

AMATH 483/583

High Performance Scientific Computing

Lecture 10:

Processes, Threads, Concurrency, Parallelism

Andrew Lumsdaine
Northwest Institute for Advanced Computing
Pacific Northwest National Laboratory
University of Washington
Seattle, WA

Overview

- Multiple cores
- Concurrency
- Processes
- Threads
- Parallelization strategies
- Correctness

Supercomputers (HPC)



NORTHWEST INSTITUTE for ADVANCED COMPUTING

3

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

Schematically

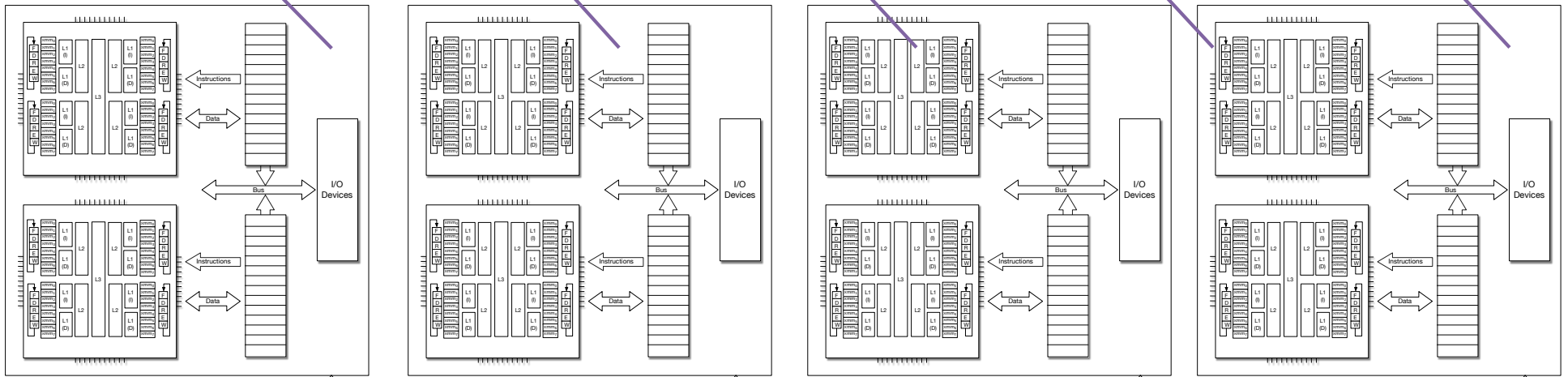
Put sockets
on a blade

Put blades
in a chassis

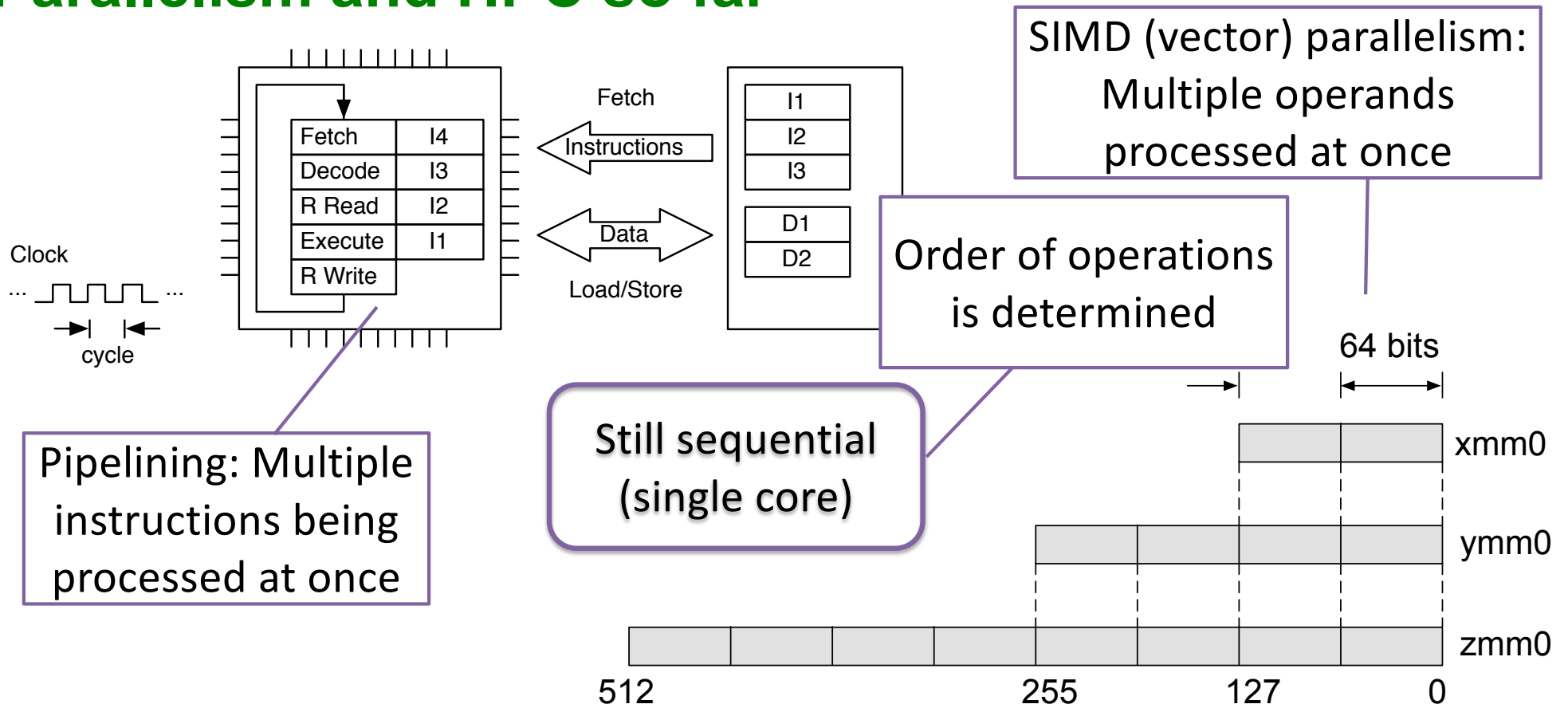
Put chassis
in a rack

Put racks in
a center

Put centers
in the cloud



Parallelism and HPC so far



General Performance Principles

- Work harder

- Faster core

Dennard scaling
(ended 2005)

- Work smarter

- Branch predictions, etc
- Better compilation
- Better algorithm
- Better implementation

What
about this?

We did this

- Get help

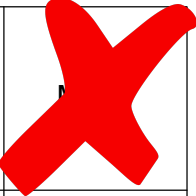
Parallel
Computing

Flynn's Taxonomy (Aside)

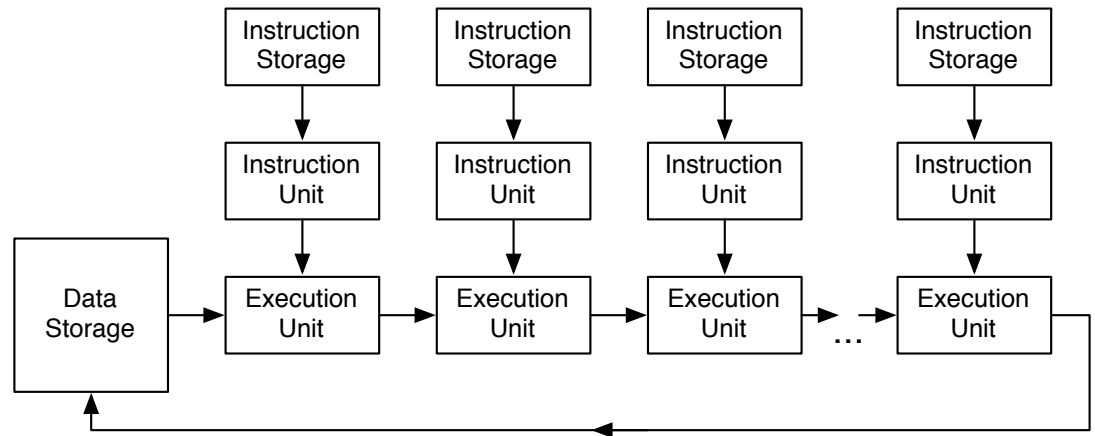
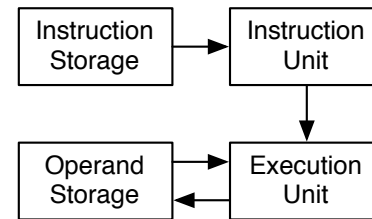
Anyone in HPC must know Flynn's taxonomy

- **Classic** classification of parallel architectures (Michael Flynn, 1966)

Plain old sequential

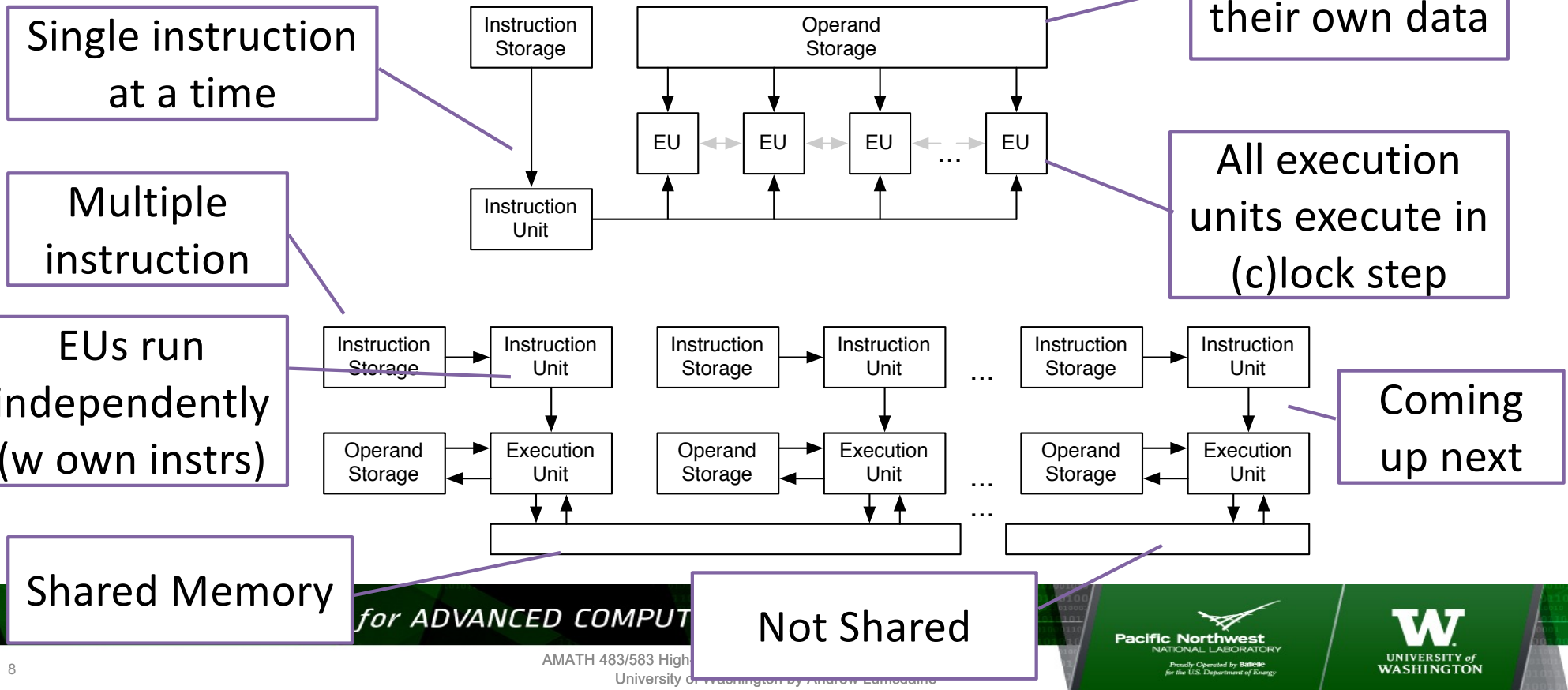
	Single Instruction	Multiple Instruction
Single Data	SISD	
Multiple Data	SIMD	MIMD

