rows());
  for (size_t i = 0; i < z.num_rows(); ++i) {
    z(i) = x(i) + y(i);
  }
}
```

# Not quite finished

```cpp
#include "Vector.hpp"

int main() {

  Vector x(100), y(100), z(100), w(100);

  z = x + y;

  return 0;
}
```

```
% c++ constness.cpp
constness.cpp:20:12: error: no matching function for call to object of type 'const Vector'
    z(i) = x(i) + y(i);
           ^
constness.cpp:7:11: note: candidate function not viable: 'this' argument has type
      'const Vector', but method is not marked const
  double& operator()(size_t i) { return storage_[i]; }
          ^
constness.cpp:20:19: error: no matching function for call to object of type 'const Vector'
    z(i) = x(i) + y(i);
                  ^
constness.cpp:7:11: note: candidate function not viable: 'this' argument has type
      'const Vector', but method is not marked const
  double& operator()(size_t i) { return storage_[i]; }
          ^
2 errors generated.
```

**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

# Constness

⚠️

```cpp
#include <vector>

class Vector {
public:
  Vector(size_t M) : num_rows_(M), storage_(num_rows_) {}

  double& operator()(size_t i) { return storage_[i]; }

  size_t num_rows() const { return num_rows_; }

private:
  size_t            num_rows_;
  std::vector<double> storage_;
};
```

**x and y are defined to be const**

**Amath583.hpp**

```cpp
#in

Vector operator+(const Vector& x, const Vector& y);
```

**"this" is not const**

**Amath583.cpp**

```cpp
#include "Vector.hpp"
#include "amath583.hpp"

Vector operator+(const Vector& x, const Vector& y) {
  Vector z(x.num_rows());
  for (size_t i = 0; i < z.num_rows(); ++i) {
    z(i) = x(i) + y(i);
  }
}
```

# Overloading

```cpp
void foo(size_t i) {
  std::cout << "foo(size_t i)" << std::endl;
}

void foo(double d) {
  std::cout << "foo(double d)" << std::endl;
}
```

Takes a size_t

Takes a double

```cpp
int main() {

  size_t a = 0;
  double b = 0.0;

  foo(a);
  foo(b);

  return 0;
}
```

```
% ./a.out
foo(size_t i)
foo(double d)
```

ting Spring 2019

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

UNIVERSITY of
WASHINGTON

# Overloading

```cpp
void foo(size_t i) {
  std::cout << "void foo(size_t i)" << std::endl;
}

size_t foo(size_t i) {
  std::cout << "size_t foo(size_t i)" << std::endl;
}
```

Returns void

Returns size_t

```
% c++ overload.cpp
overload.cpp:7:8: error: functions that differ only in their return type cannot be overloaded
size_t foo(size_t i) {
~~~~~~ ^
overload.cpp:3:6: note: previous definition is here
void foo(size_t i) {
~~~~ ^
```

```cpp
int main() {

    size_t a = 0;
    size_t b = 0;

    foo(a);
    double c = foo(a);

    return 0;
}
```

Have to pick the function then call it

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

UNIVERSITY of WASHINGTON

# No overloading on return values

```cpp
size_t foo(size_t i) {
  std::cout << "size_t foo(size_t i)" << std::

  return i;
}


int main() {

  size_t a = 0;

  foo(a);
  size_t b = foo(a);
  double c = foo(a);

  return 0;
}
```

What happens to the return value is not the concern of the function

Ignore return value

Assign to size_t

Assign to double

Pacific Northwest
NATIONAL LABORATORY

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

UNIVERSITY of
WASHINGTON

# Constness

```cpp
double parens(double& x, size_t i) {
  std::cout << "called non const parens" << std::endl;
  double y = x;
  // .. some things
  return y;
}
```

```cpp
int main() {

  double x = 5.0;
  double y = parens(x);

  const double z = 5.0;
  double w = parens(z);

  double a = parens(5.0);
  double b = parens(x + y);

  const double c = parens(x + y + z + 5.0);

  return 0;
}
```

x is a ref

```
c++ const3.cpp
const3.cpp:27:14: error: no matching function call to 'parens'
  double w = parens(z, 27);
             ^~~~~
const3.cpp:13:8: note: candidate function not viable: 1st argument ('const double') would lose const
       qualifier
double parens(double& x, size_t i) {
       ^
```

```
const3.cpp:29:14: error: no matching function for call to 'parens'
  double a = parens(5.0, 27);
             ^~~~~
const3.cpp:13:8: note: candidate function not viable: expects an l-value for 1st argument
double parens(double& x, size_t i) {
       ^
```

Not okay

```
const3.cpp:32:20: error: no matching function for call to 'parens'
  const double c = parens(x + y + 5.0, 27);
                   ^~~~~
const3.cpp:13:8: note: candidate function not viable: expects an l-value for 1st argument
double parens(double& x, size_t i) {
       ^
```

Not okay

Pacific Northwest
NATIONAL LABORATORY

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

UNIVERSITY of
WASHINGTON

# Constness

```cpp
double parens(const double& x, size_t i) {
  std::cout << "called const parens" << std::endl;
  double y = x;
  // .. some things
  return y;
}
```

x is a const ref

```cpp
int main() {

  double x = 5.0;
  double y = parens(x);

  const double z = 5.0;
  double w = parens(z);

  double a = parens(5.0);
  double b = parens(x + y);

  const double c = parens(x + y + z + 5.0);

  return 0;
}
```

okay

okay

okay

okay

```
./a.out
called const parens
called const parens
called const parens
called const parens
called const parens
```

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

UNIVERSITY of
WASHINGTON

# Conststness

x is a const ref

x is a ref

```cpp
double parens(const double& x, size_t i) {
  std::cout << "called const parens" << std::endl;
  double y = x;
  // .. some things
  return y;
}
```

```cpp
double parens(double& x, size_t i) {
  std::cout << "called non const parens" << std::endl;
  double y = x;
  // .. some things
  return y;
}
```

```cpp
int main() {

  double x = 5.0;
  double y = parens(x);

  const double z = 5.0;
  double w = parens(z);

  double a = parens(5.0);
  double b = parens(x + y);

  const double c = parens(x + y + z + 5.0);

  return 0;
}
```

x is lvalue

z marked const

5.0 is an rvalue

x + y is an rvalue

```
./a.out
called non const parens
called const parens
called const parens
called const parens
called const parens
```

Pacific Northwest
NATIONAL LABORATORY

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

UNIVERSITY of
WASHINGTON

# Why not always pass const reference?

```
double parens(const double& x, size_t i) {
  std::cout << "called const parens" << std::endl;
  double y = x;
  // .. some things
  return x;
}
```

Return double

```
int main() {
  double y = 0.5;
  double p = 3.14;

  double x = 5.0;
  parens(x, 27) = p;

  const double z = 5.0;
  parens(z, 27) = p;

  parens(5.0, 27) = p;
  parens(x + y, 27) = p;

  return 0;
}
```

```
c++ const4.cpp
const4.cpp:23:17: error: expression is not assignable
  parens(x, 27) = p;
  ~~~~~~~~~~~~~ ^
const4.cpp:26:17: error: expression is not assignable
  parens(z, 27) = p;
  ~~~~~~~~~~~~~ ^
const4.cpp:28:19: error: expression is not assignable
  parens(5.0, 27) = p;
  ~~~~~~~~~~~~~~~ ^
const4.cpp:29:21: error: expression is not assignable
  parens(x + y, 27) = p;
  ~~~~~~~~~~~~~~~~~ ^
```

UTING
High-Performance Scientific Computing Spring 2019
sity of Washington by Andrew Lumsdaine

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

UNIVERSITY of
WASHINGTON

# Before

```cpp
double parens(const double& x, size_t i) {
  std::cout << "called const parens" << std::endl;
  double y = x;
  // .. some things
  return x;
}
```

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 Spring 2019
University of Washington by Andrew Lumsdaine

# After

```cpp
double& parens(const double& x, size_t i) {
  std::cout << "called const parens" << std::endl;
  double y = x;
  // .. some things
  return x;
}
```

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*

**W** UNIVERSITY of
WASHINGTON

# Why not always pass const reference?

```cpp
double& parens(const double& x, size_t i) {
  std::cout << "called const parens" << std::endl;
  double y = x;
  // .. some things
  return x;
}
```

But x is const

Return ref to double

Can't return const

```cpp
int main() {
  double y = 0.5;
  double p = 3.14;

  double x = 5.0;
  parens(x, 27) = p;

  const double z = 5.0;
  parens(z, 27) = p;

  parens(5.0, 27) = p;
  parens(x + y, 27) = p;

  return 0;
}
```

```
c++ const5.cpp
const5.cpp:9:10: error: binding value of type 'const double' to reference to type 'double' drops
        'const' qualifier
  return x;
         ^
```

# Before

```cpp
double& parens(const double& x, size_t i) {
  std::cout << "called const parens" << std::endl;
  double y = x;
  // .. some things
  return x;
}
```

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

# After

```cpp
const double& parens(const double& x, size_t i) {
  std::cout << "called const parens" << std::endl;
  double y = x;
  // .. some things
  return x;
}
```

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

# Why not always pass const reference?

```cpp
const double& parens(const double& x, size_t i) {
  std::cout << "called const parens" << std::endl;
  double y = x;
  // .. some things
  return x;
}
```

```cpp
int main() {
  double y = 0.5;
  double p = 3.14;

  double x = 5.0;
  parens(x, 27) = p;

  const double z = 5.0;
  parens(z, 27) = p;

  parens(5.0, 27) = p;
  parens(x + y, 27) = p;

  return 0;
}
```

```
c++ const5.cpp
const5.cpp:26:17: error: cannot assign to return value because function 'parens' returns a const value
  parens(x, 27) = p;
  ~~~~~~~~~~~~~ ^
const5.cpp:5:7: note: function 'parens' which returns const-qualified type 'const double &' declared
      here
const double& parens(const double& x, size_t i) {
      ^~~~~~
const5.cpp:29:17: error: cannot assign to return value because function 'parens' returns a const value
  parens(z, 27) = p;
  ~~~~~~~~~~~~~ ^
const5.cpp:5:7: note: function 'parens' which returns const-qualified type 'const double &' declared
      here
const double& parens(const double& x, size_t i) {
      ^~~~~~
const5.cpp:31:19: error: cannot assign to return value because function 'parens' returns a const value
  parens(5.0, 27) = p;
  ~~~~~~~~~~~~~~~ ^
const5.cpp:5:7: note: function 'parens' which returns const-qualified type 'const double &' declared
      here
const double& parens(const double& x, size_t i) {
      ^~~~~~
const5.cpp:32:21: error: cannot assign to return value because function 'parens' returns a const value
  parens(x + y, 27) = p;
  ~~~~~~~~~~~~~~~~~ ^
const5.cpp:5:7: note: function 'parens' which returns const-qualified type 'const double &' declared
      here
const double& parens(const double& x, size_t i) {
      ^~~~~~~
```

# Before

```cpp
double& parens(const double& x, size_t i) {
  std::cout << "called const parens" << std::endl;
  double y = x;
  // .. some things
  return x;
}
```

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*

W
UNIVERSITY of
WASHINGTON

# After

```cpp
double& parens(double& x, size_t i) {
  std::cout << "called const parens" << std::endl;
  double y = x;
  // .. some things
  return x;
}
```

NORTHWEST INSTITUTE for ADVANCED COMPUTING

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

UNIVERSITY of
WASHINGTON

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

# How about no const at all?

```cpp
double& parens(double& x, size_t i) {
  std::cout << "called const parens" << std::endl;
  double y = x;
  // .. some things
  return x;
}
```

```cpp
int main() {
  double y = 0.5;
  double p = 3.14;

  double x = 5.0;
  parens(x, 27) = p;

  const double z = 5.0;
  parens(z, 27) = p;

  parens(5.0, 27) = p;
  parens(x + y, 27) = p;

  return 0;
}
```

```
c++ const5.cpp
const5.cpp:30:3: error: no matching function for call to 'parens'
  parens(z, 27) = p;
  ^~~~~~
const5.cpp:14:9: note: candidate function not viable: 1st argument ('const double') would lose const
        qualifier
double& parens(double& x, size_t i) {
        ^
const5.cpp:32:3: error: no matching function for call to 'parens'
  parens(5.0, 27) = p;
  ^~~~~~
const5.cpp:14:9: note: candidate function not viable: expects an l-value for 1st argument
double& parens(double& x, size_t i) {
        ^
const5.cpp:33:3: error: no matching function for call to 'parens'
  parens(x + y, 27) = p;
  ^~~~~~
const5.cpp:14:9: note: candidate function not viable: expects an l-value for 1st argument
double& parens(double& x, size_t i) {
        ^
```

UTING
High-Performance Scientific Computing Spring 2019
...sity of Washington by Andrew Lumsdaine

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

UNIVERSITY of
WASHINGTON

# How about no const at all?

```cpp
int main() {
    double y = 0.5;
    double p = 3.14;

    double x = 5.0;
    parens(x, 27) = p;

    const double z = 5.0;
    parens(z, 27) = p;

    parens(5.0, 27) = p;
    parens(x + y, 27) = p;

    return 0;
}
```

This makes sense

This **should** be an error

This **should** be an error

This **should** be an error

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

# More sensible

This makes sense

This makes sense

This makes sense

This makes sense

```cpp
int main() {
  double y = 0.5;
  double p = 3.14;

  double x = 5.0;
  parens(x, 27) = p;

  const double z = 5.0;
  double q = parens(z, 27);

  double r = parens(5.0, 27);
  double s = parens(x + y, 27);

  return 0;
}
```

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

# More sensible

```cpp
double& parens(double& x, size_t i) {
  std::cout << "called non const parens" << std::endl;
  double y = x;
  // .. some things
  return x;
}

int main() {
  double y = 0.5;
  double p = 3.14;

  double x = 5.0;
  parens(x, 27) = p;

  const double z = 5.0;
  double q = parens(z, 27);

  double r = parens(5.0, 27);
  double s = parens(x + y, 27);

  return 0;
}
```

```
c++ const6.cpp
const6.cpp:30:14: error: no matching function for call to 'parens'
  double q = parens(z, 27);
             ^~~~~~
const6.cpp:14:9: note: candidate function not viable: 1st argument ('const double') would lose const
      qualifier
double& parens(double& x, size_t i) {
        ^
const6.cpp:32:14: error: no matching function for call to 'parens'
  double r = parens(5.0, 27);
             ^~~~~~
const6.cpp:14:9: note: candidate function not viable: expects an l-value for 1st argument
double& parens(double& x, size_t i) {
        ^
const6.cpp:33:14: error: no matching function for call to 'parens'
  double s = parens(x + y, 27);
             ^~~~~~
const6.cpp:14:9: note: candidate function not viable: expects an l-value for 1st argument
double& parens(double& x, size_t i) {
        ^
```

Oops, need to be const

Going in circles?

# More sensible

```cpp
const double& parens(const double& x, size_t i) {
  std::cout << "called non const parens" << std::endl;
  double y = x;
  // .. some things
  return x;
}

int main() {
  double y = 0.5;
  double p = 3.14;

  double x = 5.0;
  parens(x, 27) = p;

  const double z = 5.0;
  double q = parens(z, 27);

  double r = parens(5.0, 27);
  double s = parens(x + y, 27);

  return 0;
}
```

```
c++ const6.cpp
const6.cpp:27:17: error: cannot assign to return value because function 'parens' returns a const value
  parens(x, 27) = p;
  ~~~~~~~~~~~~~ ^
const6.cpp:6:7: note: function 'parens' which returns const-qualified type 'const double &' declared
        here
const double& parens(const double& x, size_t i) {
     ^~~~~~~
```

Oops, need to be non  const

Going in circles?