Based on multiplicity of instruction streams, data storage



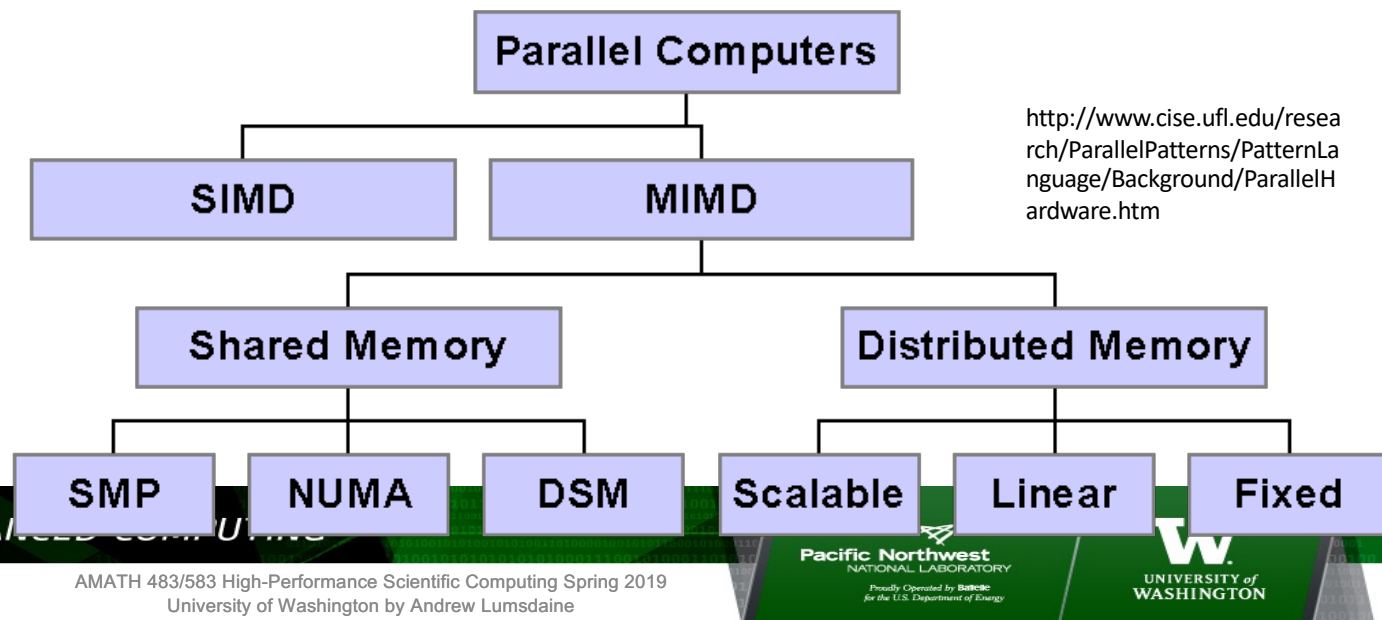
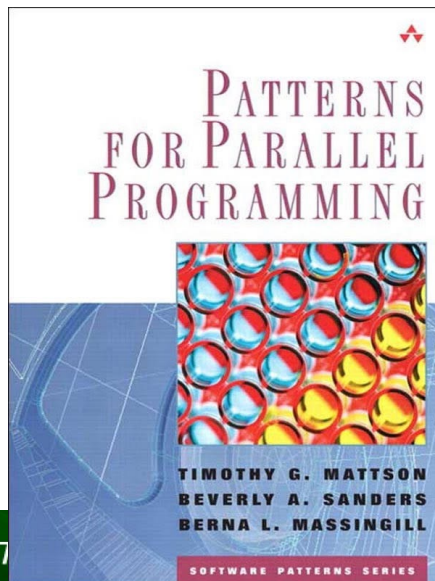
SIMD and MIMD

- Two principal parallel computing paradigms (multiple CPUs) But each have their own data



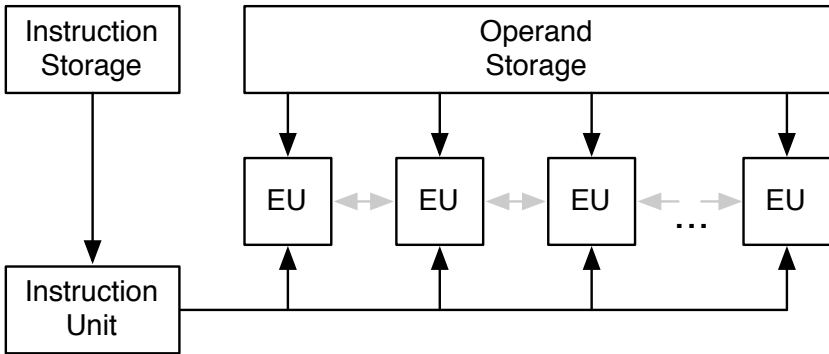
A More Refined (Programmer-Oriented) Taxonomy

- Three major modes: SIMD, Shared Memory, Distributed Memory
- Different programming approaches are generally associated with different modes of parallelism (threads for shared, MPI for distributed)
- A modern supercomputer will have all three major modes present