# Overloading to the rescue

```cpp
const double& parens(const double& x, size      double& parens(double& x, size_t i) {
  std::cout << "called non const parens"          std::cout << "called non const parens" << std::endl;
  double y = x;                                    double y = x;
  // .. some things                                // .. some things
  return x;                                        return x;
}                                                }
```

const

const

Not const

Not const

```cpp
int main() {
  double y = 0.5;
  double p = 3.14;

  double x = 5.0;
  parens(x, 27) = p;

  const double z = 5.0;
  double q = parens(z, 27);

  double r = parens(5.0, 27);
  double s = parens(x + y, 27);

  return 0;
}
```

```
./a.out
called non const parens
called const parens
called const parens
called const parens
```

# What does this have to do with operator()

```
const double& parens(const double& x, siz
  std::cout << "called non const parens" <
  double y = x;
  // .. some things
  return x;
}
```

```
double& parens(double& x, size_t i) {
  std::cout << "called non const parens" << std::endl;
  double y = x;
  // .. some things
  return x;
}
```

const

Not const

const

Not const

```
class Vector {
public:
  Vector(size_t M) : num_rows_(M), storage_(num_rows_) {}

  double& operator()(size_t i) { return storage_[i]; }

private:
  size_t              num_rows_;
  std::vector<double> storage_;
};
```

Where is the const or non-const thing to overload on?

# What does this have to do with operator()

```cpp
const double& parens(const double& x, size  ...
  std::cout << "called non const parens" ...
  double y = x;
  // .. some things
  return x;
}
```

```cpp
double& parens(double& x, size_t i) {
  std::cout << "called non const parens" << std::endl;
  double y = x;
  // .. some things
  return x;
}
```

const

const

Not const

Not const

```cpp
class Vector {
public:
  Vector(size_t M) : num_rows_(M), storage_(num_rows_) {}

        double& operator()(size_t i) { return storage_[i]; }
  const double& operator()(size_t i) { return storage_[i]; }

private:
  si                    um_rows_;
  st                    torage_;
};
```

Only differing by return type

Where is the const or non-const thing to overload on?

# There is a secret argument

```cpp
const double& parens(const double& x, size_
  std::cout << "called non const parens" 
  double y = x;
  // .. some things
  return x;
}
```

```cpp
double& parens(double& x, size_t i) {
  std::cout << "called non const parens" << std::endl;
  double y = x;
  // .. some things
  return x;
}
```

const

const

Not const

Not const

```cpp
class Vector {
public:
  Vector(size_t M) : num_rows_(M), storage_(num_rows_) {}

    double& operator()(size_t i) { return storage_[i]; }
  const double& operator()(size_t i) { return storage_[i]; }



                        num_rows_;
  std::vector<double> storage_;
};
```

Called "this"

There is a secret argument

There is a secret argument

# There is a secret argument

```
const double& parens(const double& x, siz|        double& parens(double& x, size_t i) {
  std::cout << "called non const parens"          std::cout << "called non const parens" << std::endl;
  double y = x;                                    double y = x;
  // .. some things                                // .. some things
  return x;                                        return x;
}                                                }
```

const

Not const

const

Not const

```
class Vector {
public:
  Vector(size_t M) : num_rows_(M), storage_(num_rows_) {}

        double& operator()(Vector *this, size_t i) { return storage_[i]; }
  const double& operator()(Vector *this, size_t i) { return storage_[i]; }

private:
  size_t              num_rows_;
  std::vector<double> storage_;
};
```

How would we fix our const problem?

# Before

```cpp
class Vector {
public:
  Vector(size_t M) : num_rows_(M), storage_(num_rows_) {}

        double& operator()(Vector *this, size_t i) { return storage_[i]; }
  const double& operator()(Vector *this, size_t i) { return storage_[i]; }

private:
  size_t                num_rows_;
  std::vector<double> storage_;
};
```

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

# After

```cpp
class Vector {
public:
  Vector(size_t M) : num_rows_(M), storage_(num_rows_) {}

        double& operator()(Vector *this, size_t i) { return storage_[i]; }
  const double& operator()(const Vector *this, size_t i) { return storage_[i]; }

private:
  size_t              num_rows_;
  std::vector<double> storage_;
};
```

const "this"

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

# After After

```
class Vector {
public:
  Vector(size_t M) : num_rows_(M), storage_(num_rows_) {}

        double& operator()(size_t i)       { return storage_[i]; }
  const double& operator()(size_t i) const { return storage_[i]; }

private:
  size_t              num_rows_;
  std::vector<double> storage_;
};
```

const "this"

NORTHWEST INSTITUTE for ADVANCED COMPUTING

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

UNIVERSITY of
WASHINGTON

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

# Matrix Representation

- Two issues
  - Interface (what is the abstraction we want to present?)
  - Implementation (how is the abstraction realized?)
- Sometimes there are tradeoffs
  - Evaluate relative to end user
  - In HPC – performance is most important
  - Elsewhere – safety, ease of use, standards compliance, etc

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

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

# Matrix Representation

Notionally 2Dimensional

Use a 2D data structure

$$C_{ij} = \sum_k A_{ik} B_{kj}$$

$M$

$C$

$N$

$M$

$A$

$K$

$B$

$K$

$N$

NORTHWEST INSTITUTE for ADVANCED COMPUTING

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

UNIVERSITY of
WASHINGTON

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

# Matrix Representation

$$C_{ij} = \sum_k A_{ik} B_{kj}$$

$i \longrightarrow$

$j \downarrow$

$\begin{bmatrix} A_{ij} & \\ & \end{bmatrix}$

Use a doubly indexed data structure

column $\longrightarrow$

row $\downarrow$

A matrix is a doubly indexed set

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

Use a doubly indexed accessor

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

# Matrix Representation

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

column

row

Use a doubly indexed accessor

To create one memory address

addr

CPU memory is a linear address space

NORTHWEST INSTITUTE *for* ADVANCED COMPUTING

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

UNIVERSITY *of*
WASHINGTON

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

# Matrix Representation

- To translate double index to single address

```
double **storage_;
```

Array of arrays

```
storage_[i][j]
```

Use inner pointer to get data element

Lookup inner pointer from outer pointer

```
std::vector<std::vector<double> > storage_;    storage_[i][j]
```

Vector of vectors

Use inner vector to get data element

Lookup inner vector from outer vector

```
std::vector<double> storage_;                    storage_[k]
```

Use single vector to get data element

Need to compute this

# Matrix Representation

column

row

addr

Number of columns

Number of columns

NORTHWEST INSTITUTE for ADVANCED COMPUTING

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

UNIVERSITY of
WASHINGTON

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

# Matrix Representation

First element

column

Second element

addr

row

First element

Second element

Number of columns

NORTHWEST INSTITUTE for ADVANCED COMPUTING

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

UNIVERSITY of
WASHINGTON

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

# Matrix Representation

C(0,0)

addr ↓

C(0,1)

column →

row ↓

C(1,0)

vC[0]

vC[1]

vC[N]

```
Matrix C(M, N);

vector<double> vC(M*N);
```

Number of columns (N)

NORTHWEST INSTITUTE for ADVANCED COMPUTING

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

UNIVERSITY of
WASHINGTON

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

# Matrix Representation



**Row-Oriented**

column

row

C(i,j)

addr

vC[i*N + j]

```
Matrix C(M, N);

vector<double> vC(M*N);
```

Increment fastest along row

Number of columns (N)

NORTHWEST INSTITUTE for ADVANCED COMPUTING

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

UNIVERSITY of
WASHINGTON

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

# Matrix Representation

**column** →

row ↓

C(i,j)

addr ↓

**Column-Oriented**

vC[j*M + i]

```
Matrix C(M, N);

vector<double> vC(M*N);
```

Increment fastest along column

Number of rows (M)

NORTHWEST INSTITUTE for ADVANCED COMPUTING

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

UNIVERSITY of
WASHINGTON

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

# Matrix Implementation

**Write this** → **Do this**

column

row

addr

C(i,j)

C(i,j) → vC[i*N + j]

vC[i*N + j]

```
Matrix C(M, N);

vector<double> vC(M*N);
```

Increment fastest along row

Number of columns (N)

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

# Matrix Implementation

Fortran did this transparently

column →

row ↓

C(i,j)

addr ↓

C(i,j)        vC[i*N + j]

vC[i*N + j]

```
Matrix C(M, N);

vector<double> vC(M*N);
```

Increment fastest along row

Number of columns (N)

NORTHWEST INSTITUTE for ADVANCED COMPUTING

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

UNIVERSITY of
WASHINGTON

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

# Matrix in C++

- Two dimensional accessor  `C(i,j)`
- One dimensional access  `vC[i*N + j]`

- Simultaneously (and safely) need matrix with
  - (i,j) two dimensional accessor
  - Transparent translation to one dimensional accessor

- Preprocessor?  `#define C(i,j) vC[i*N+j]`

Only works for C and vC

Where does N come from

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

# Matrix in C++

- Two dimensional accessor `C(i,j)`
- One dimensional access `vC[i*N + j]`

<br>

- A `Matrix` needs to
  - Have its "own" `vector<double> vC`
  - Have its "own" `N`
  - Have a doubly indexed accessor

NORTHWEST INSTITUTE for ADVANCED COMPUTING

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

UNIVERSITY of
WASHINGTON

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

# Class Matrix

```cpp
#include <vector>

class Matrix {
public:
  Matrix(size_type M, size_type N)
      : num_rows_(M), num_cols_(N), storage_(num_rows_ * num_cols_) {}
  Matrix(size_type M, size_type N, double init)
      : num_rows_(M), num_cols_(N), storage_(num_rows_ * num_cols_, init) {}

  double &operator()(size_type i, size_type j) {
    return storage_[i * num_cols_ + j];
  }
  const double &operator()(size_type i, size_type j) const {
    return storage_[i * num_cols_ + j];
  }

  size_type num_rows() const { return num_rows_; }
  size_type num_cols() const { return num_cols_; }

private:
  size_type num_rows_, num_cols_;
  std::vector<double> storage_;
};
```

# Class Matrix

```cpp
class Matrix {
public:
  Matrix(size_type M, size_type N)
      : num_rows_(M), num_cols_(N), storage_(num_rows_ * num_cols_) {}

  Matrix(size_type M, size_type N, double init)
      : num_rows_(M), num_cols_(N), storage_(num_rows_ * num_cols_, init) {}

private:
  size_type num_rows_, num_cols_;
  std::vector<double> storage_;
};
```

NORTHWEST INSTITUTE for ADVANCED COMPUTING

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

UNIVERSITY of
WASHINGTON

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

# Class Matrix

```cpp
class Matrix {
public:
  double &operator()(size_type i, size_type j) {
    return storage_[i * num_cols_ + j];
  }



  size_type num_rows() const { return num_rows_; }
  size_type num_Cols() const { return num_cols_; }

private:
  size_type num_rows_, num_cols_;
  std::vector<double> storage_;
};
```

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

# Class Matrix

```cpp
class Matrix {
public:
  double &operator()(size_type i, size_type j) {
    return storage_[i * num_cols_ + j];
  }
  const double &operator()(size_type i, size_type j) const {
    return storage_[i * num_cols_ + j];
  }

  size_type num_rows() const { return num_rows_; }
  size_type num_Cols() const { return num_cols_; }

private:
  size_type num_rows_, num_cols_;
  std::vector<double> storage_;
};
```

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

# Class Matrix

```cpp
class Matrix {
public:
  double &operator()(size_type i, size_type j) {
    return storage_[i * num_cols_ + j];
  }
  const double &operator()(size_type i, size_type j) const {
    return storage_[i * num_cols_ + j];
  }

  size_type num_rows() const { return num_rows_, }
  size_type num_Cols() const { return num_cols_; }

private:
  size_type num_rows_, num_cols_;
  std::vector<double> storage_;
};
```

How would we write the other orientation?

What is the orientation?

Does it matter which?