SIMD in SSE/AVX

Flynn's original conceptual model

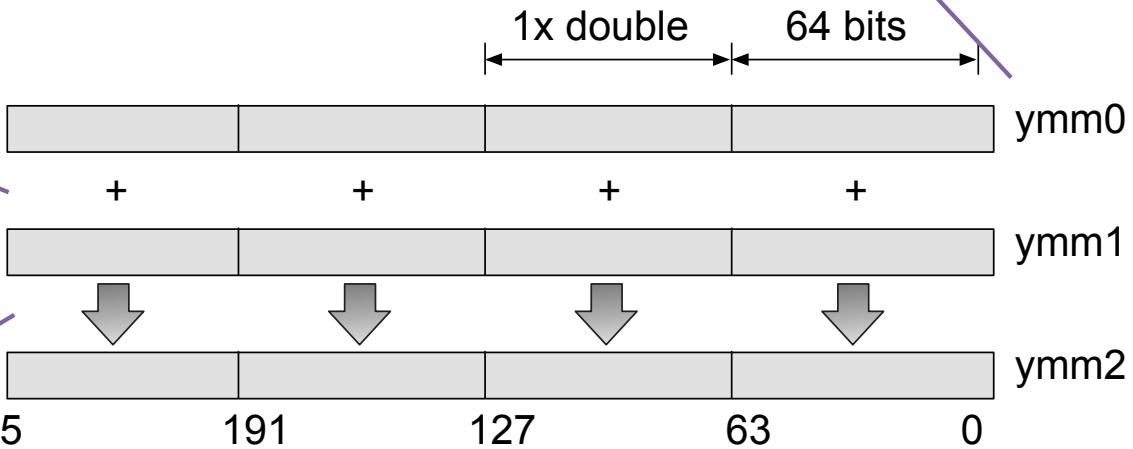


ymm are 256 bit registers

```
vfadd231pd %ymm0, %ymm1, %ymm2
```

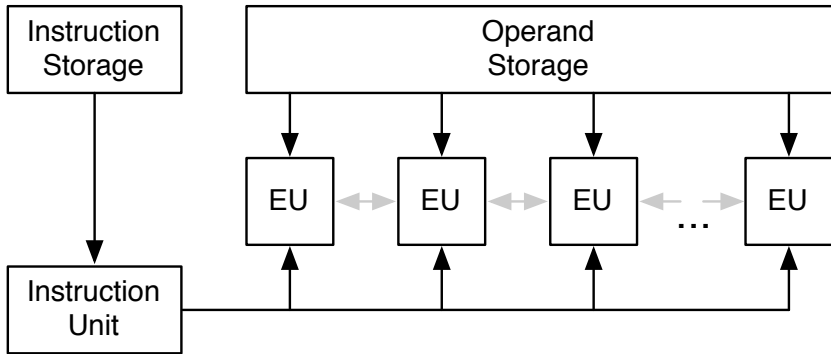
One machine instruction

Adds all four doubles *simultaneously*



SIMD in SSE/AVX

Flynn's original conceptual model

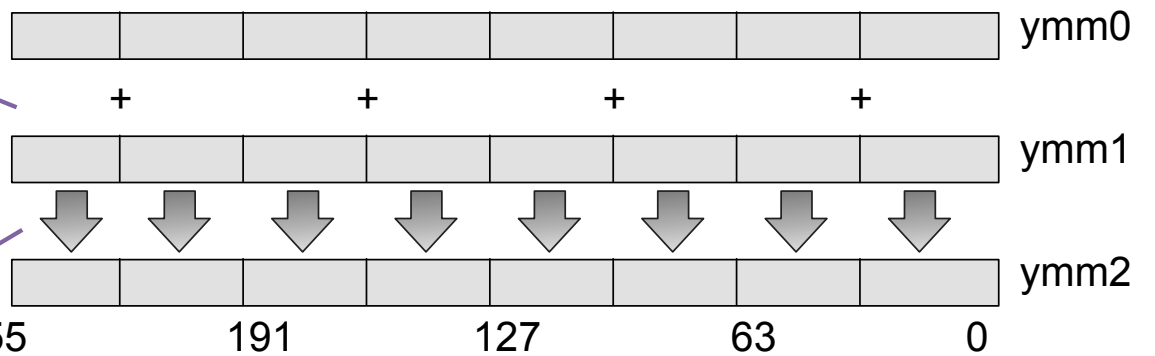


ymm are 256 bit registers



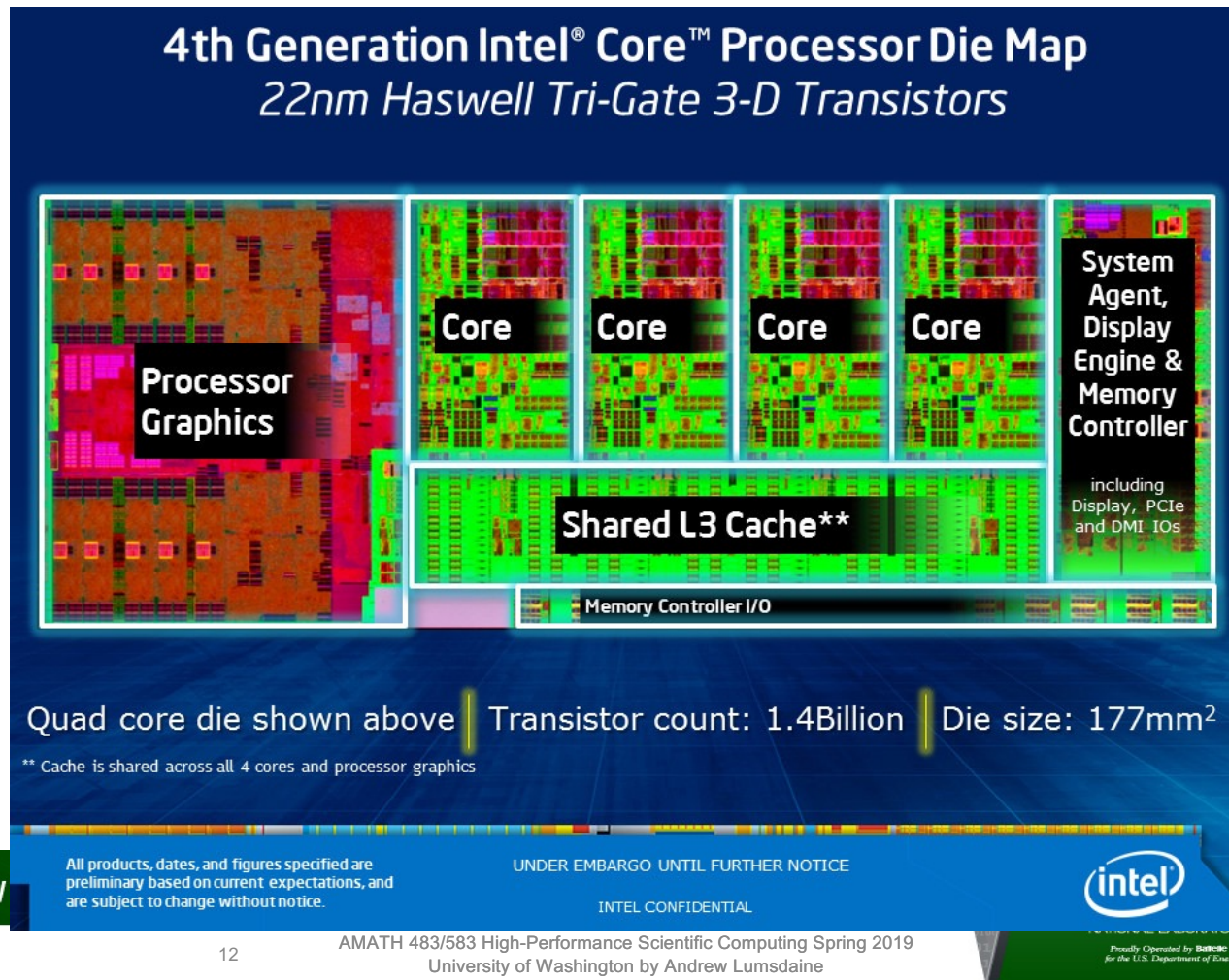
```
vfadd231ps %ymm0, %ymm1, %ymm2
```

One machine instruction



Adds all eight floats *simultaneously*

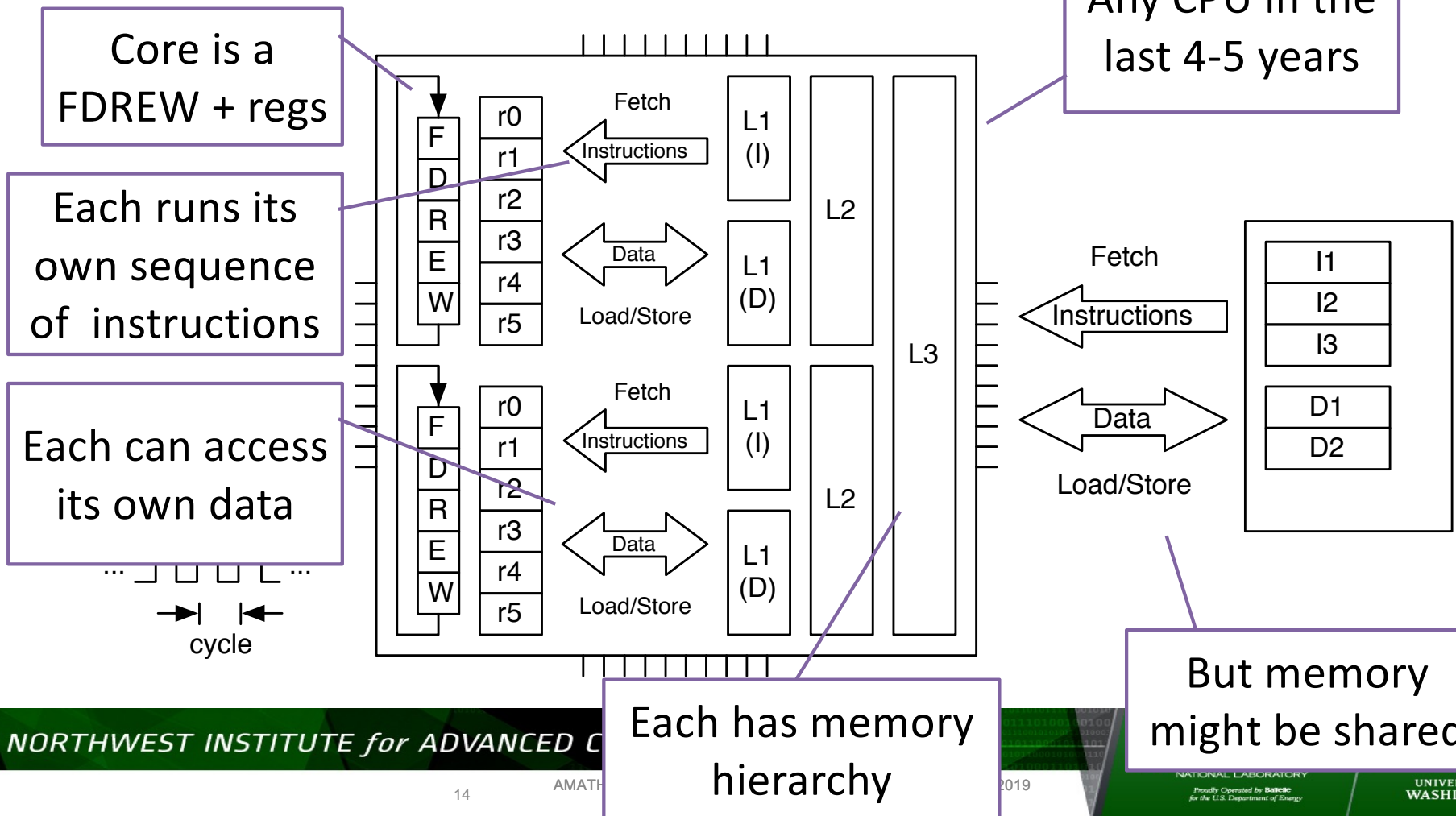
Multicore Architecture



Multicore for HPC

- How do multicore chips operate (how does the hardware work)?
- How do they get high performance?
- How does the software exploit the hardware (how do we write our software to exploit the hardware)?
- What are the abstractions that we need to use to reason about multicore systems?
- What are the programming abstractions and mechanisms?
- Terminology: Program, process, thread
- More terminology: Parallel, concurrent, asynchronous

Multicore Architecture



Core is a FDREW + regs

Each runs its own sequence of instructions

Each can access its own data

Any CPU in the last 4-5 years

Each has memory hierarchy

But memory might be shared

Parallelization Example

- You are the TA for CSE 142 and have to grade 22 exams
 - The exam has 8 questions on it
 - It takes 3 minutes to grade one question
-
- How long will it take you to grade all of the exams?



Parallelization Example

- You are the TA for CSE 142 and have to grade 22 exams
 - The exam has 8 questions on it
 - It takes 3 minutes to grade one question
 - You ask 21 friends who agree to help you
-
- How long will it take the 22 of you to grade all of the exams?
-
- Describe your approach
 - List your assumptions



Parallelization Example

- You are the TA for CSE 142 and have to grade 1012 exams ($1012 = 46 * 22$)
- The exam has 8 questions on it
- It takes 3 minutes to grade one question
- You ask 21 friends who agree to help you
- How long will it take the 22 of you to grade all of the exams?
- Describe your approach
- Describe another approach
- List your assumptions



Parallelization Example

- You are the TA for CSE 142 and have to grade 8 exams
 - The exam has 22 questions on it
 - It takes 3 minutes to grade one question
 - You ask 21 friends who agree to help you
-
- How long will it take the 22 of you to grade all of the exams?
-
- Describe your approach



Parallelization Example

- You are the TA for CSE 142 and have to grade 368 exams ($368 = 46 * 8$)
- The exam has 22 questions on it
- It takes 3 minutes to grade one question
- You ask 21 friends who agree to help you
- How long will it take the 22 of you to grade all of the exams?
- What if you had 368 friends? $368 * 22$?

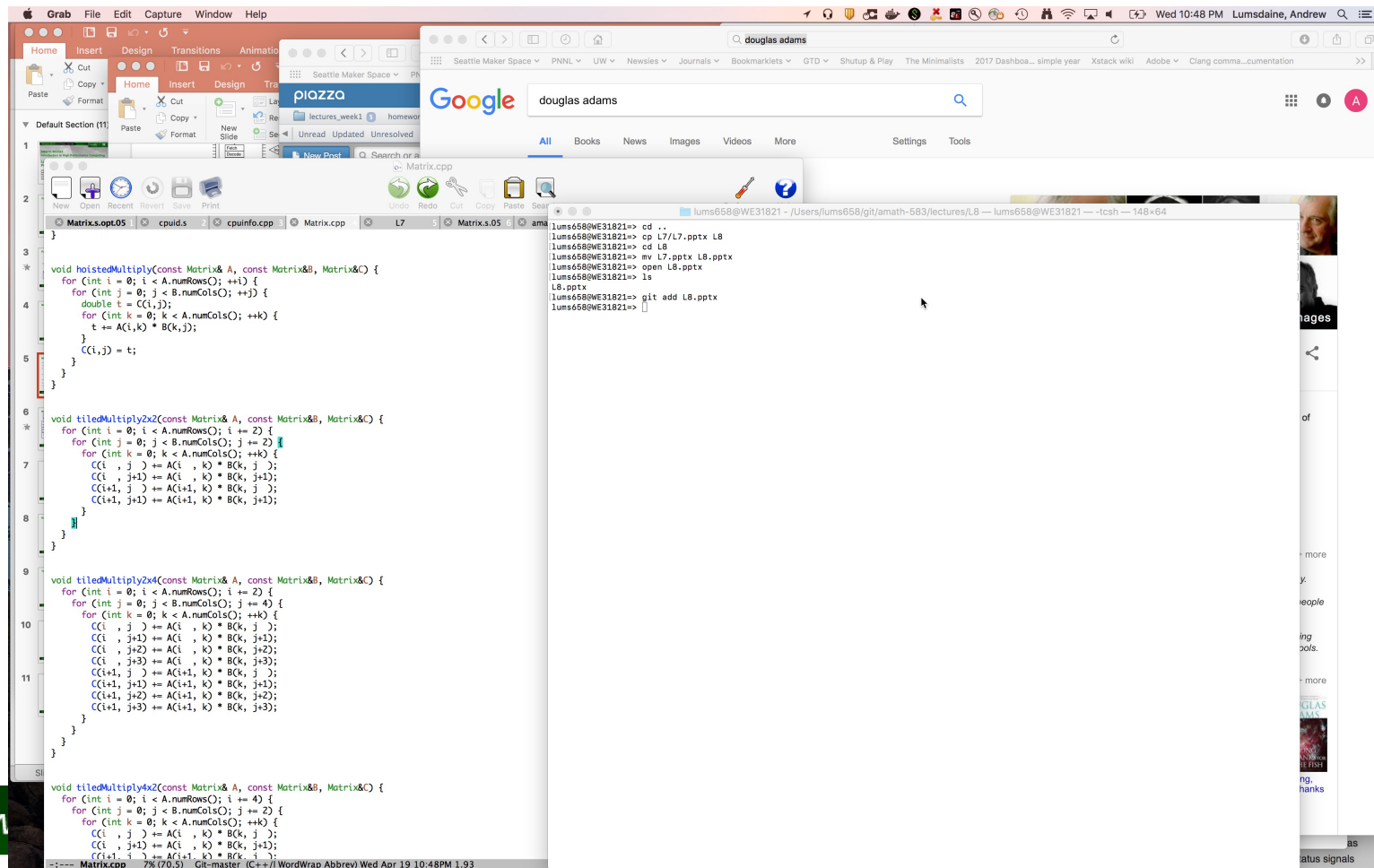


Compare And Contrast

- Time for everyone grades one exam
- Time for everyone grades one question

- How (why) did you use the approaches you did?

How Do We Run Many Programs at the Same Time?

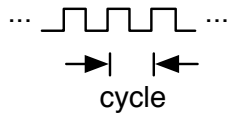


Running a Program

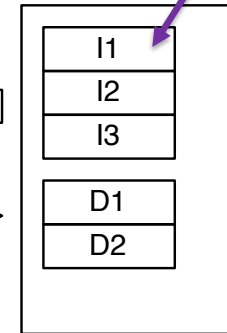
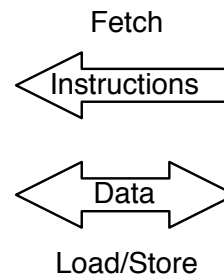
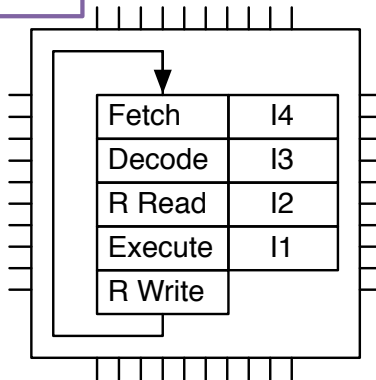
When a CPU is executing bytes from one program

Bytes from program stored in memory

It isn't executing bytes from another



Including from the OS (just another program)



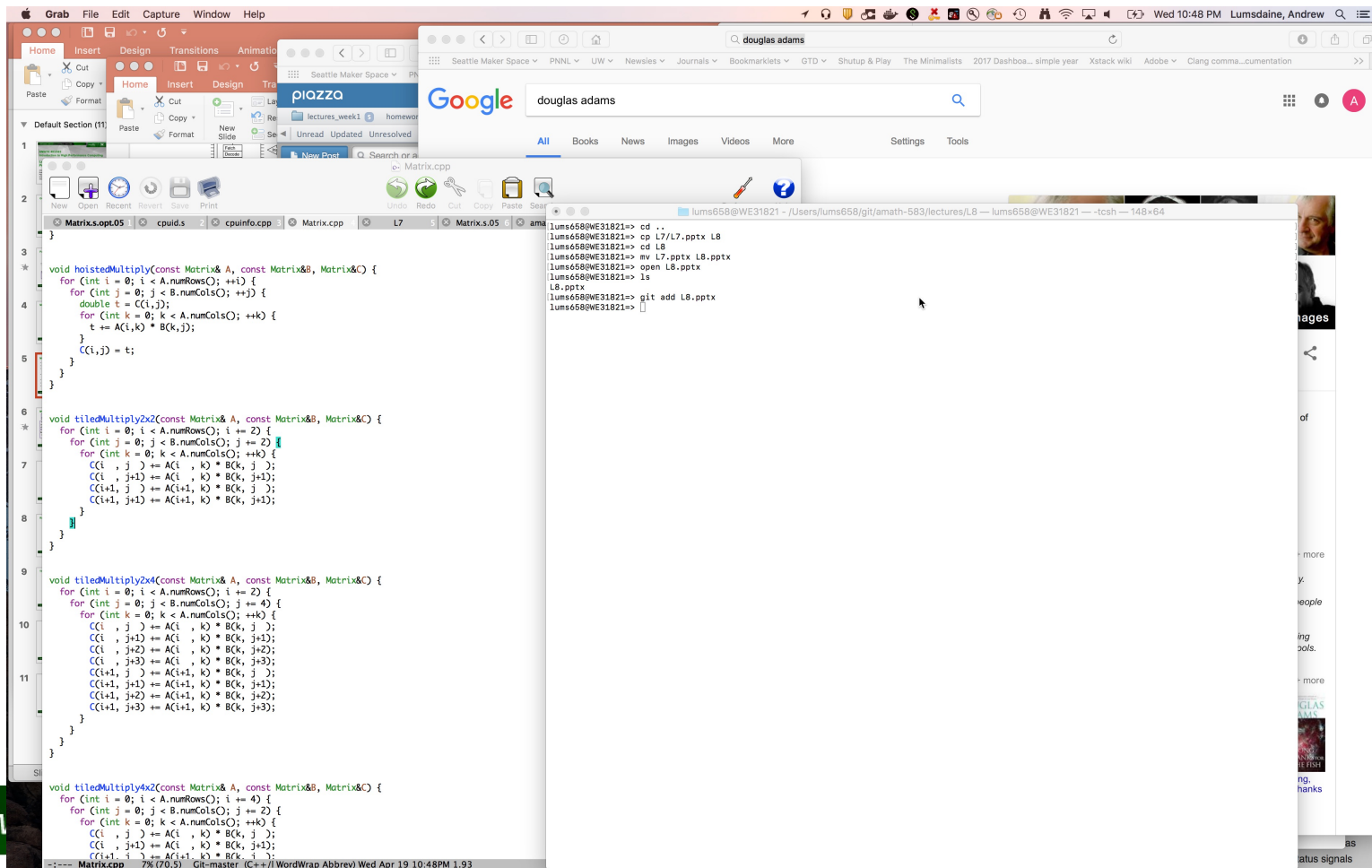
```

.global    __Z15hoistedMultiplyRKMatrixS1_RS_
.p2align  4, 0x30
__Z15hoistedMultiplyRKMatrixS1_RS_:    ## @__Z15hoistedMultiplyRKMatrixS1_RS_
.cfi_startproc
## BB#0:
pushq    %rbp
Ltmp16:  .cfi_def_cfa_offset 16
Ltmp17:  .cfi_offset 4;rbp, -16
movq    %rsp, %rbp
Ltmp18:  .cfi_def_cfa_register 4;rbp
pushq    %r15
pushq    %r14
pushq    %r13
pushq    %r12
pushq    %rbx
Ltmp19:  .cfi_offset 4;rbx, -56
Ltmp20:  .cfi_offset 4;rbx, -48
Ltmp21:  .cfi_offset 4;rbx, -40
Ltmp22:  .cfi_offset 4;rbx, -32
Ltmp23:  .cfi_offset 4;rbx, -24
movq    8(%rsi), %rax
testq   %rax, %rax    ## 8-byte Spill
je      LBB2_9
## BB#1:
movq    16(%rsi), %r12
movq    8(%rdx), %rax
movq    %rax, -104(%rbp)    ## 8-byte Spill
movq    16(%rdx), %rdx
movq    8(%rsi), %rax
movq    16(%rdi), %r13
leaq   -1(%rcx), %rsi    ## 8-byte Spill
movq    %rsi, -88(%rbp)
movl   %ecx, %esi
    
```

How does another program run?