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

# Before

```cpp
class Matrix {
public:
  double &operator()(size_type i, size_type j) {
    return storage_[i * num_cols_ + j];
  }
  const double &operator()(size_type i, size_type j) const {
    return storage_[i * num_cols_ + j];
  }

  size_type num_rows() const { return num_rows_; }
  size_type num_Cols() const { return num_cols_; }

private:
  size_type num_rows_, num_cols_;
  std::vector<double> storage_;
};
```

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

# After

```cpp
class Matrix {
public:
  double &operator()(size_type i, size_type j) {
    return storage_[j * num_rows_ + i];
  }
  const double &operator()(size_type i, size_type j) const {
    return storage_[j * num_rows_ + i];
  }

  size_type num_rows() const { return num_rows_; }
  size_type num_Cols() const { return num_cols_; }

private:
  size_type num_rows_, num_cols_;
  std::vector<double> storage_;
};
```

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

# Finally

```cpp
#include <vector>

class Vector {
public:
  Vector(size_t M) : num_rows_(M), storage_(num_rows_) {}

        double& operator()(size_t i)       { return storage_[i]; }
  const double& operator()(size_t i) const { return storage_[i]; }

  size_t num_rows() { return num_rows_; }

private:
  size_t            num_rows_;
  std::vector<double> storage_;
};
```

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

# Example: Matrix-Matrix Product

- For matrices $\mathbf{A} \in \mathbb{R}^{M \times K}$ and $\mathbf{B} \in \mathbb{R}^{K \times N}$, compute $\mathbf{C} \in \mathbb{R}^{M \times N}$

$$\mathbf{C} \leftarrow \mathbf{A} \times \mathbf{B}$$ Defined according to $C_{ij} = \sum_k A_{ik} B_{kj}$

Locality!
(Data reuse)

- Workhorse computational kernel (underlying LINPACK, HPL)
- Compute-intensive: $O(N^3)$ work with $O(N^2)$ data

Every element is accessed N times

- Basic algorithm in C++

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

Maximize
locality

NORTHWEST INSTITUTE for ADVANCED COMPUTING

74

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 FLO

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

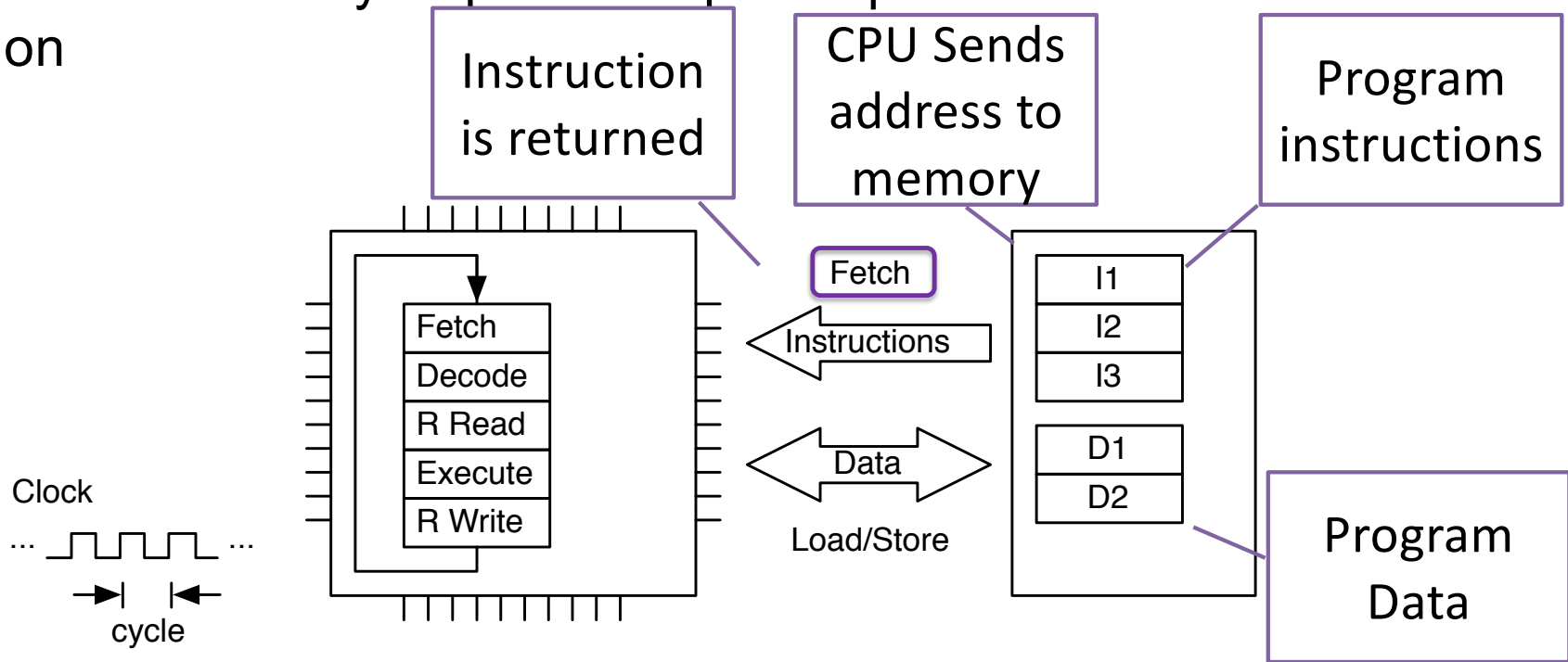Memory

Program instructions and data

# 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

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

# 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

Instructions

I1

I2

I3

Fetch
Decode
R Read
Execute
R Write

Data

D1

D2

Clock

... ⊓_⊓_⊓ ...

→| |← cycle

Load/Store

Program Data

NORTHWEST INSTITUTE for ADVANCED COMPUTING

77

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
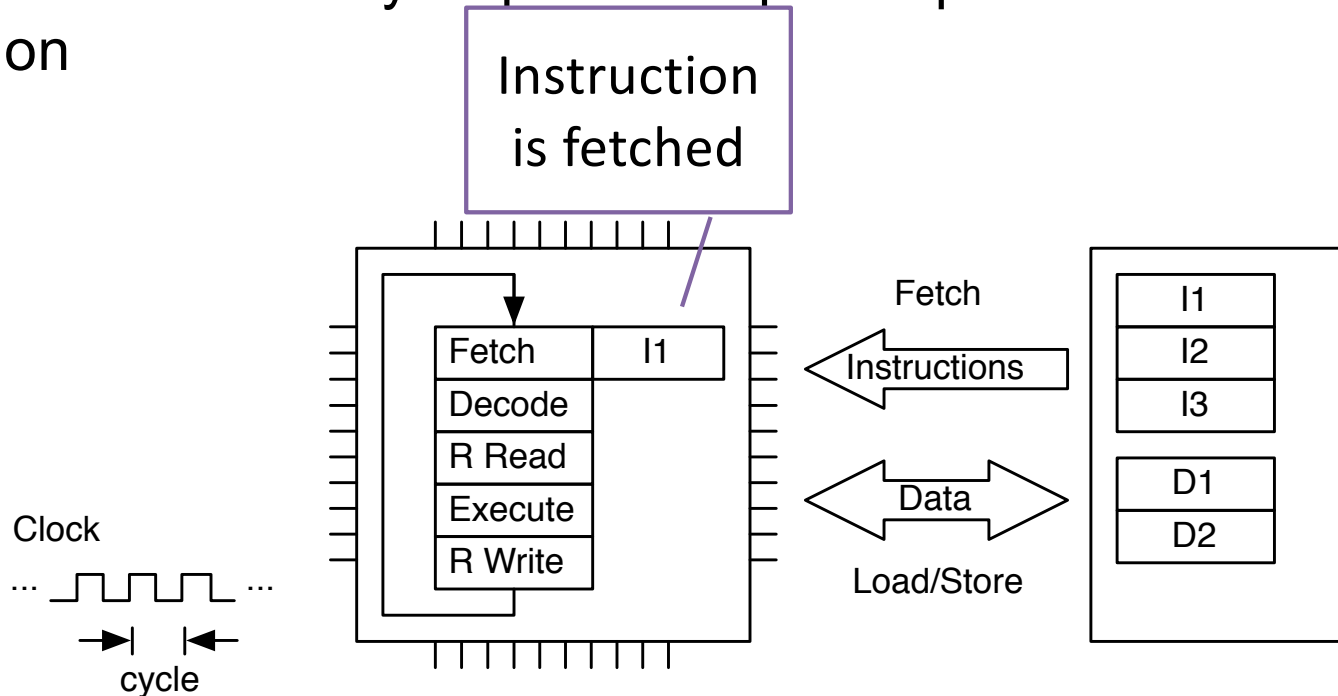
UNIVERSITY of
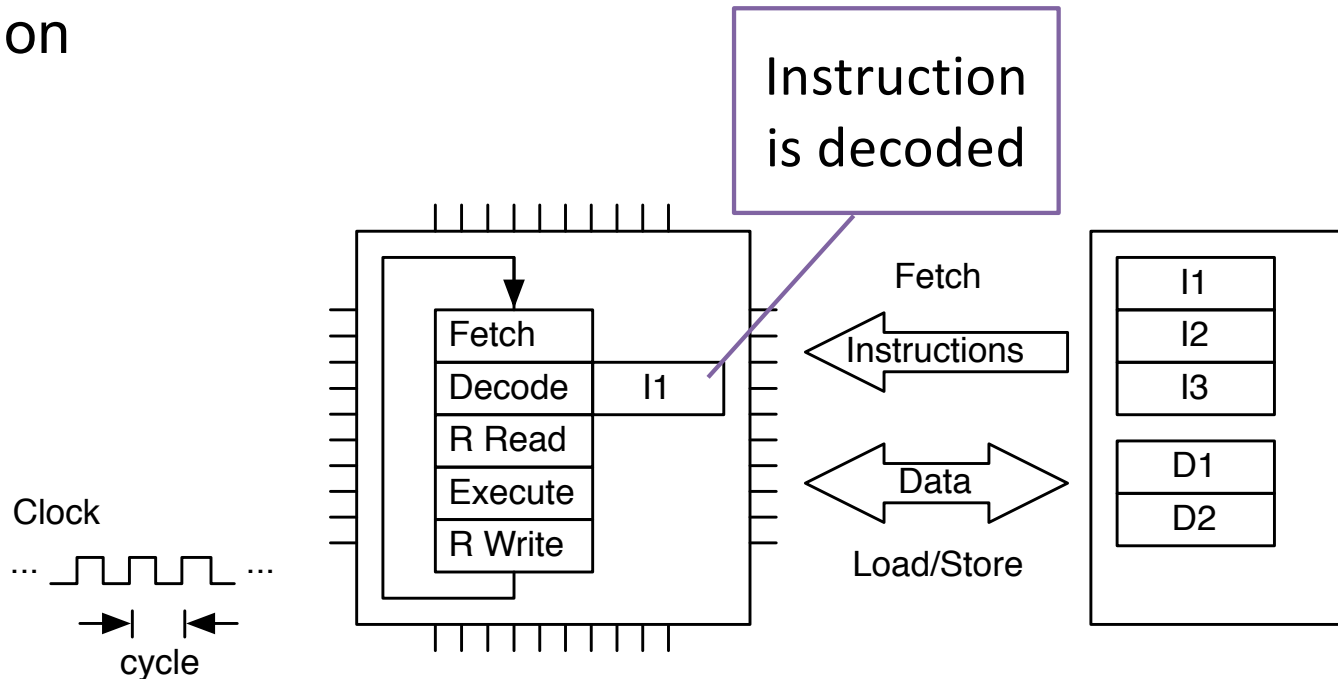WASHINGTON

# 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 decoded

Fetch

Instructions

I1

Fetch
Decode    I1
R Read
Execute
R Write

Data

Load/Store

Clock

...⎍⎍⎍...

cycle

I1
I2
I3

D1
D2

# Processor Core Instruction Handling

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

Registers are read

Fetch

Fetch
Decode
R Read    I1
Execute
R Write

Instructions

Data

Load/Store

I1
I2
I3

D1
D2

Clock

... ⊓⊔⊓⊔⊓⊔ ...

cycle

NORTHWEST INSTITUTE for ADVANCED COMPUTING

80

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

# Processor Core Instruction Handling

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



Clock

... ⊓⊔⊓⊔⊓ ...

cycle

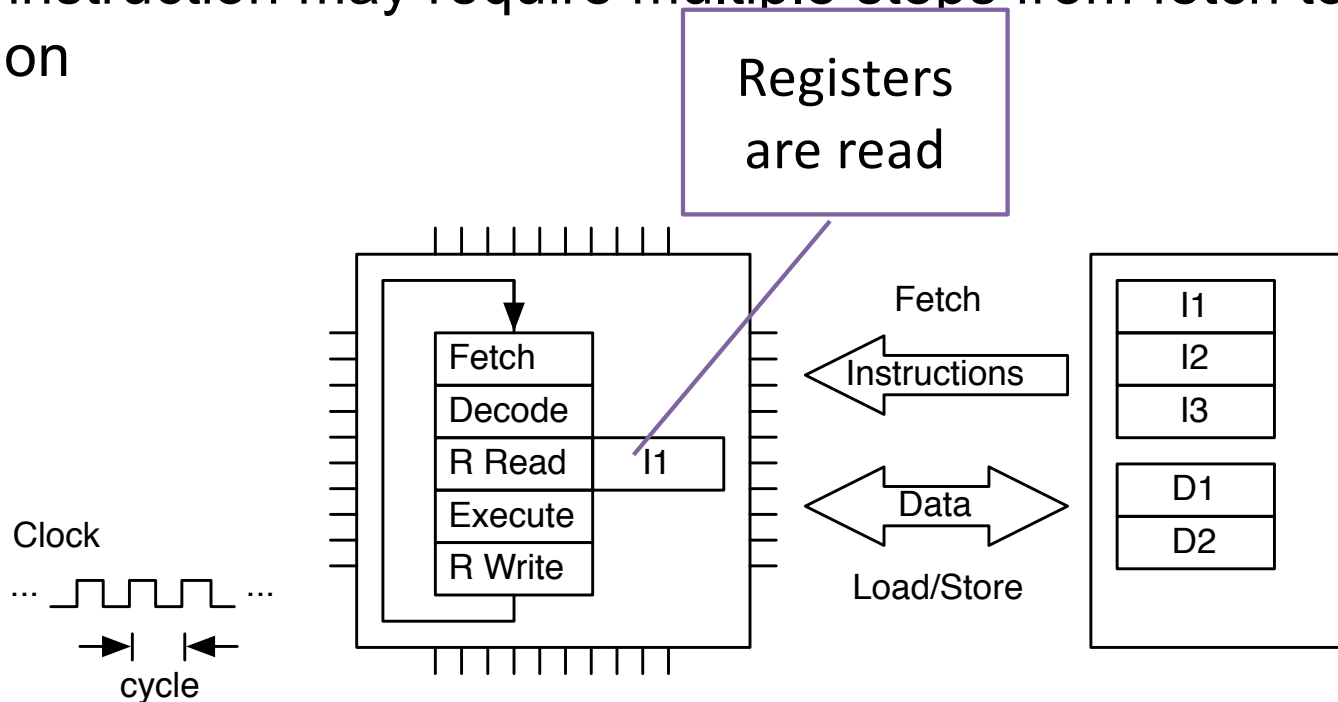Fetch
Decode
R Read
Execute    I1
R Write

Fetch
Instructions

Data

Load/Store

I1
I2
I3

D1
D2

Instruction is executed

NORTHWEST INSTITUTE for ADVANCED COMPUTING

81

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

# Processor Core Instruction Handling

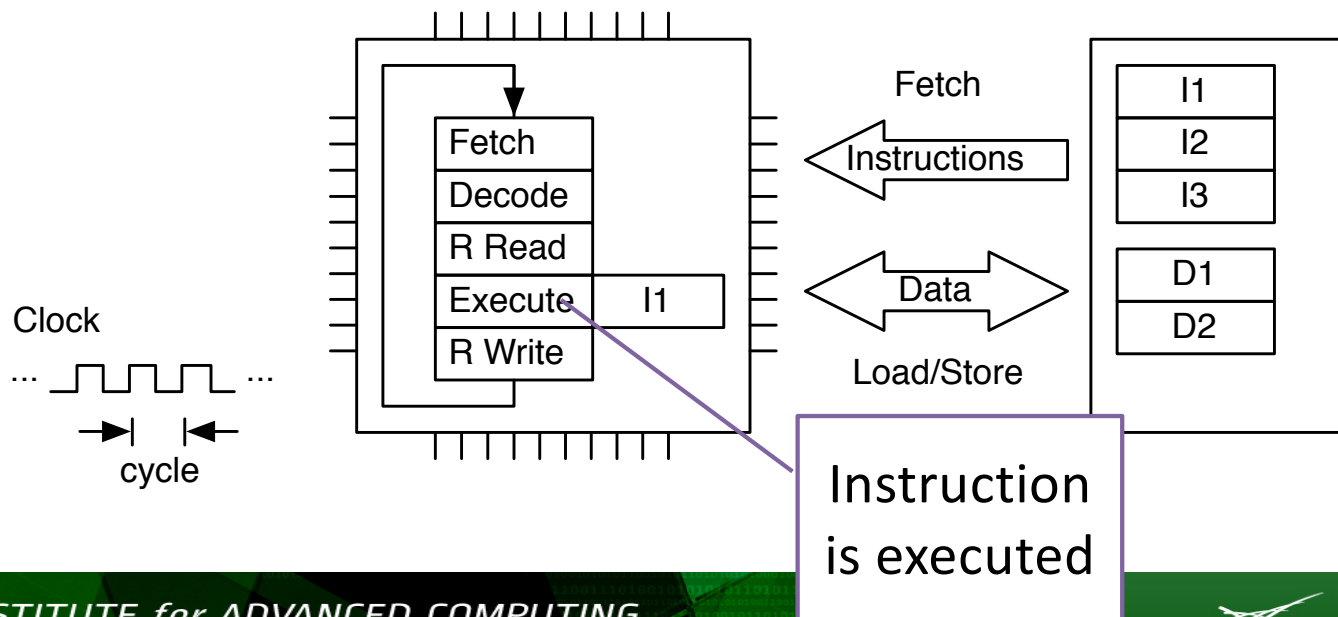- A single instruction may require multiple steps from fetch to completion

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

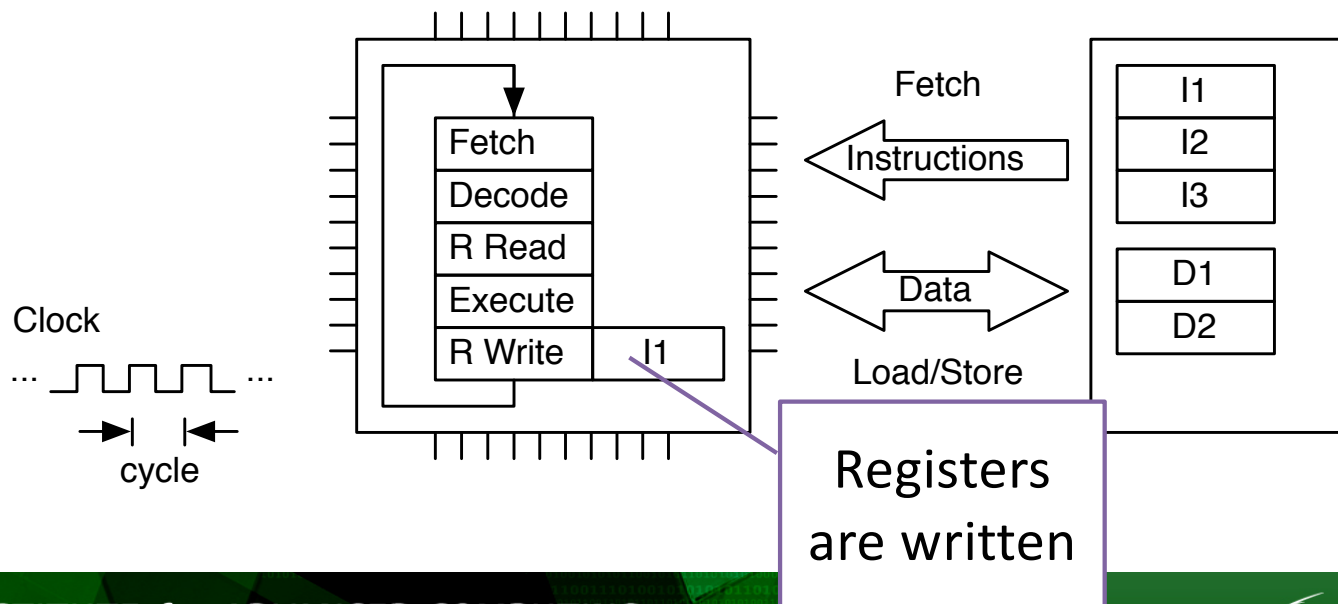- A single instruction may require multiple steps from fetch to completion

Fetch
Decode
R Read
Execute
R Write

Clock

... cycle

Fetch

Instructions
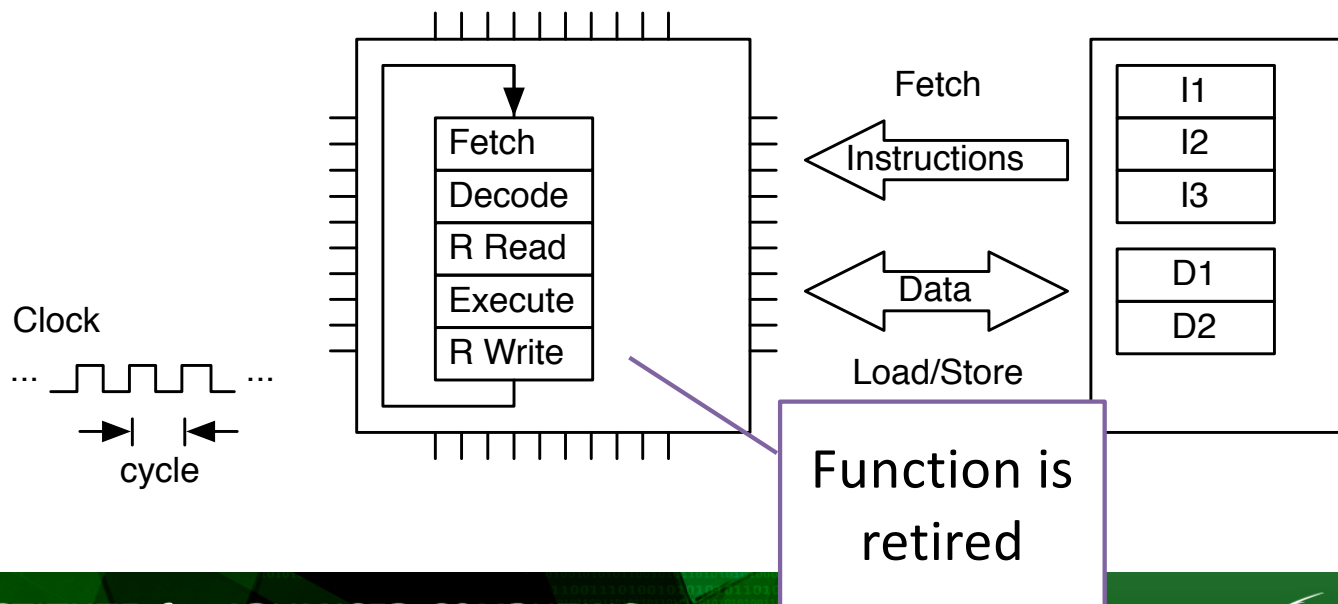
Data

Load/Store

I1
I2
I3

D1
D2

Function is retired
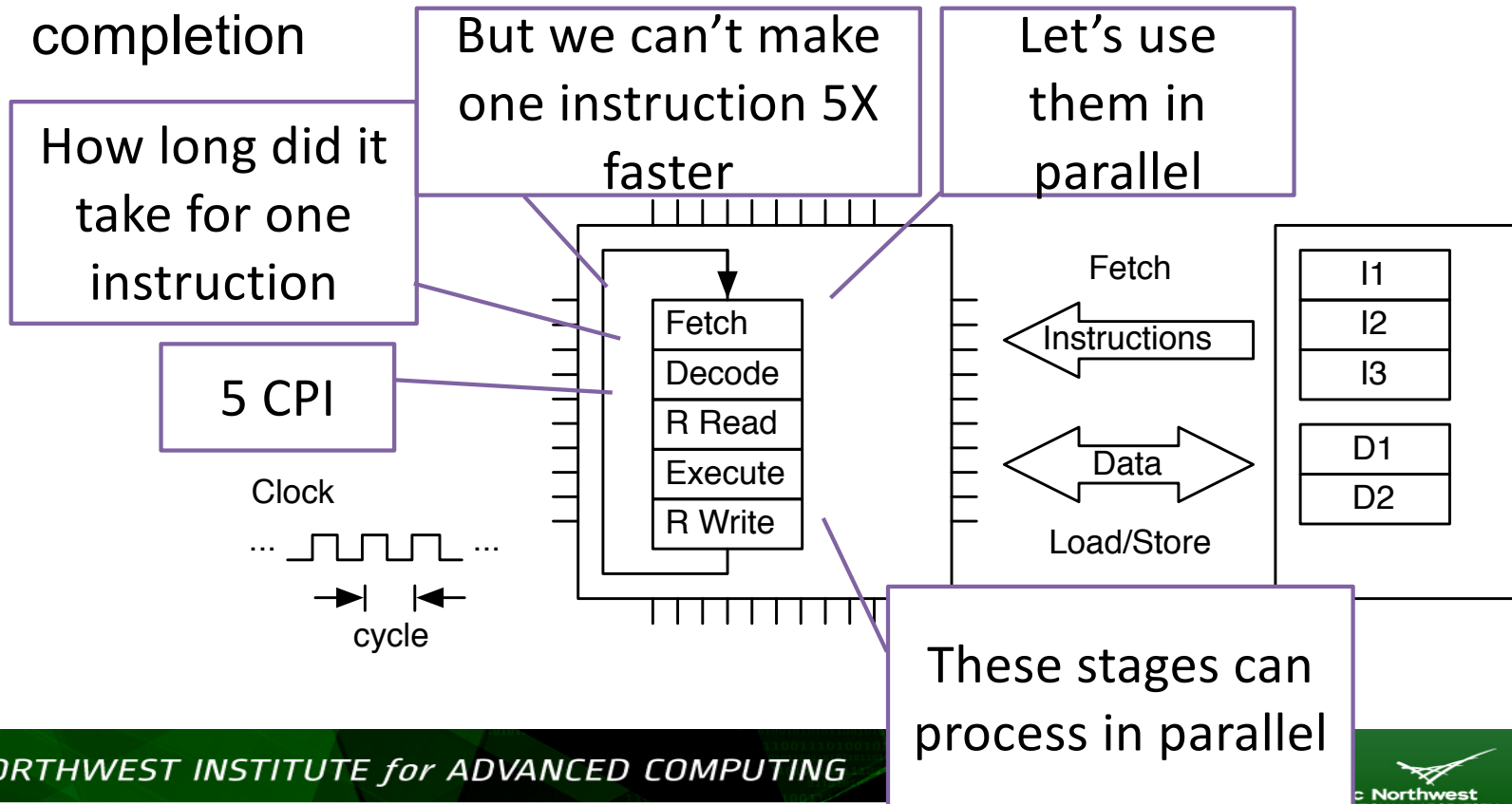
# Processor Core Instruction Handling

- 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

# Processor Core Instruction Handling

- By pipelining, multiple instructions can be executed at each clock cycle

- Form of instruction-level parallelism (ILP)

Fetch first instruction

Fetch

Instructions

Data

Load/Store

Fetch
Decode
R Read
Execute
R Write

I1

I1
I2
I3

D1
D2

Clock

... cycle

# 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

Fetch | I2
Decode | I1
R Read
Execute
R Write

Instructions

Data

I1
I2
I3

D1
D2

Clock

... ⊓⊔⊓⊔⊓ ...

cycle

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
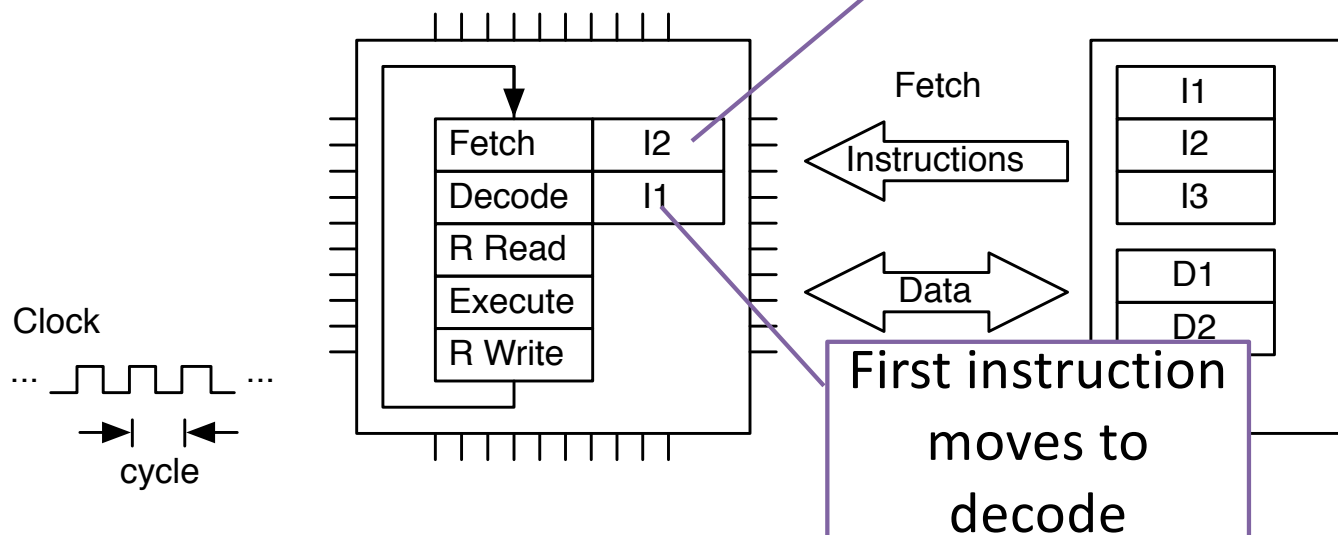
UNIVERSITY of
WASHINGTON