How did the bytes get here?

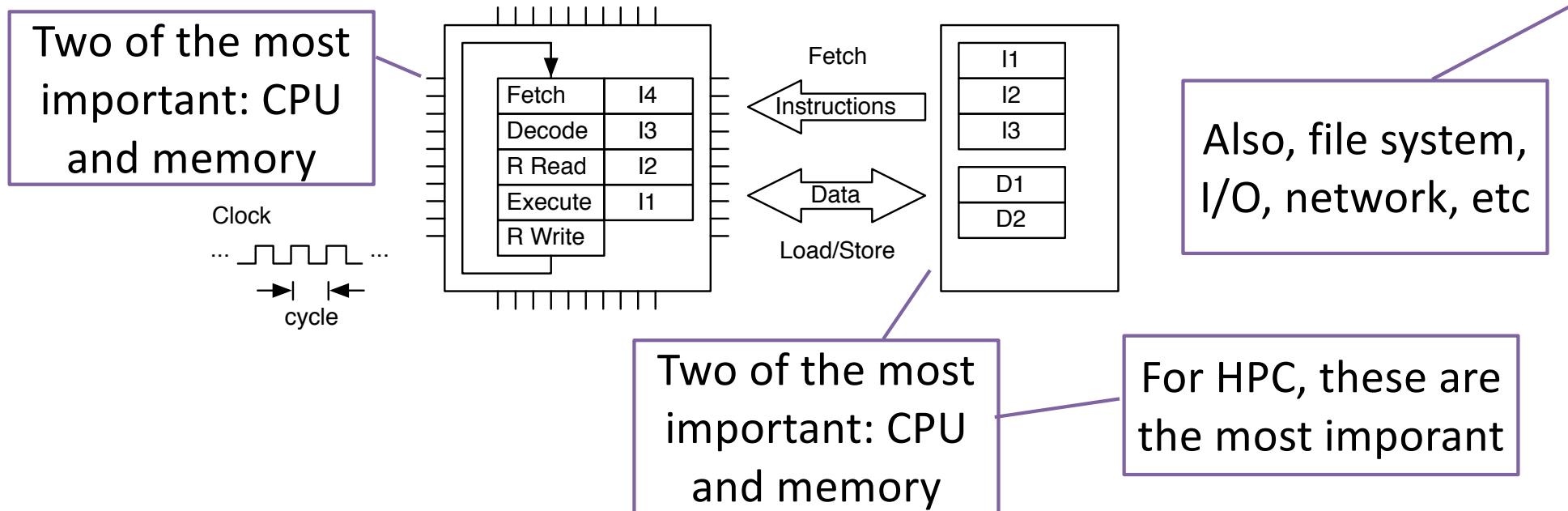
How Do We Run Many Programs at the Same Time?



NORTH

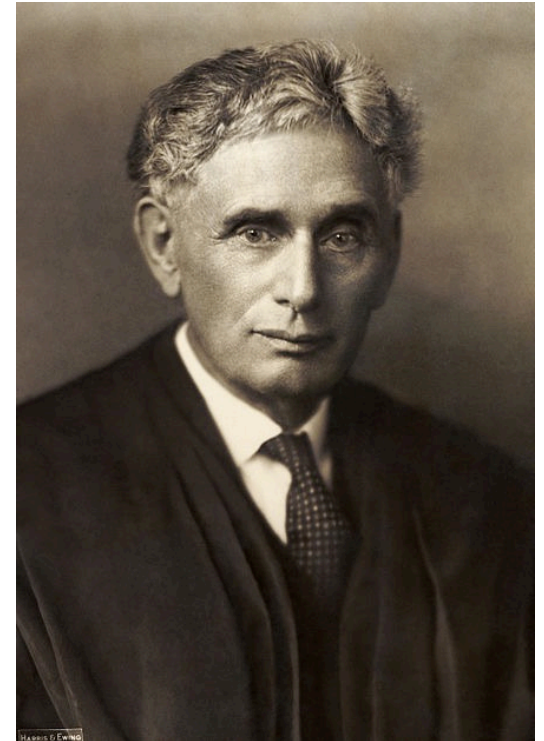
A Word About Operating Systems

- An operating system is **a program** that provides a standard interface between the resources of a computer and the users of the computer

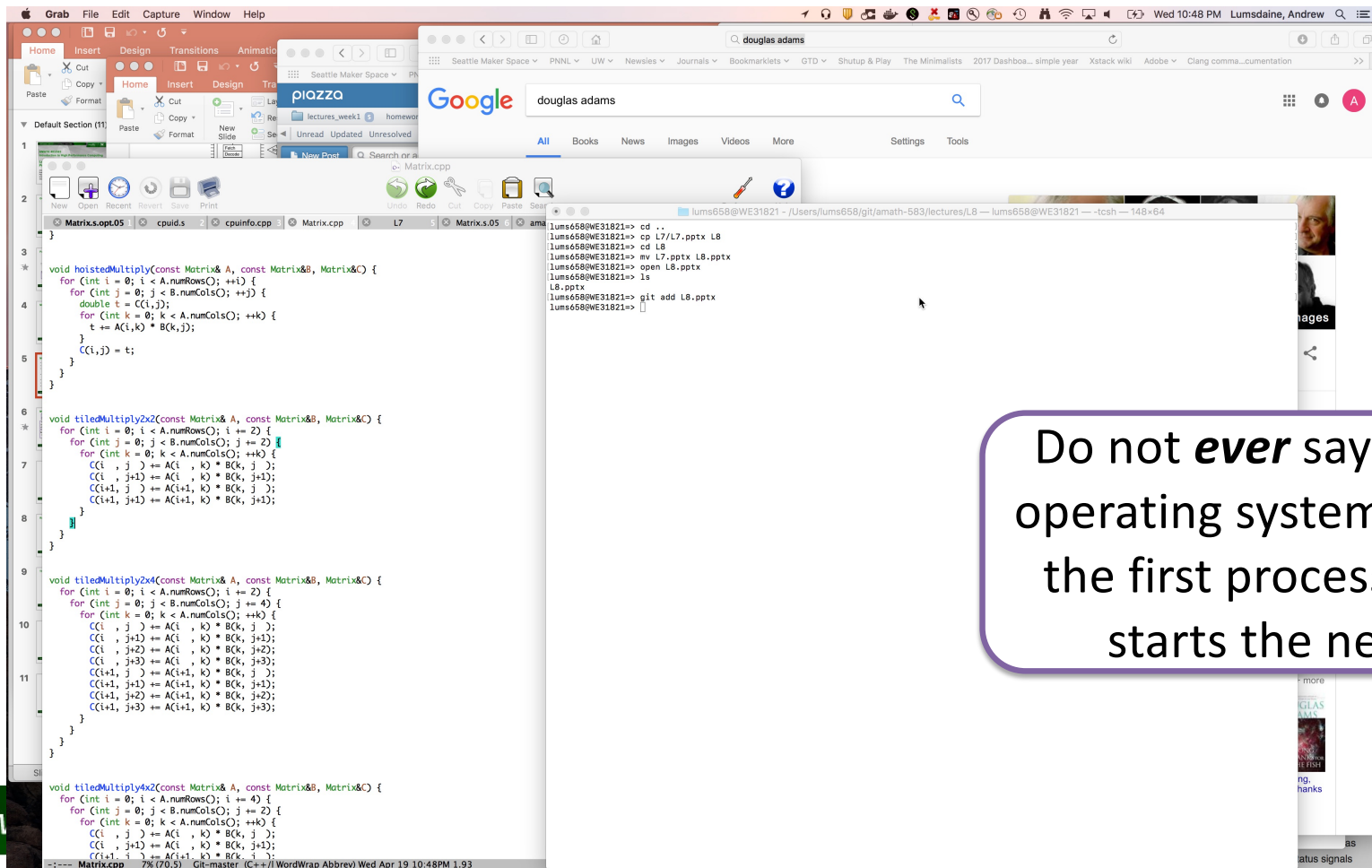


Processes and Threads

- A process is an abstraction for a collection of resources to represent a (running) program
 - CPU
 - Memory
 - Address space
- A thread is an abstraction of execution (using the resources within a process)
 - Can share an address space



How Do We Run Many Programs ~~at the same time?~~ ^{at the same time?}



Do not *ever* say: "the operating system stops the first process and starts the next"

The Operating System Can Run When...

- The process whose instructions are being executed by the CPU (the running process) requests a service from the OS (makes a **system call**)
- In response to a hardware interrupt
- It does not spontaneously run
- It is not somehow running in the background
- Again, when the CPU is executing instructions for one program, it is not executing instructions for another program
- The only way anything happens on the computer is if the CPU executes instructions that make it happen

Process Abstraction

Stored in Process Control Block (PCB)

Set of information about process resources

Sufficient to be able to start a process after stopped

Also for accounting / administrative purposes

Process management

Registers
 Program counter
 Program status word
 Stack pointer
 Process state
 Priority
 Scheduling parameters
 Process ID
 Parent process
 Process group
 Signals
 Time when process started
 CPU time used
 Children's CPU time
 Time of next alarm

Memory management

Pointer to text segment
 Pointer to data segment
 Pointer to stack segment

File management

Root directory
 Working directory
 File descriptors
 User ID
 Group ID

What does program counter represent?

The Process Concept

\$ top -u

Process ID

```
lums658@WE31821 - /Users/lums658
Processes: 419 total, 2 running, 417 sleeping, 1988 threads
Load Avg: 1.93, 1.88, 1.87 CPU usage: 3.45% user, 3.69% sys,
MemRegions: 156549 total, 7076M resident, 141M private, 3629M
VM: 4328G vsize, 627M framework vsize, 71344832(64) swapins,
Disks: 57070556/1524G read, 36025949/792G written.
```

PID	COMMAND	%CPU	TIME	#TH	#WQ	#PORT	MEM	PUR
0	kernel_task	12.6	29:59:12	177/9	0	2	1809M+	0B
114	hid	4.4	01:46:55	6	3	381+	3024K+	0B
8333	top	4.0	00:00.72	1/1	0	21	5016K	0B
8334	screencaptur	3.9	00:00.06	4	3	57	2500K+	20K
91791	LaTeXiT	2.3	09:45.97	6	2	255	42M	0B
67565	Terminal	2.0	01:50.53	13	8	346+	72M	0B
3288	Calendar	1.6	09:54.07	3	1	292	95M	185K
1234	com.docker.h	1.1	02:02:24	18	1	38	763M	0B
846	usernoted	1.1	03:13.97	5	4	139+	11M+	896K
83998	Slack Helper	1.0	01:40.81	19	2	149	189M+	0B
71742	splunkd	0.8	40:02.25	35	0	48	85M	0B
63334	Slack Helper	0.6	01:19.70	5	2	124	7780K	0B
184	mDNSResponde	0.5	22:51.68	5	1	103	5628K	0B
111-	NetworkMonit	0.4	12:37.75	28	27	49+	22M+	0B
883	CalNCService	0.3	19:18.74	5	3	182+	39M+	0B
853	SystemUIServ	0.2	02:45.35	5	3	371	33M+	28K-
63333	Slack	0.2	04:36.66	33	1	390	73M	0B

How much CPU

How many threads

PID	COMMAND	%CPU	TIME	#TH	#WQ	#PORT	MEM	PUR
162	WindowServer	13.8	07:48:22	6	2	702+	537M+	93M
0	kernel_task	12.6	29:59:12	177/9	0	2	1809M+	0B
4	hid	4.4	01:46:55	6	3	381+	3024K+	0B
33	top	4.0	00:00.72	1/1	0	21	5016K	0B
34	screencaptur	3.9	00:00.06	4	3	57	2500K+	20K
791	LaTeXiT	2.3	09:45.97	6	2	255	42M	0B
67565	Terminal	2.0	01:50.53	13	8	346+	72M	0B
3288	Calendar	1.6	09:54.07	3	1	292	95M	185K
1234	com.docker.h	1.1	02:02:24	18	1	38	763M	0B
846	usernoted	1.1	03:13.97	5	4	139+	11M+	896K
83998	Slack Helper	1.0	01:40.81	19	2	149	189M+	0B
71742	splunkd	0.8	40:02.25	35	0	48	85M	0B
63334	Slack Helper	0.6	01:19.70	5	2	124	7780K	0B
184	mDNSResponde	0.5	22:51.68	5	1	103	5628K	0B
111-	NetworkMonit	0.4	12:37.75	28	27	49+	22M+	0B
883	CalNCService	0.3	19:18.74	5	3	182+	39M+	0B
853	SystemUIServ	0.2	02:45.35	5	3	371	33M+	28K-
63333	Slack	0.2	04:36.66	33	1	390	73M	0B

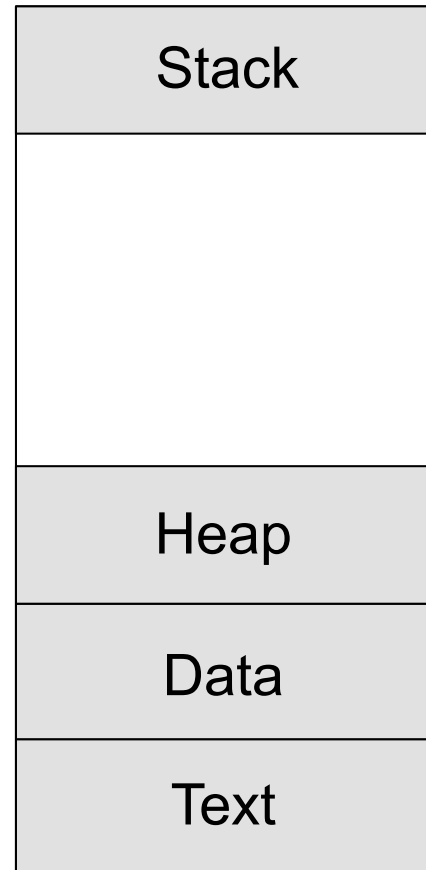
Process address space

Memory resources
for each process

All 32/48/64 bits

Address
Space

How can each
process use all the
address space?