# 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

87

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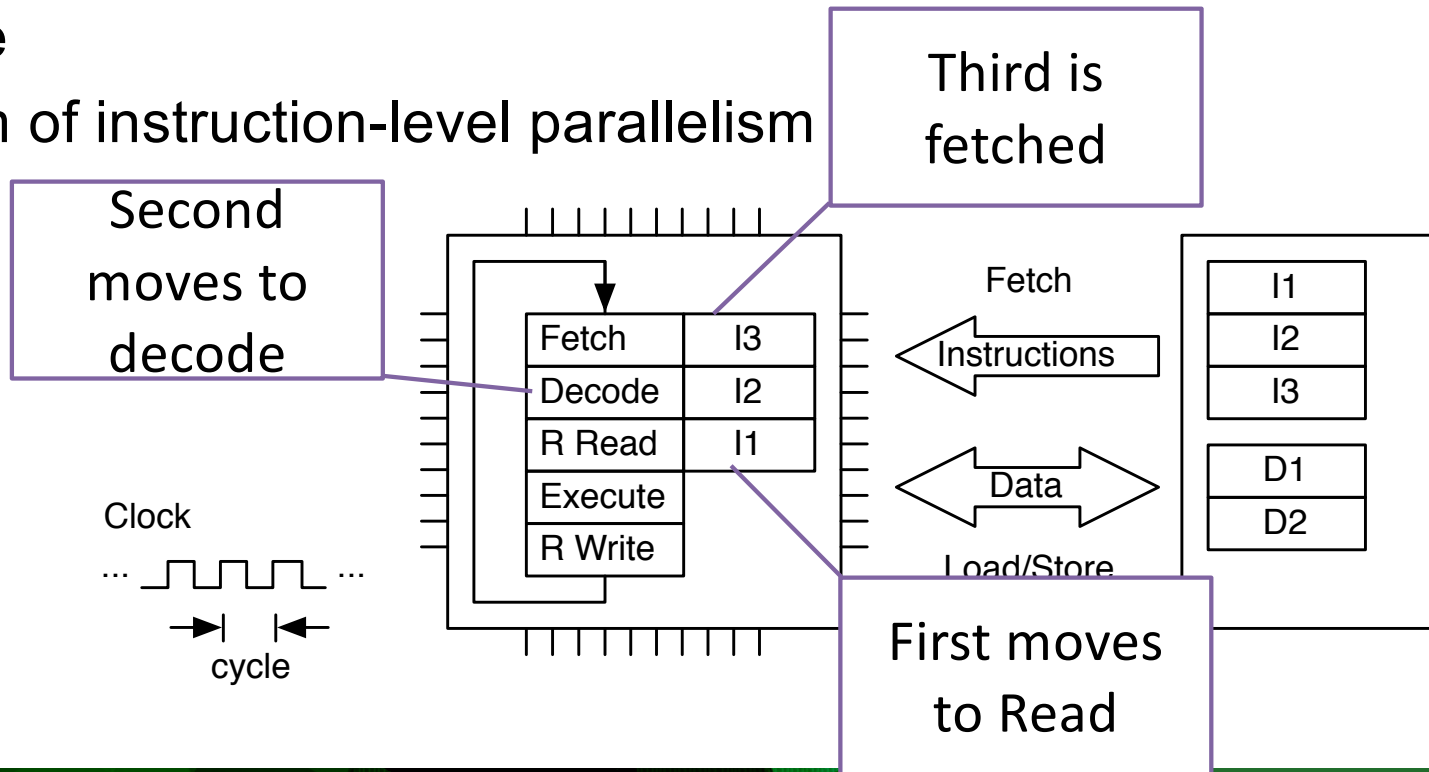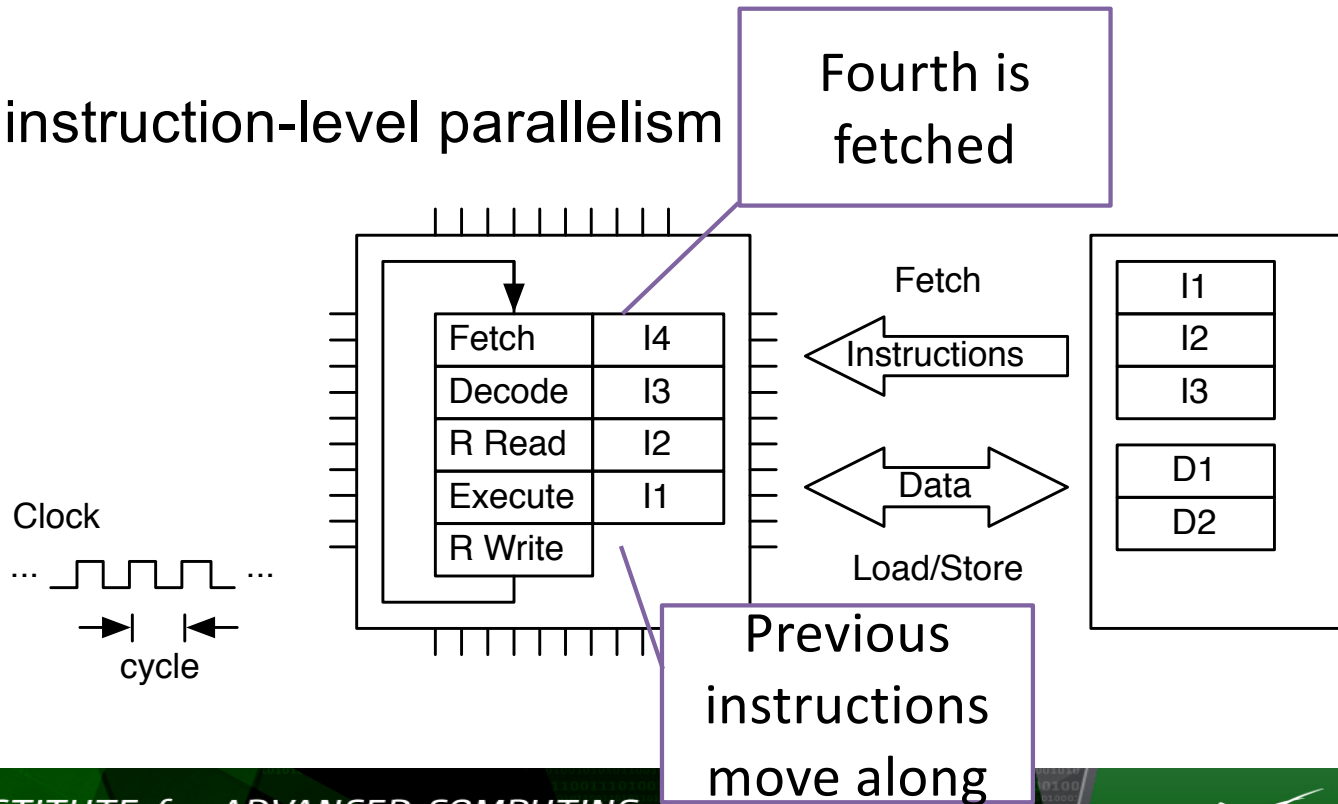University of Washington by Andrew Lumsdaine

# Processor Core Instruction Handling

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

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

Previous instructions move along

NORTHWEST INSTITUTE for ADVANCED COMPUTING

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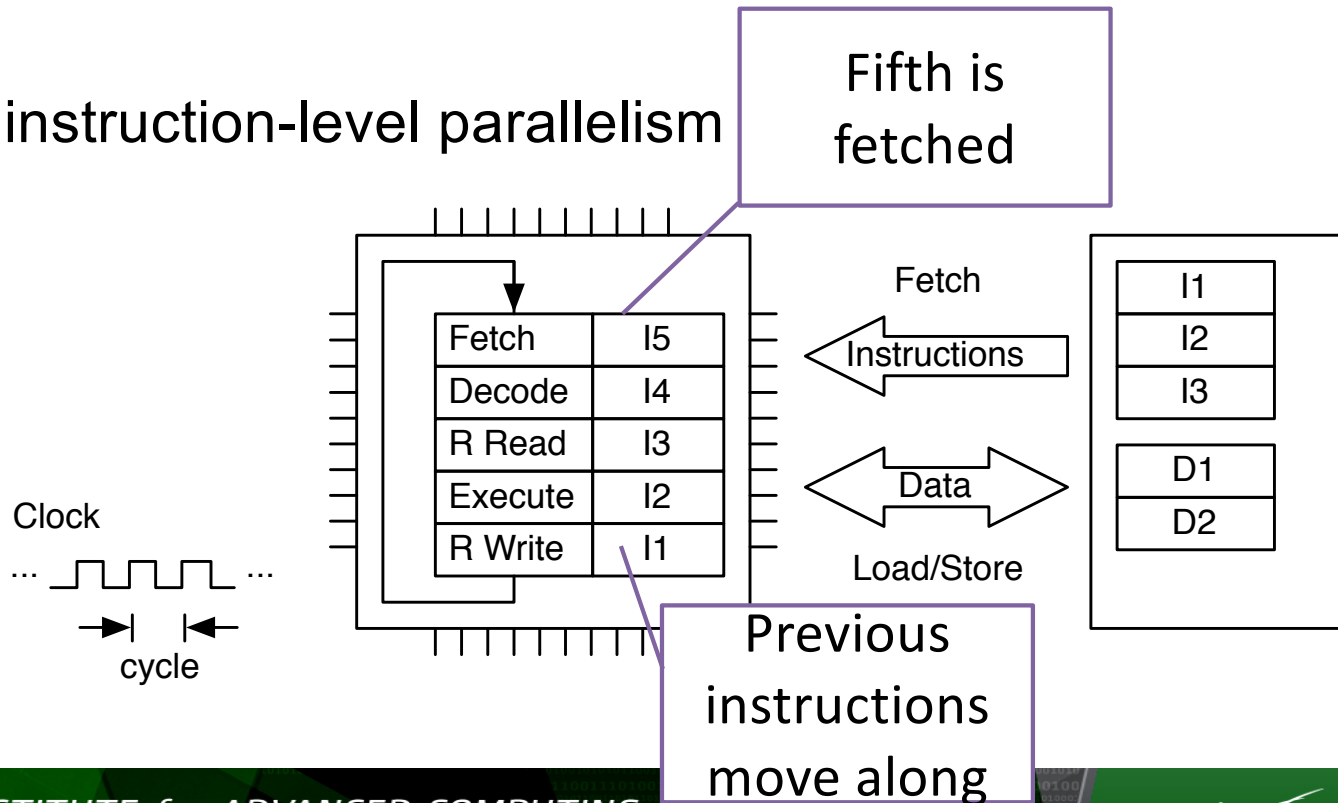
UNIVERSITY of
WASHINGTON

# Processor Core Instruction Handling

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

Fifth is fetched

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

Fetch
Instructions

Data

Load/Store

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

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

Clock
... ⊓⊔⊓⊔ ...
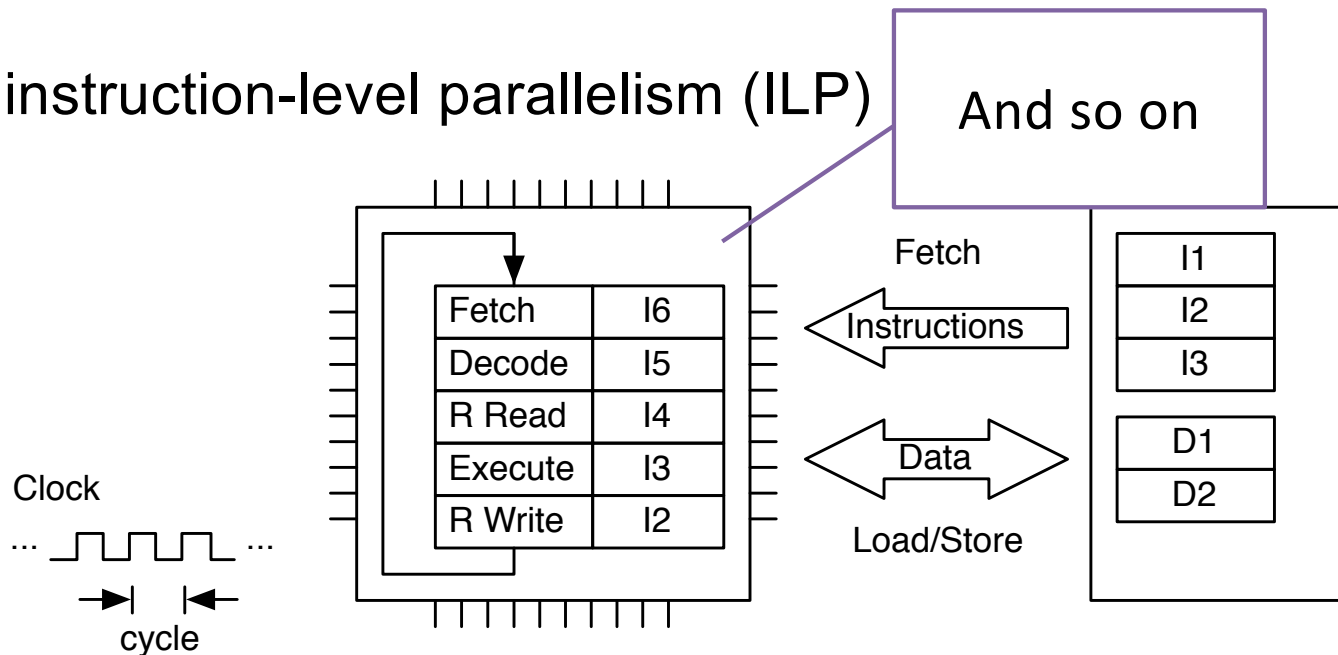cycle

Previous instructions move along

# Processor Core Instruction Handling

- By pipelining, multiple instructions can be executed at each clock cycle

- Form of instruction-level parallelism (ILP)