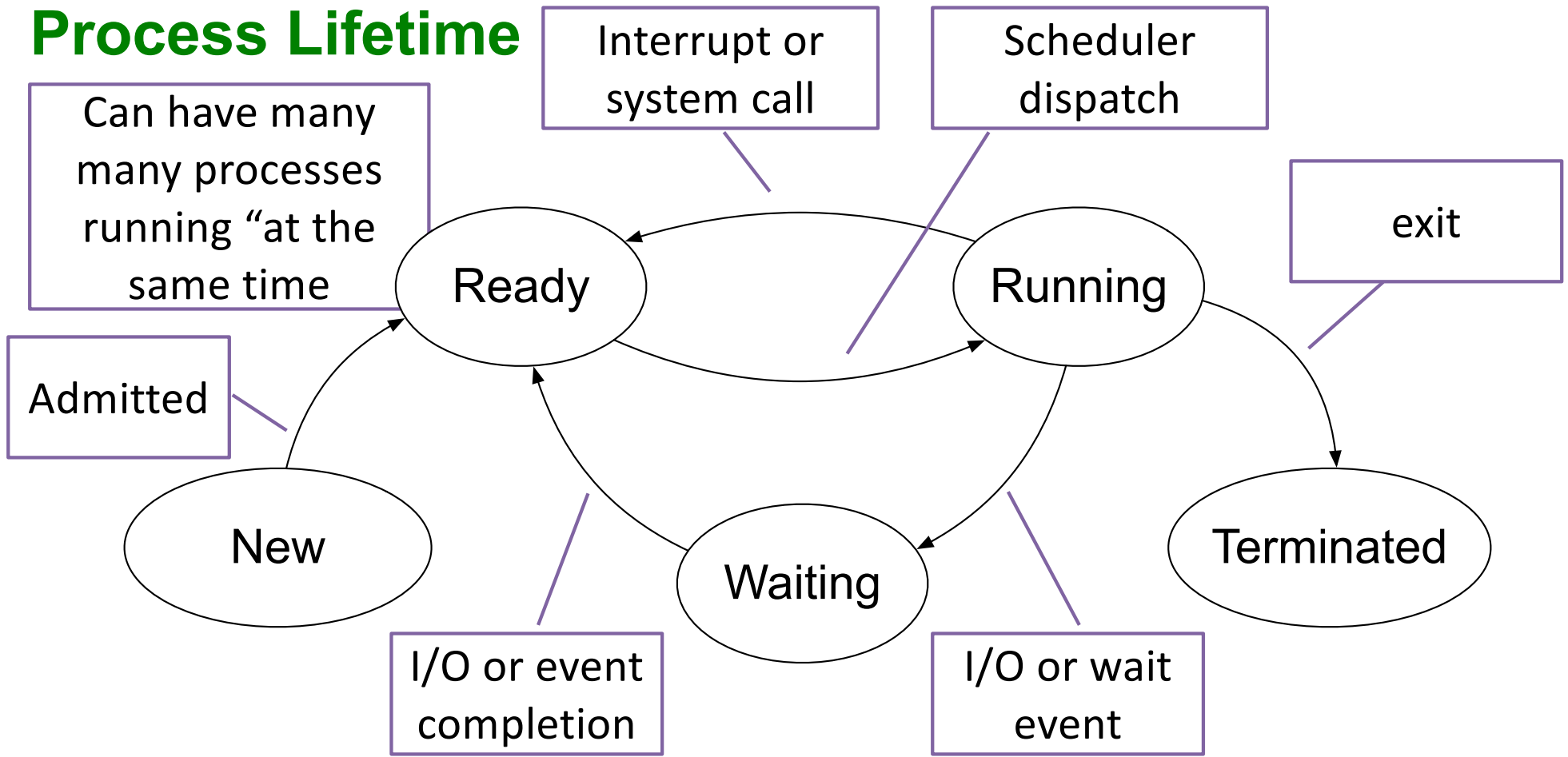
Created and
managed at run time

Created and
managed at run time

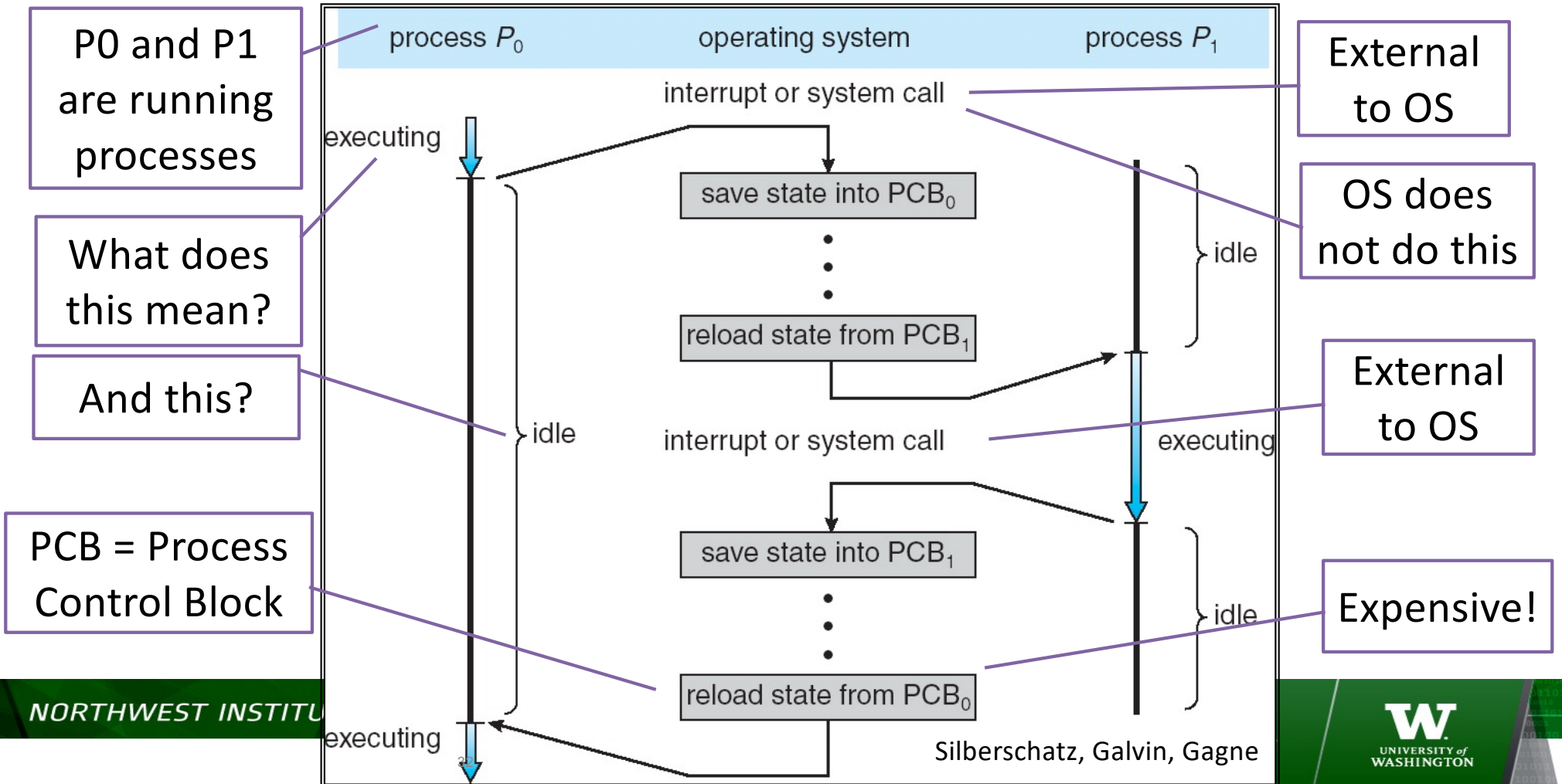
Compiled /
Linked

Stored
Program

Process Lifetime

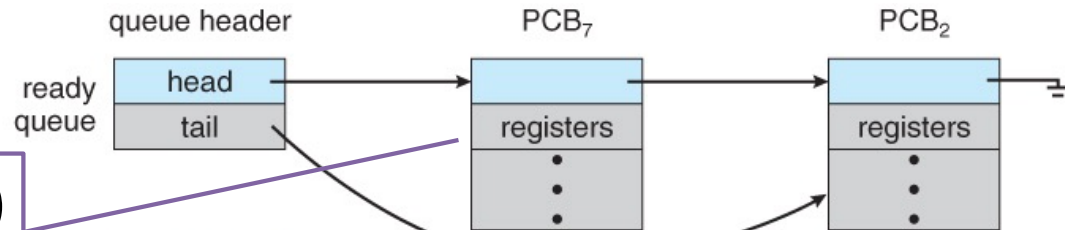


Context Switch

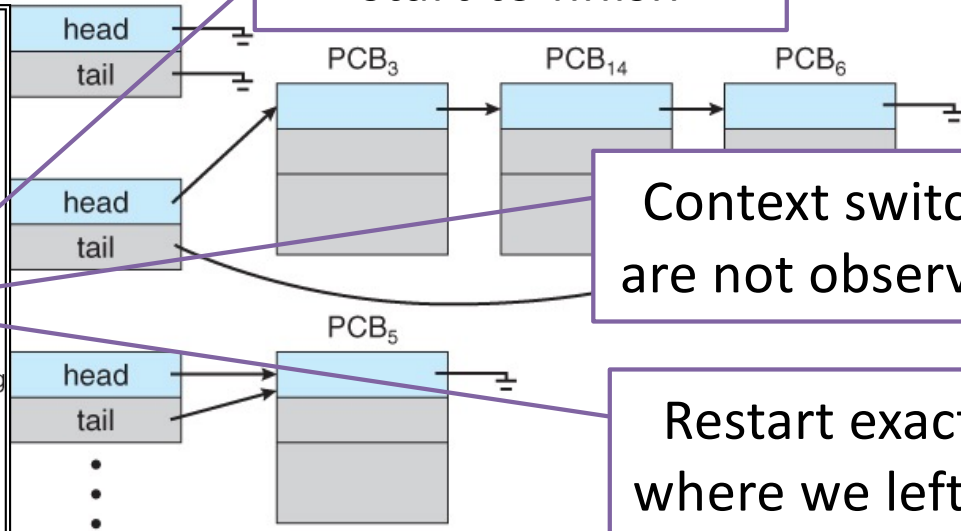
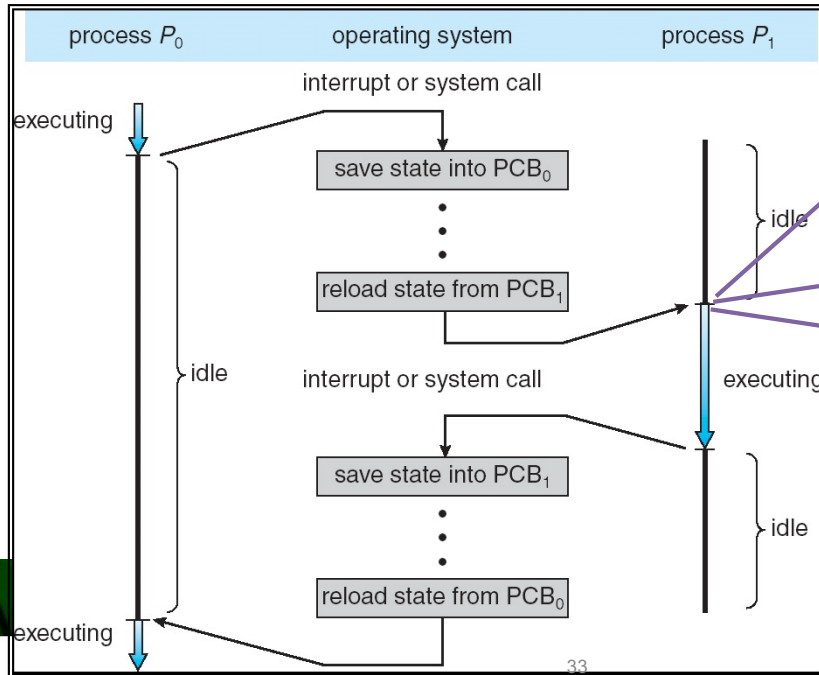


Process Queues

A process control block (PCB) has all information necessary to manage a process



Program runs from start to finish



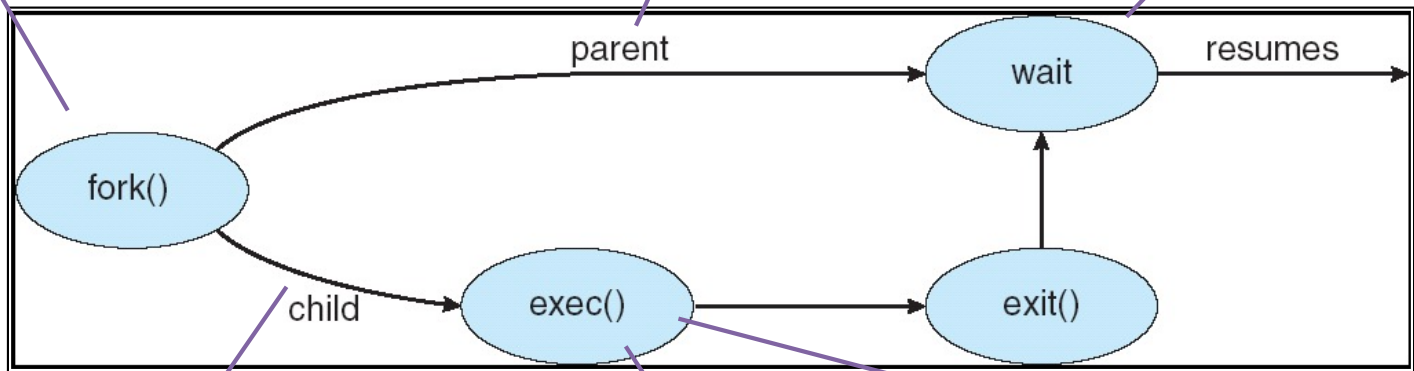
Context switches are not observable

Restart exactly where we left off

Process invokes fork()

The other process (the "parent") keeps executing

Can wait for other process to complete



The OS makes a copy of the original process and makes it runnable

One of the processes (the "child") runs exec()

Which pulls in new program bits to run

You see this fork/exec/wait almost all the time with one particular program you run (which?)