And so on

Fetch

| Fetch | I6 |
| Decode | I5 |
| R Read | I4 |
| Execute | I3 |
| R Write | I2 |

Clock

... cycle

Instructions

Data

Load/Store

| I1 |
| I2 |
| I3 |

| D1 |
| D2 |

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
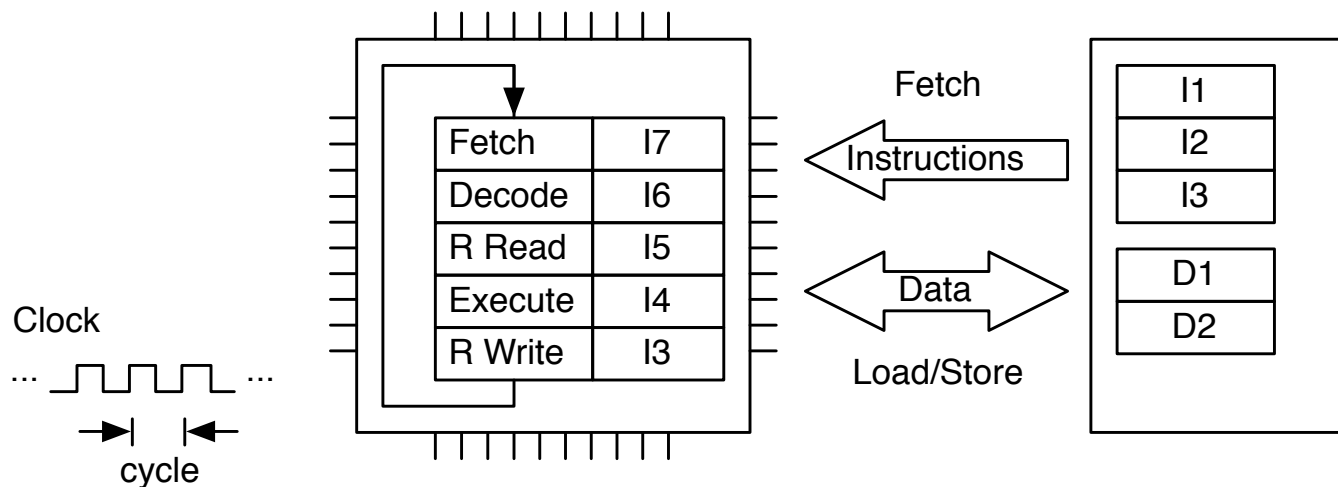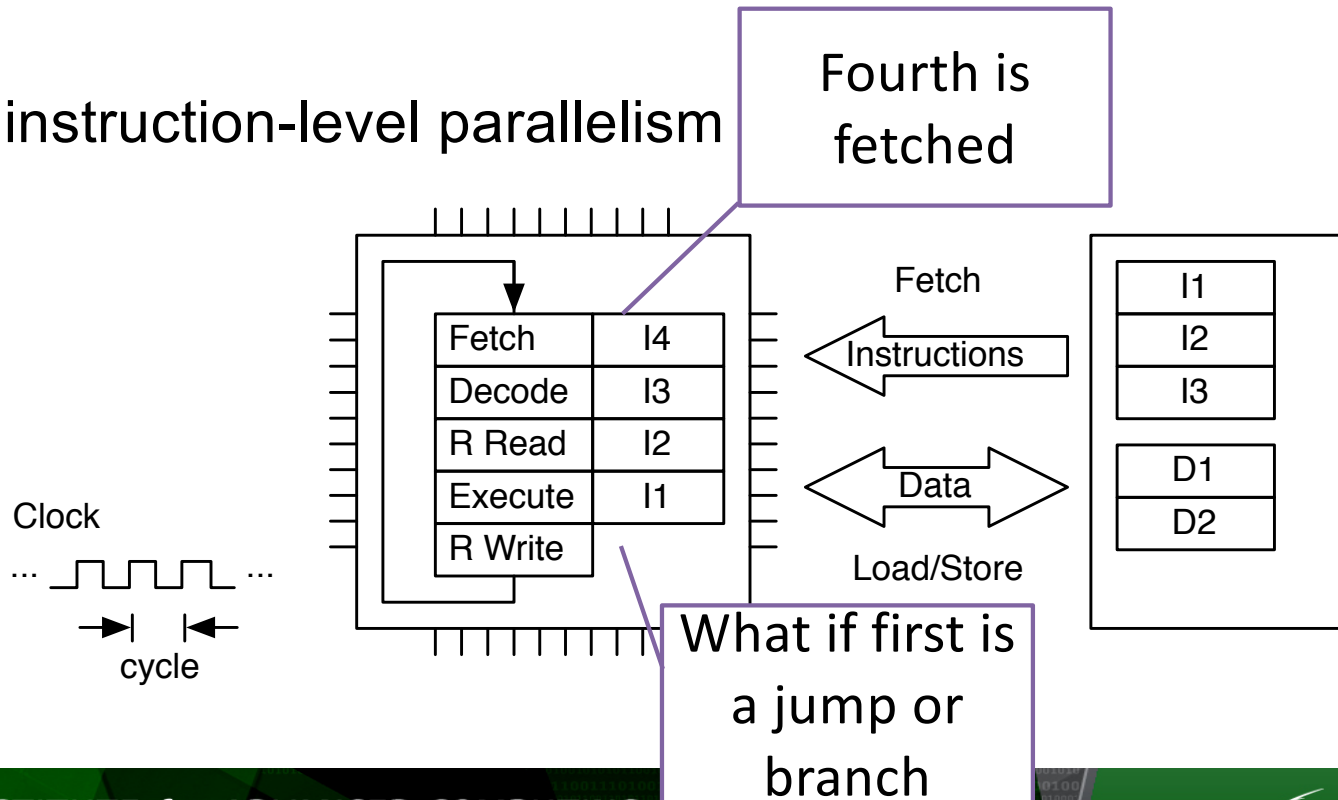
UNIVERSITY of
WASHINGTON

# Processor Core Instruction Handling

- By pipelining, multiple instructions can be executed at each clock cycle

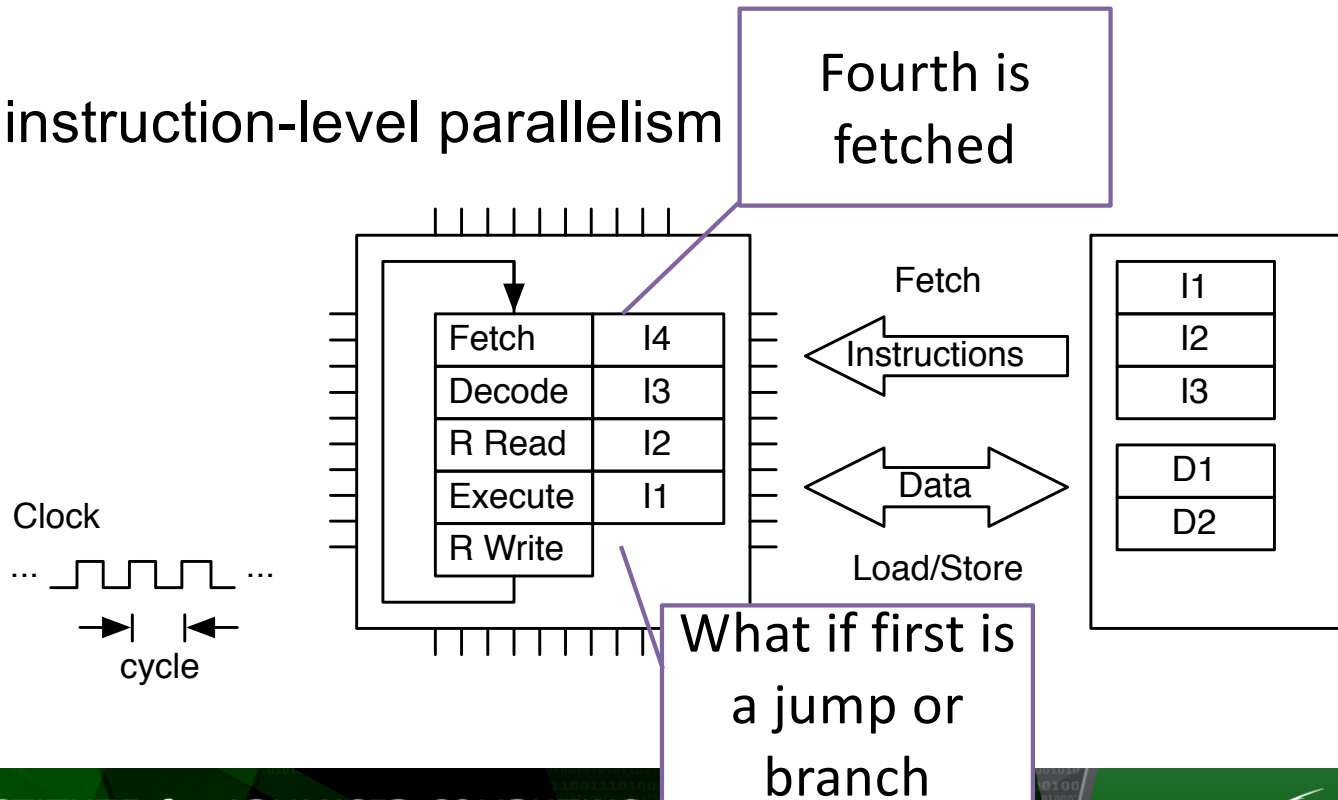- Form of instruction-level parallelism (ILP)

# 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 | |

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

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

# Pipeline Stall

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

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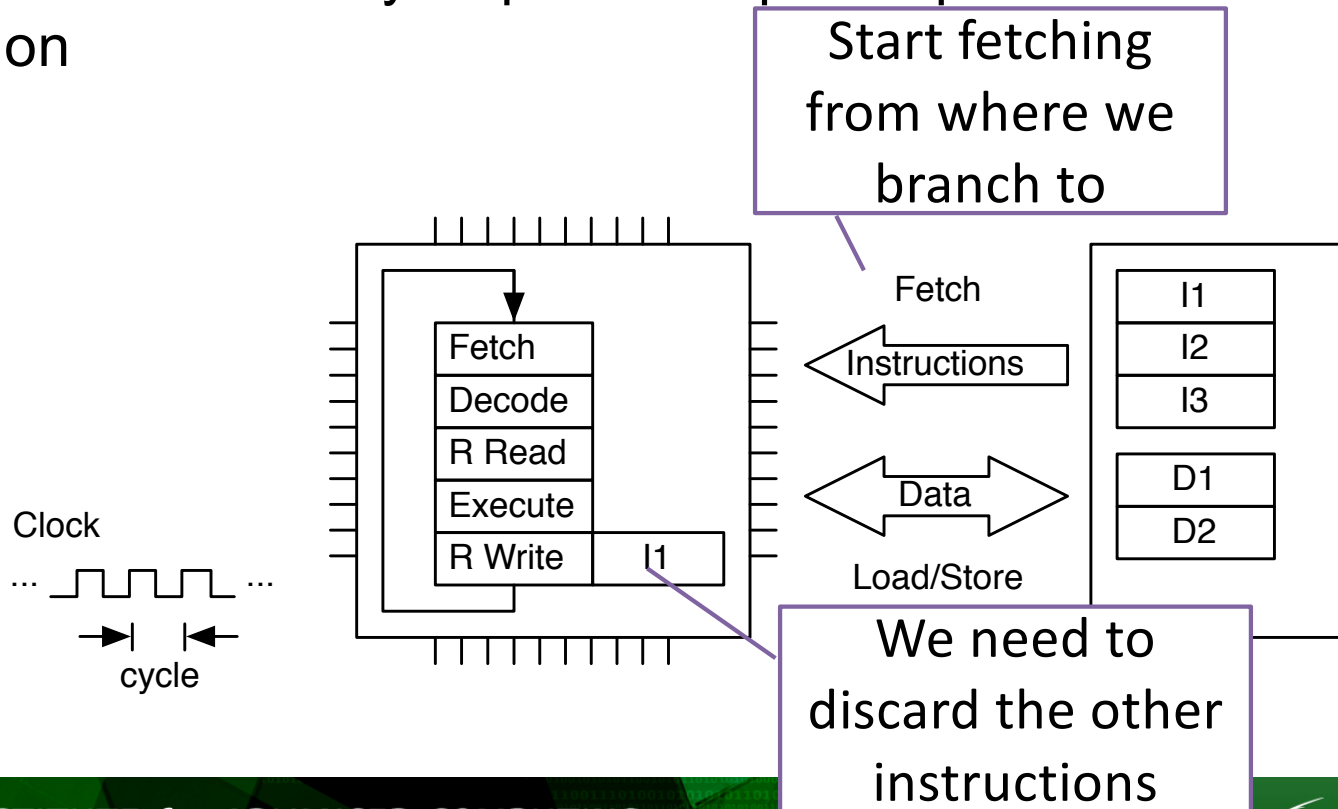
| Fetch | I4 |
|-------|----|
| Decode | I3 |
| R Read | I2 |
| Execute | I1 |
| R Write | |

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

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

I1
I2
I3

D1
D2

Data

Load/Store

Fetch
Decode
R Read
Execute
R Write | I1

Clock

... cycle

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
Decode
R Read
Execute
R Write

Clock

... cycle

Fetch
Instructions

Data
Load/Store

I1
I2
I3

D1
D2

# Branch Prediction

- Load the instructions we think will be branched to

This is the instruction the CPU predicts will branch to

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

Fetch

Instructions

Data

| I1 |
| I2 |
| I3 |

| D1 |
| D2 |

Clock

… ⊓⊔⊓⊔ …

cycle

First instruction moves to decode

NORTHWEST INSTITUTE for ADVANCED COMPUTING

96

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

# 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
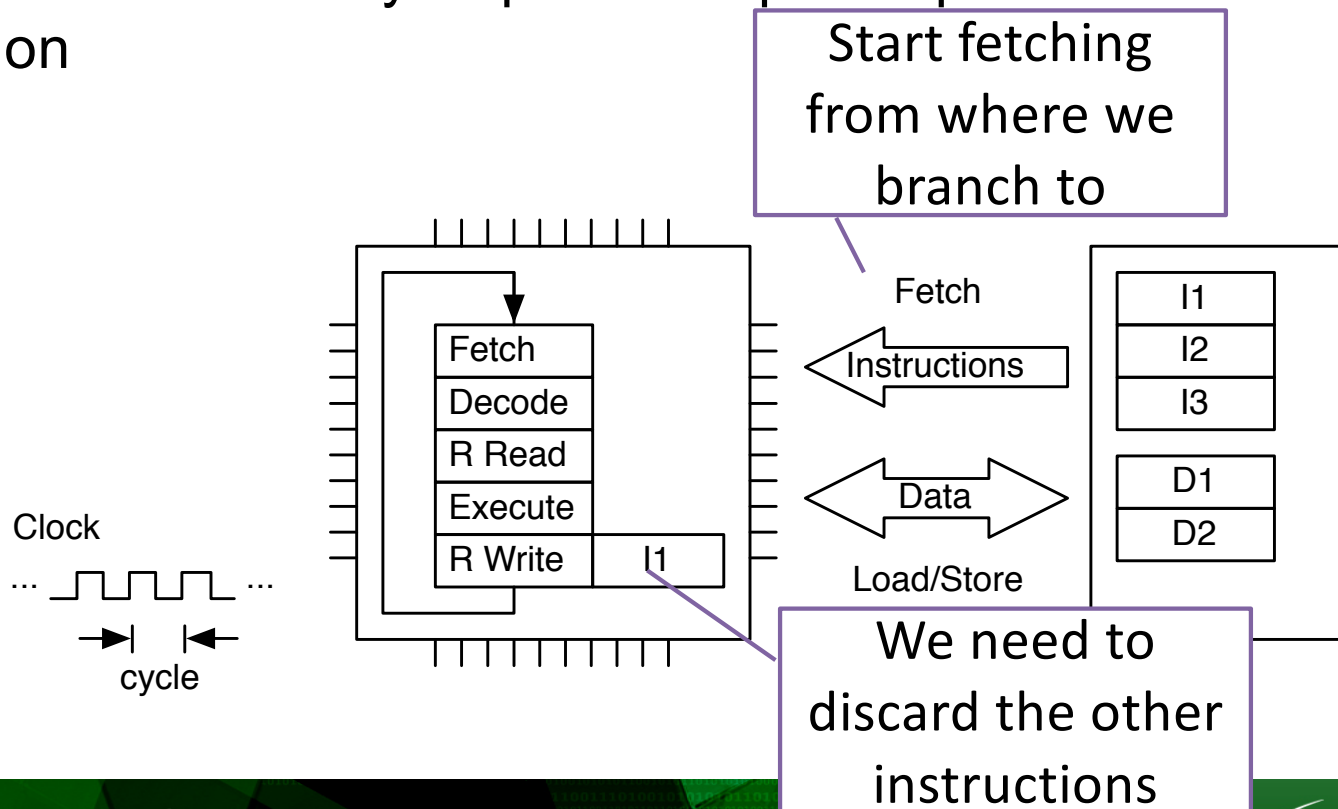Instructions

Data

Load/Store

I1
I2
I3

D1
D2

Clock
... ⊓⊔⊓⊔ ...
cycle

# 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 | |

Clock

... cycle

Data

Load/Store

| I1 |
| I2 |
| I3 |

| D1 |
| D2 |

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

Clock

... ⊔⊔⊔ ...