Example: process creation in UNIX

One process calls fork()

```
#include <unistd.h>

int main () {
    fork();

    return 0;
}
```

Each process "thinks" it called fork() and returned

Two processes return from fork()

```
#include <unistd.h>

int main () {
    fork();
    return 0;
}
```

Two processes return from fork()

```
#include <unistd.h>

int main () {
    fork();
    return 0;
}
```

fork() make an exact copy

Example

```
int main() {  
  {  
    int pids[20];  
  
    for (int i = 0; i < 20; ++i) {  
      pids[i] = fork();  
    }  
  
    return 0;  
  }  
}
```

fork() returns a
PID identifier

Loop 20 times

Call fork() 20
times

How many processes
get created?

Example

How deep is the tree?

$i == 0$

```
int main() {  
  {  
    int pids[20];  
  
    for (int i = 0; i < 20; ++i) {  
      pids[i] = fork();  
    }  
  
    return 0;  
  }  
}
```

Don't do this (ever)!

How many processes?

$i == 1$

```
int main() {  
  {  
    int pids[20];  
  
    for (int i = 0; i < 20; ++i) {  
      pids[i] = fork();  
    }  
  
    return 0;  
  }  
}
```

```
int main() {  
  {  
    int pids[20];  
  
    for (int i = 0; i < 20; ++i) {  
      pids[i] = fork();  
    }  
  
    return 0;  
  }  
}
```

$i == 2$

```
int main() {  
  {  
    int pids[20];  
  
    for (int i = 0; i < 20; ++i) {  
      pids[i] = fork();  
    }  
  
    return 0;  
  }  
}
```

```
int main() {  
  {  
    int pids[20];  
  
    for (int i = 0; i < 20; ++i) {  
      pids[i] = fork();  
    }  
  
    return 0;  
  }  
}
```

```
int main() {  
  {  
    int pids[20];  
  
    for (int i = 0; i < 20; ++i) {  
      pids[i] = fork();  
    }  
  
    return 0;  
  }  
}
```

```
int main() {  
  {  
    int pids[20];  
  
    for (int i = 0; i < 20; ++i) {  
      pids[i] = fork();  
    }  
  
    return 0;  
  }  
}
```

man fork()

```
#include <unistd.h>
pid_t fork();
```

The child process has a unique id

Upon successful completion, fork() returns a value of 0 to the child process and the returns the process ID of the child process to the parent process

```
lums658@WE31821 - /Users/lums658/git/amath-583/lectures/L8 — lums658@WE31821 — less · man fork — 135x52
FORK(2)                                BSD System Calls Manual                                FORK(2)
NAME
    fork -- create a new process
SYNOPSIS
    #include <unistd.h>
    pid_t
    fork(void);
DESCRIPTION
    fork() causes creation of a new process. The new process (child process) is an exact copy of the calling process (parent process) except for the following:
    o The child process has a unique process ID.
    o The child process has a different parent process ID (i.e., the process ID of the parent process).
    o The child process has its own copy of the parent's descriptors. These descriptors reference the same underlying objects, so that, for instance, file pointers in file objects are shared between the child and the parent, so that an lseek(2) on a descriptor in the child process can affect a subsequent read or write by the parent. This descriptor copying is also used by the shell to establish standard input and output for newly created processes as well as to set up pipes.
    o The child processes resource utilizations are set to 0; see setrlimit(2).
RETURN VALUES
    Upon successful completion, fork() returns a value of 0 to the child process and returns the process ID of the child process to the parent process. Otherwise, a value of -1 is returned to the parent process, no child process is created, and the global variable errno is set to indicate the error.
ERRORS
    fork() will fail and no child process will be created if:
    [EAGAIN]      The system-imposed limit on the total number of processes under execution would be exceeded. This limit is configuration-dependent.
    [EAGAIN]      The system-imposed limit MAXUPRC (<sys/param.h>) on the total number of processes under execution by a single user would be exceeded.
    [ENOMEM]      There is insufficient swap space for the new process.
SEE ALSO
    #include <sys/types.h>
    #include <unistd.h>
    The include file <sys/types.h> is necessary.
SEE ALSO
```


Example Revisited

```
int main() {  
  {  
    pid_t pids[20];  
  
    for (int i = 0; i < 20; ++i) {  
      pids[i] = fork();  
      if (pids[i] == 0)  
        break;  
    }  
  
    return 0;  
  }  
}
```

Get return
value of fork()

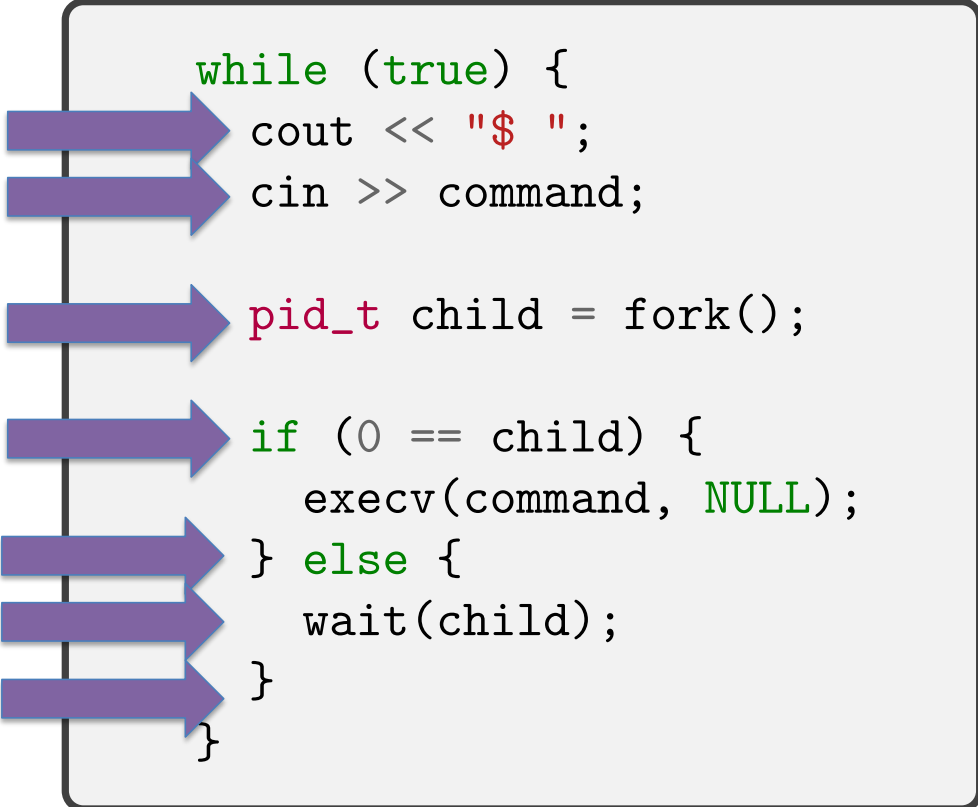
How many
processes
now?

If zero, the
process is a child

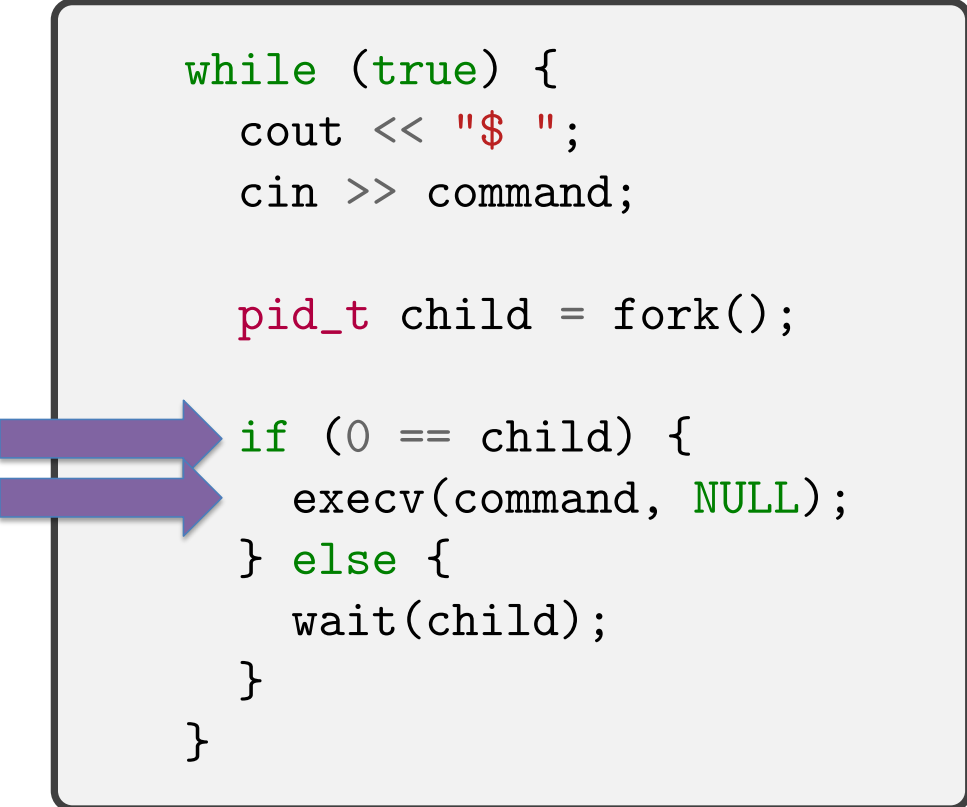
If no, the process
is the parent,
keep going

Process creation in UNIX (fork / exec pattern)

```
while (true) {  
    cout << "$ ";  
    cin >> command;  
    pid_t child = fork();  
    if (0 == child) {  
        execv(command, NULL);  
    } else {  
        wait(child);  
    }  
}
```



```
while (true) {  
    cout << "$ ";  
    cin >> command;  
    pid_t child = fork();  
    if (0 == child) {  
        execv(command, NULL);  
    } else {  
        wait(child);  
    }  
}
```



How Do We Run Multiple Programs Concurrently?

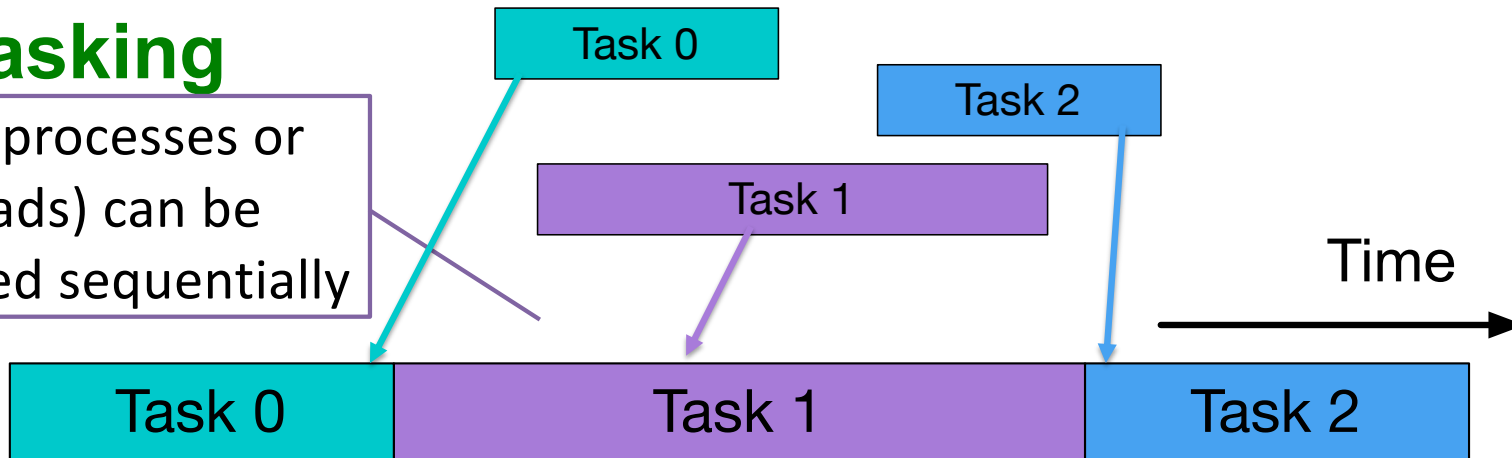
The screenshot displays a macOS desktop environment. In the foreground, a code editor window shows C++ code for matrix multiplication. The code includes a function `hoistedMultiply` and three tiled multiplication functions: `tiledMultiply2x2`, `tiledMultiply2x4`, and `tiledMultiply4x2`. The `hoistedMultiply` function uses a standard row-major loop. The tiled functions use nested loops to process the matrix in blocks of size 2, 4, and 4 respectively. A terminal window in the background shows the following commands and output:

```
lums658@WE31821:~$ cd ..
lums658@WE31821:~/lums658$ cp L7/L7.pptx L8
lums658@WE31821:~/lums658$ cd L8
lums658@WE31821:~/lums658/L8$ mv L7.pptx L8.pptx
lums658@WE31821:~/lums658/L8$ open L8.pptx
lums658@WE31821:~/lums658/L8$ ls
L8.pptx
lums658@WE31821:~/lums658/L8$ git add L8.pptx
lums658@WE31821:~/lums658/L8$
```