cycle

Fetch
Decode
R Read
Execute
R Write    I1

Fetch
Instructions

Data

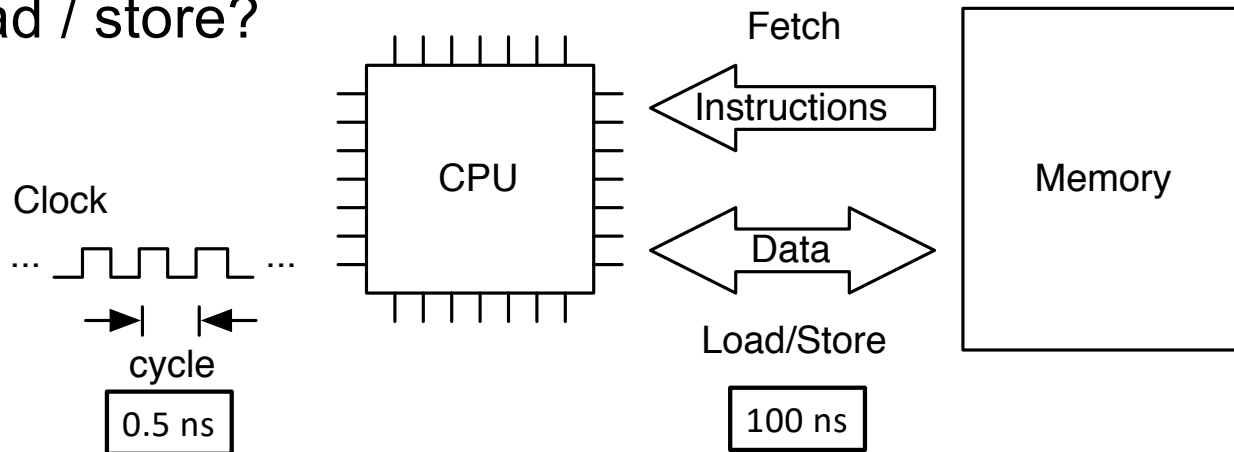Load/Store

I1
I2
I3

D1
D2

# 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

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

# Memory Access
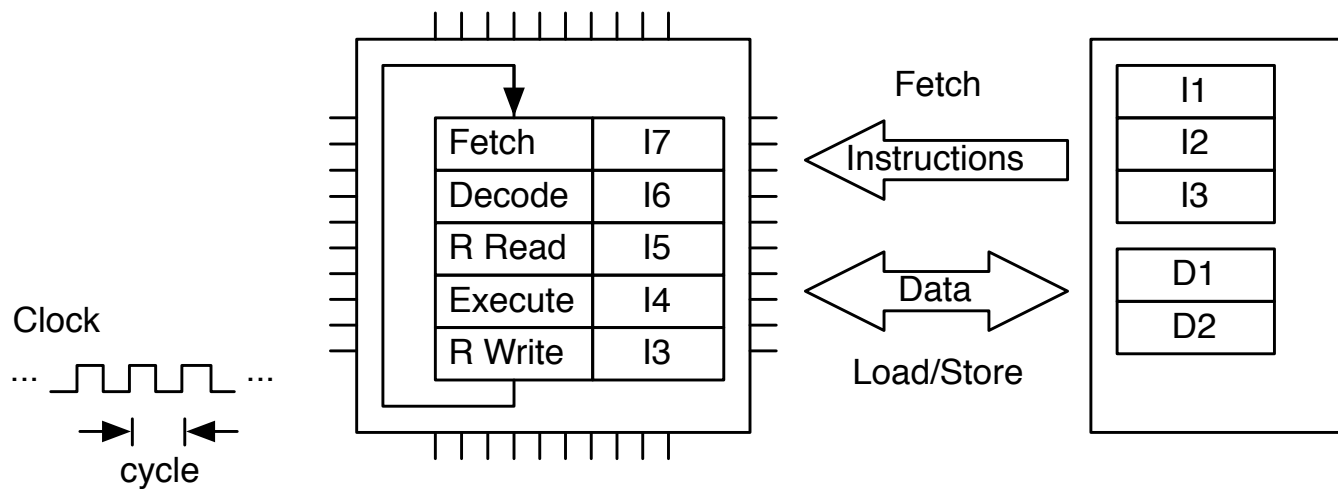
The next one may be cheaper

- 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

101

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
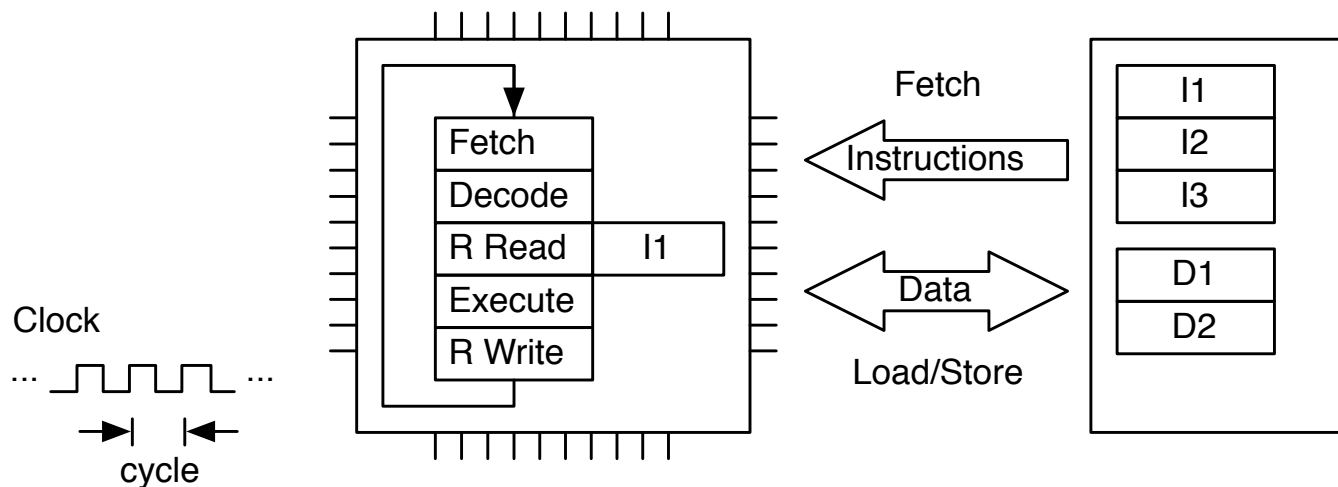
UNIVERSITY of
WASHINGTON

# Memory Access Costs
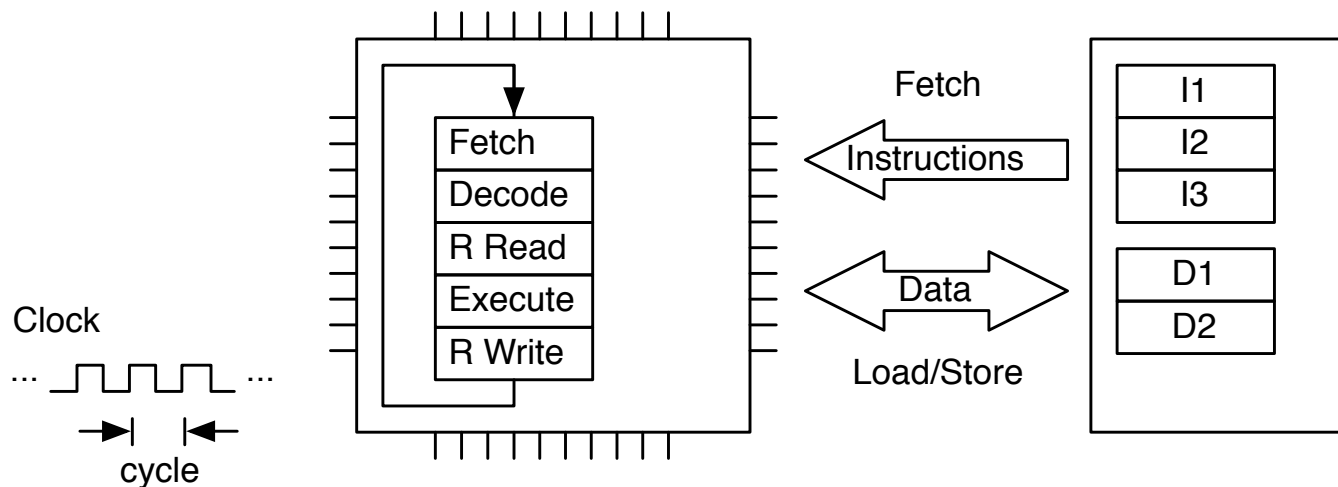
- Access to main memory has huge impact on performance

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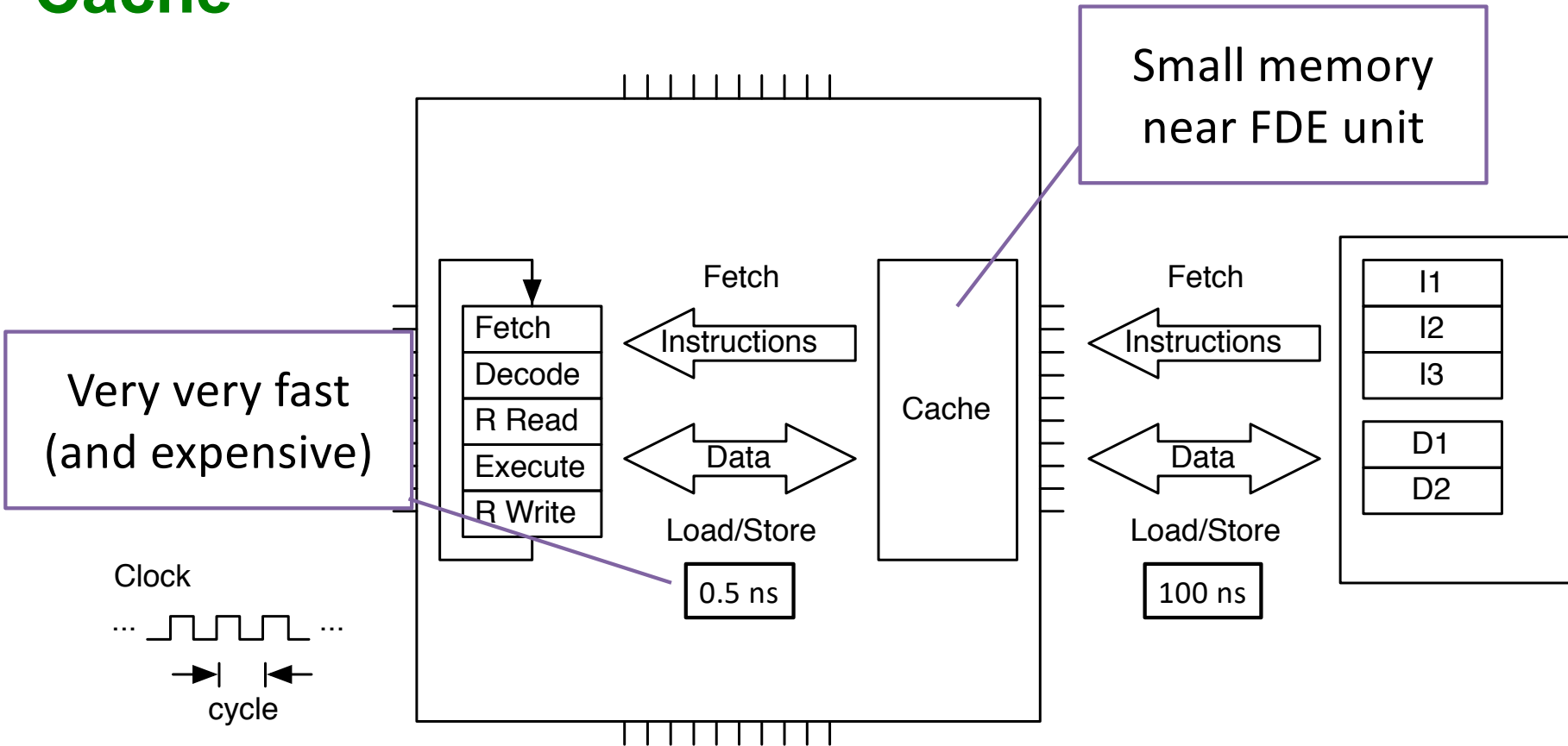
- Access to main memory has huge impact on performance (600X)
- Processor would be idle almost all the time



Clock

... cycle

Fetch
Decode
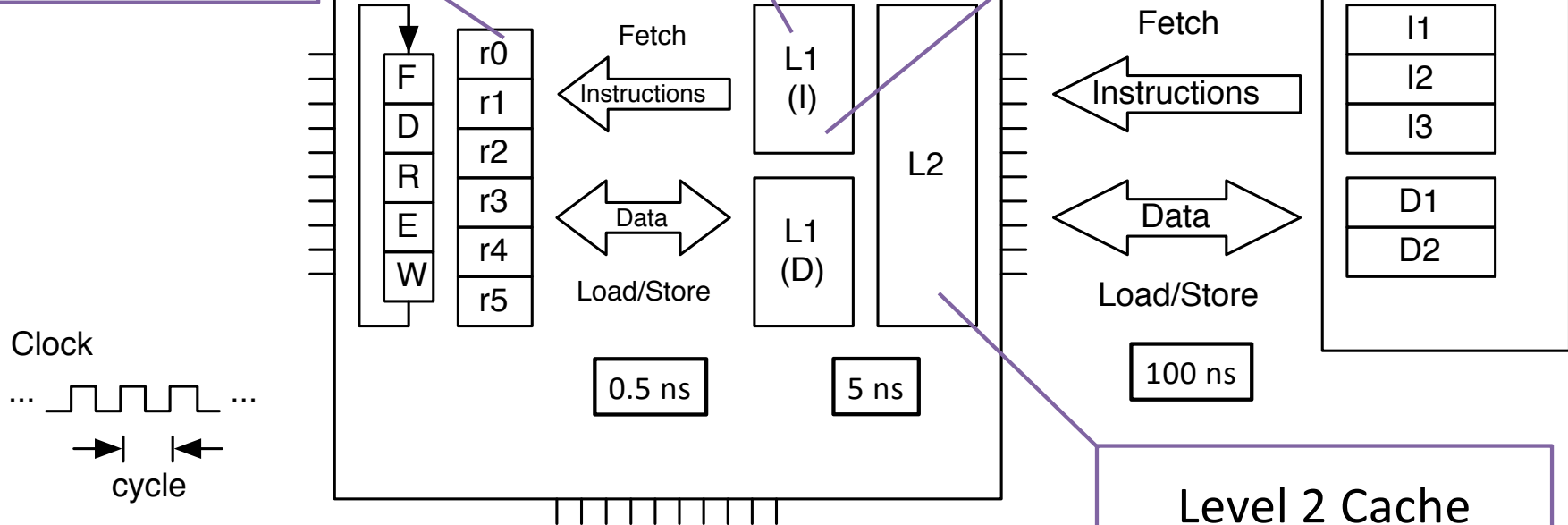R Read
Execute
R Write

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

# Cache

Small memory near FDE unit

Very very fast (and expensive)

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

Fetch

← Instructions

Cache

← Data →

Load/Store

0.5 ns

Fetch

← Instructions

← Data →

Load/Store

100 ns

| I1 |
| I2 |
| I3 |

| D1 |
| D2 |

Clock

… ⊓⊔⊓⊔ …

→| |←

cycle

NORTHWEST INSTITUTE for ADVANCED COMPUTING

105

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

# Hierarchical Memory

Level 1 Cache
(very very fast)

Separate L1 for
instuctions/data

Registers
(immediately fast)

Fetch

Instructions

Data

Load/Store

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

L1
(I)

L1
(D)

L2

0.5 ns

5 ns

Fetch

Instructions

Data

Load/Store

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

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

100 ns

Level 2 Cache
(pretty fast)

Clock

...  ⊓⊔⊓⊔  ...

cycle

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

UNIVERSITY of
WASHINGTON

# Hierarchical Memory

There is also an MMU and TLB

Data goes from L2 to L1

FDE works with data in registers

Data goes from L1 to registers

Fetch
Instructions

L1 (I)

L2

Fetch
Instructions

I1
I2
I3

D1
D2

F D R E W

r0 r1 r2 r3 r4 r5

Data
Load/Store

L1 (D)

Data
Load/Store

Clock
... ⎍⎍⎍ ...
cycle

0.5 ns        5 ns        100 ns

Data goes from main memory to L2

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

UNIVERSITY of WASHINGTON
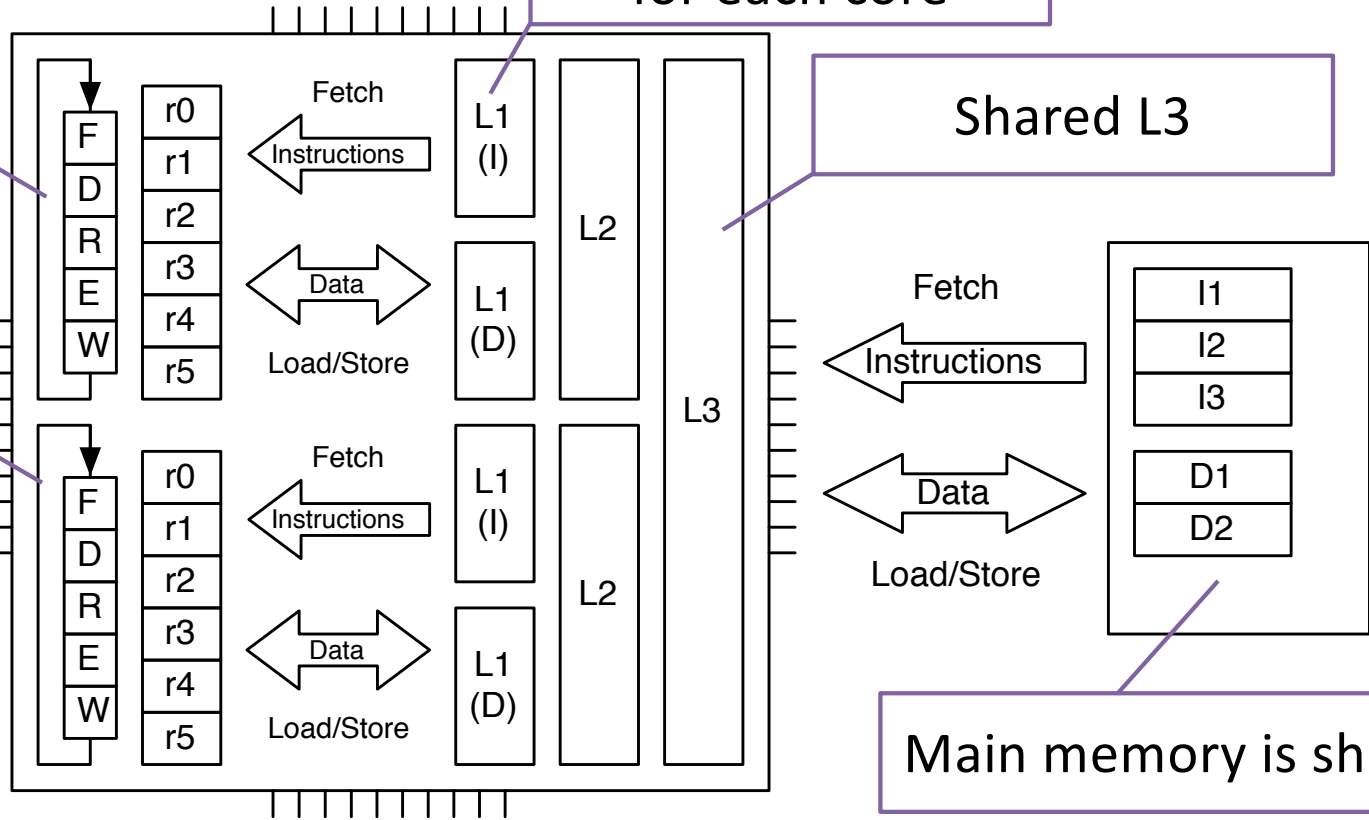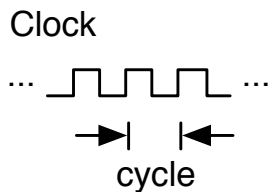
# Cache and Multicore



Separate L1 and L2 for each core

Cores work on separate register sets and instrs

Shared L3

Cores work on separate register sets and instrs

Main memory is shared

NORTHWEST INSTITUTE for ADVANCED COMPUTING

108

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



Highest performance with data here

Higher performance with data here

Not so good performance with data here

Clock

... cycle

Fetch
Instructions

Data
Load/Store

F
D
R
E
W

r0
r1
r2
r3
r4
r5

L1
(I)

L1
(D)

L2

I1
I2
I3

D1
D2

NORTHWEST INSTITUTE for ADVANCED COMPUTING

109

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



Keep as much data here as possible

Keep as much data here as possible

Keep as much data here as possible

Fetch Instructions

Data Load/Store

Clock

... ⊓⊔⊓⊔ ...

cycle

F D R E W

r0 r1 r2 r3 r4 r5

L1 (I)

L1 (D)

L2

Fetch Instructions

Data Load/Store

I1 I2 I3

D1 D2

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
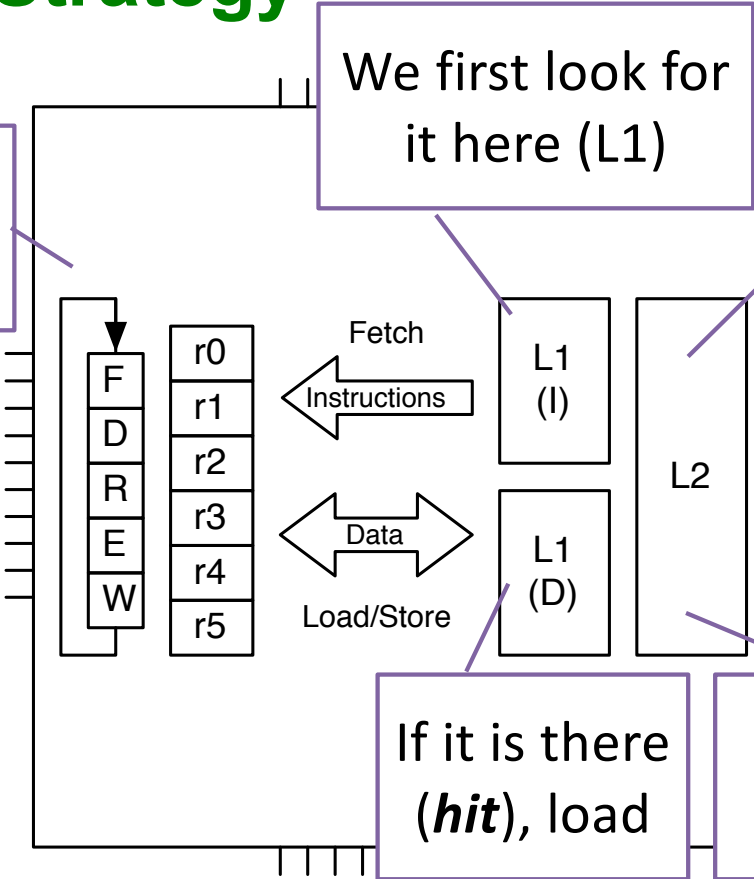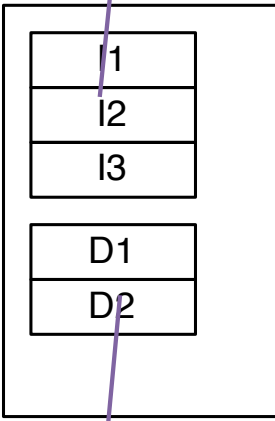
# Locality → Strategy



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?
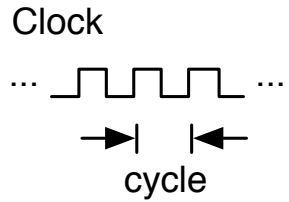
Fetch

Instructions

Data

Load/Store

Fetch

Instructions

Data

Load/Store

L1 (I)

L1 (D)

L2

F D R E W

r0 r1 r2 r3 r4 r5

I1 I2 I3

D1 D2

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

NORTHWEST INSTITUTE for ADVANCED COMPUTING

111

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

When we need the next operand

We want it to be here

Or here



Fetch
Instructions

L1 (I)

L2

Fetch
Instructions

I1
I2
I3

Data
Load/Store

L1 (D)

Data
Load/Store

D1
D2

F
D
R
E
W

r0
r1
r2
r3
r4
r5

Clock

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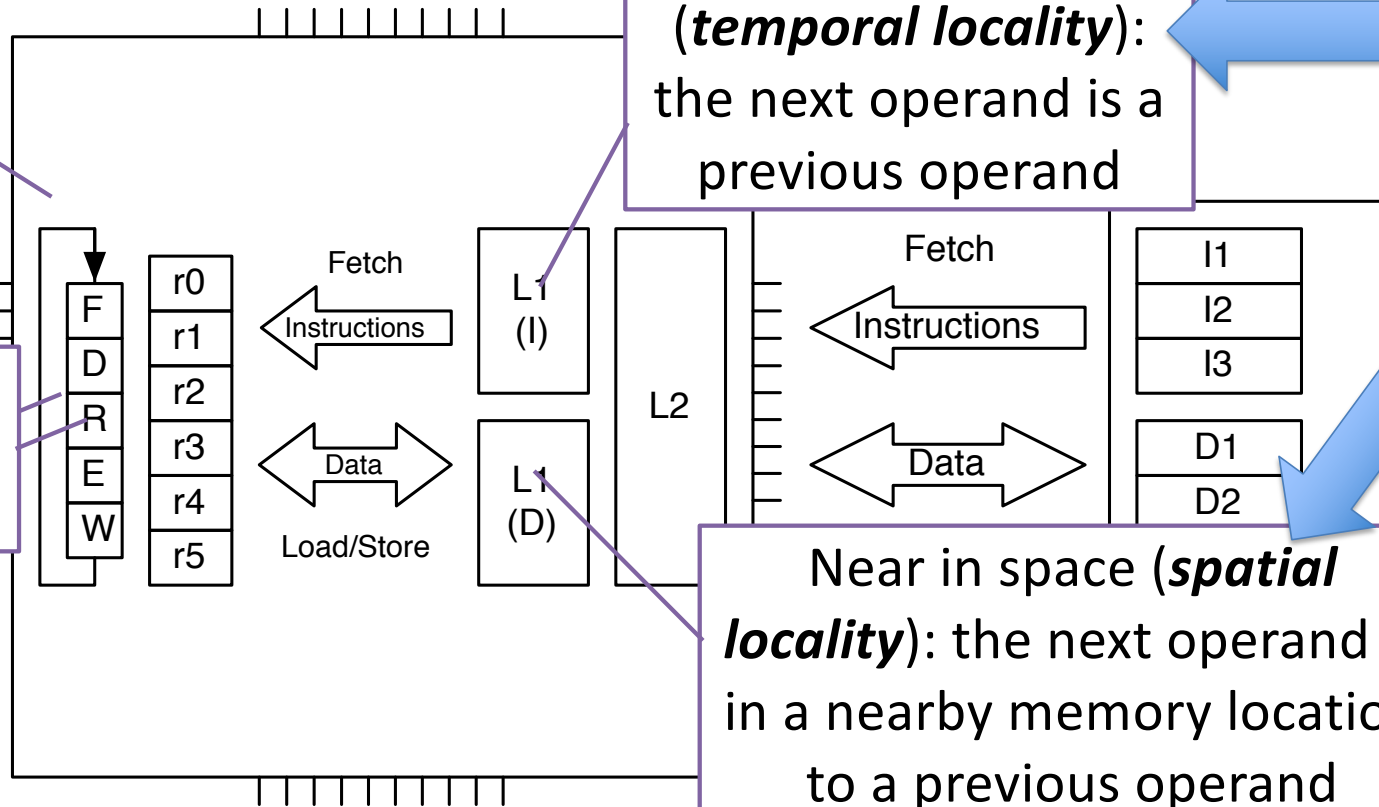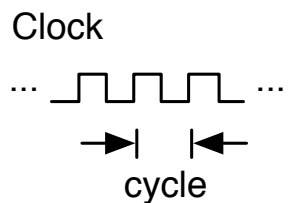
On a miss, copy the data we want *and its neighbors*

# 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)

L2

L1
(D)

Fetch
Instructions

Data

I1
I2
I3

D1
D2

Clock

... ⊓⊔⊓⊔⊓ ...

cycle

NORTHWEST INSTITUTE *for* ADVANCED COMPUTING

113

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

114

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

# Thank you!

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*

**W**
UNIVERSITY *of*
WASHINGTON

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 Spring 2019
University of Washington by Andrew Lumsdaine