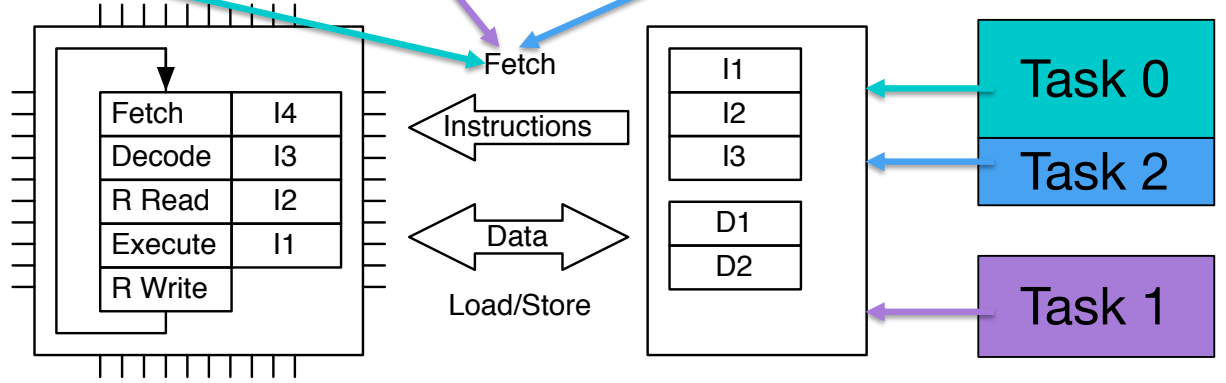
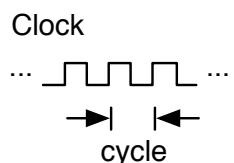
A web browser window in the background shows a Google search for "douglas adams". The search results are partially visible, showing a profile picture and some text.

Multitasking

Tasks (processes or threads) can be scheduled sequentially



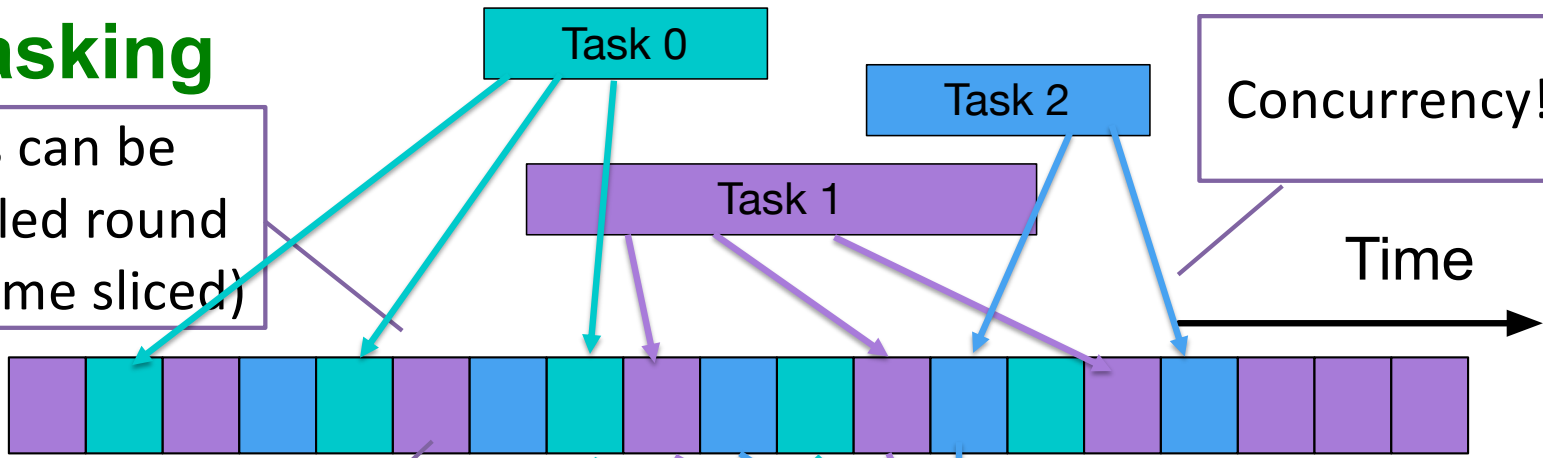
Run to completion



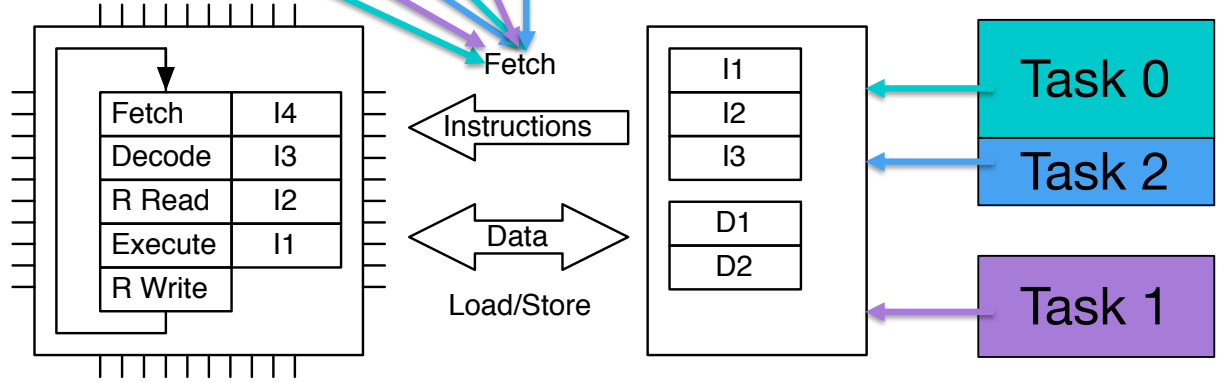
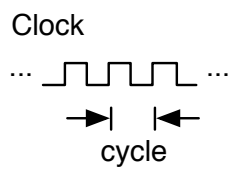
Multitasking

Tasks can be scheduled round robin (time sliced)

Concurrency!



Run to context switch (system call or interrupt)



Multitasking on Multicore

Concurrency!

Time sliced and mapped to separate cores

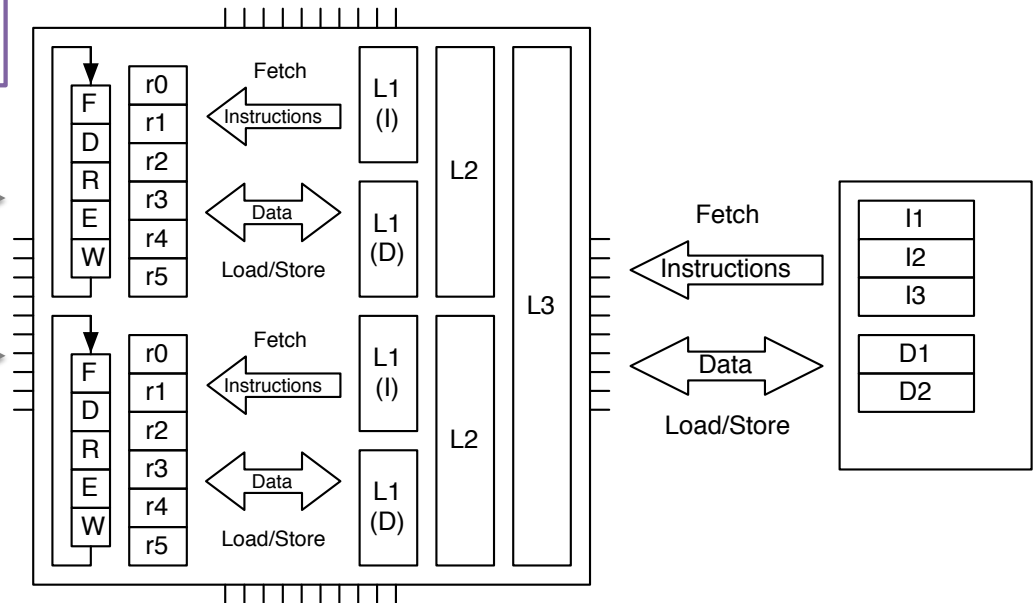
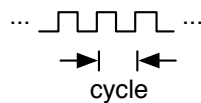
A single threaded task can only use one core at a time



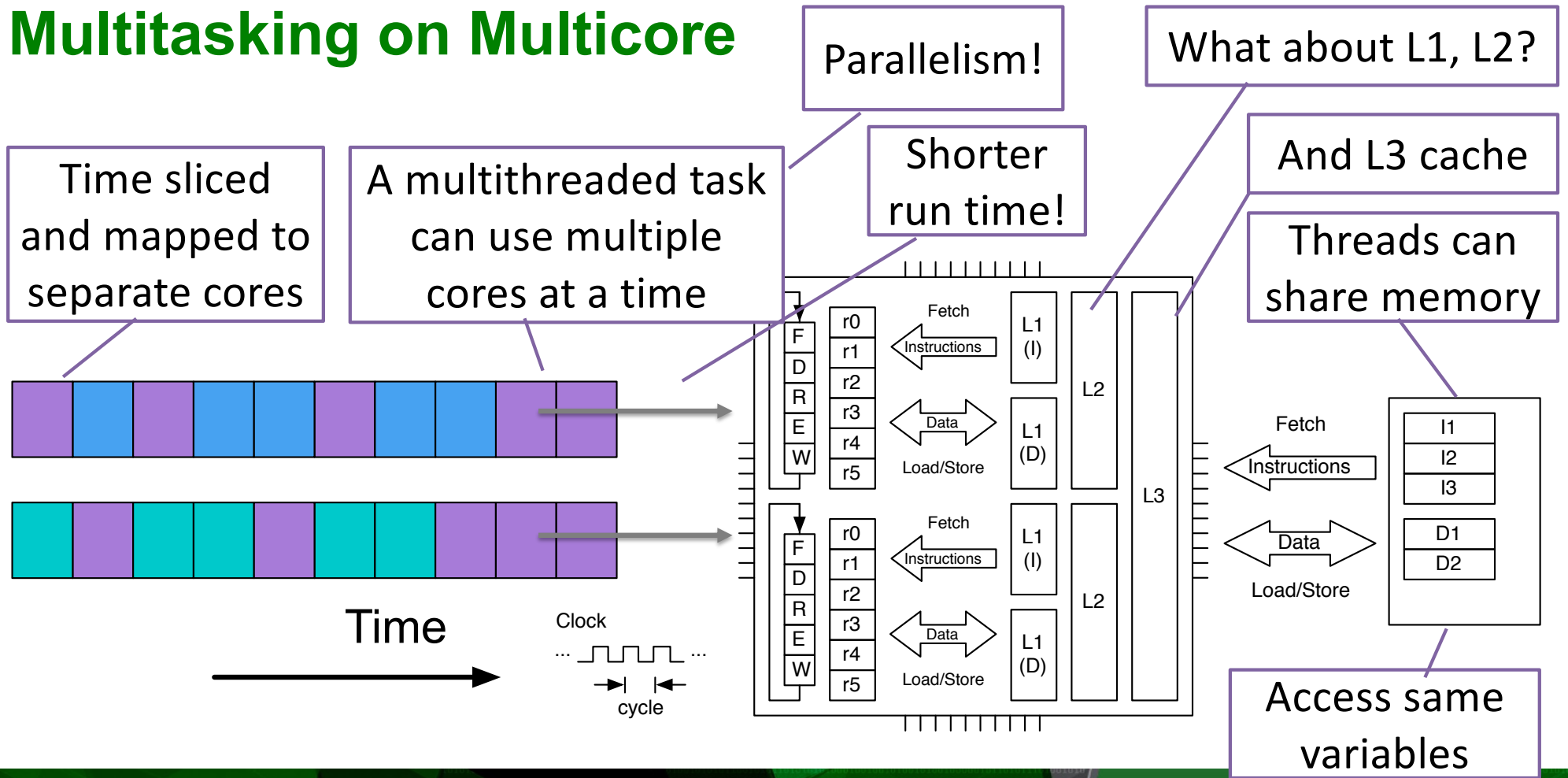
Time



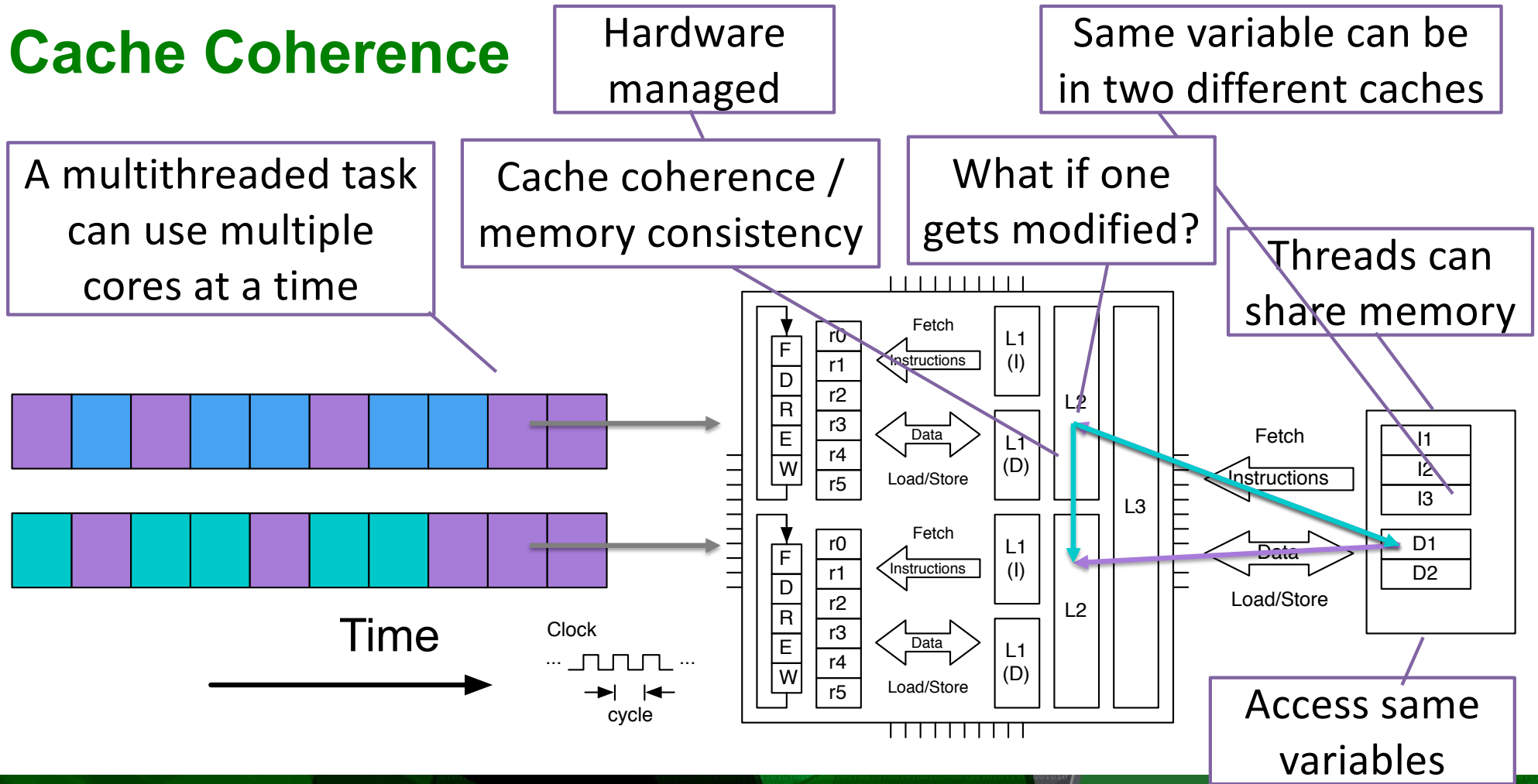
Clock



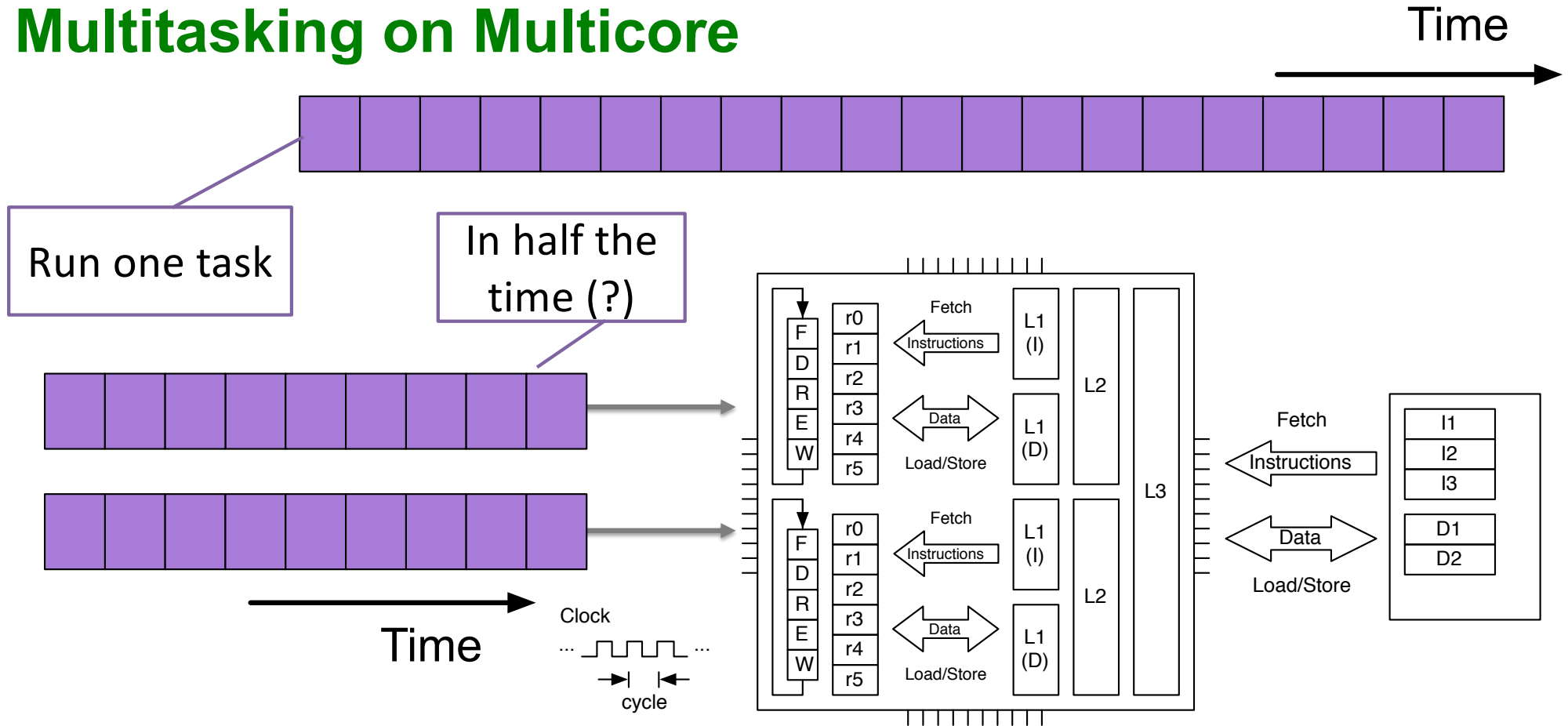
Multitasking on Multicore



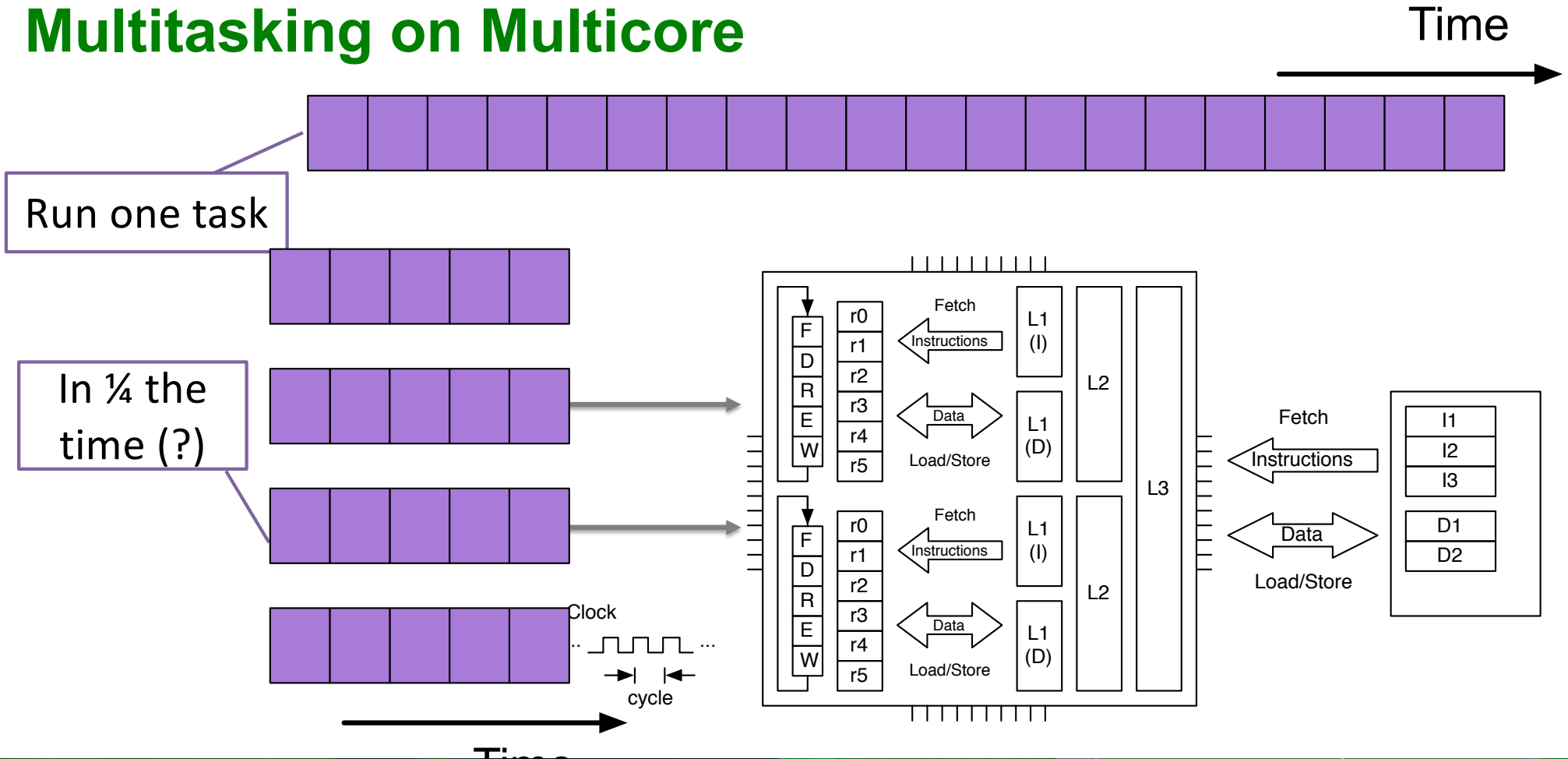
Cache Coherence



Multitasking on Multicore



Multitasking on Multicore



Multitasking on Multicore

Nonetheless, this is the essence of *parallel* computing

Parallel computation isn't done until all cores are done

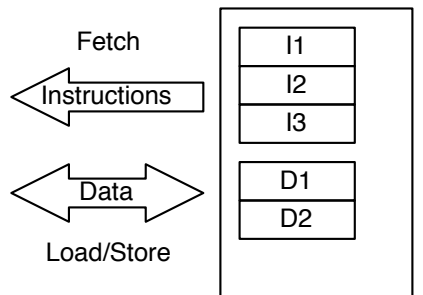
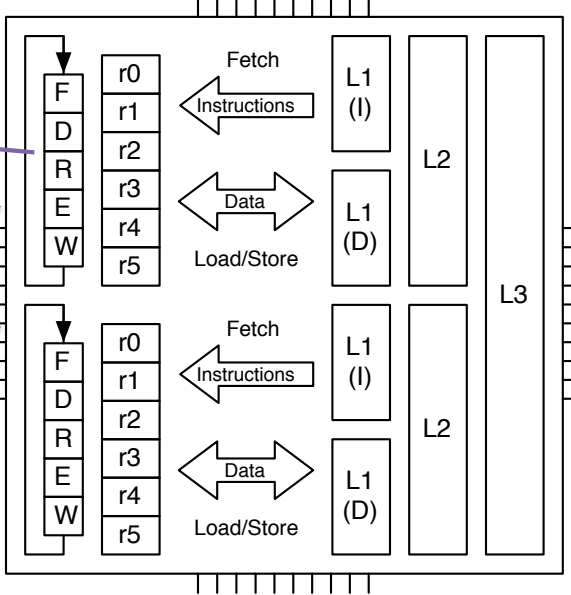
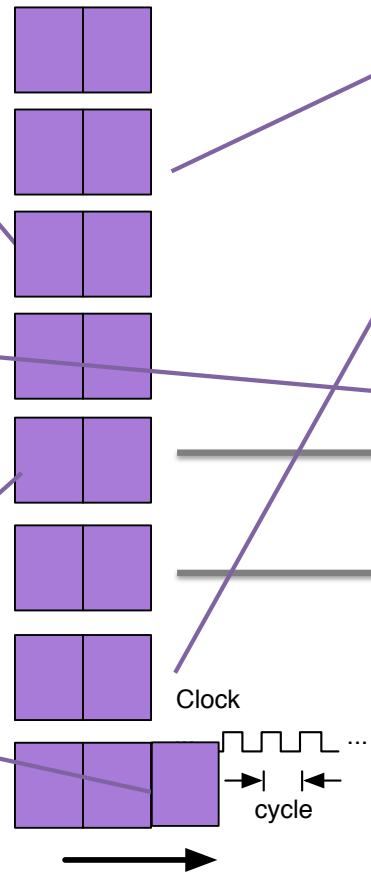
Not the same as concurrent

In 1/8 the time (?)

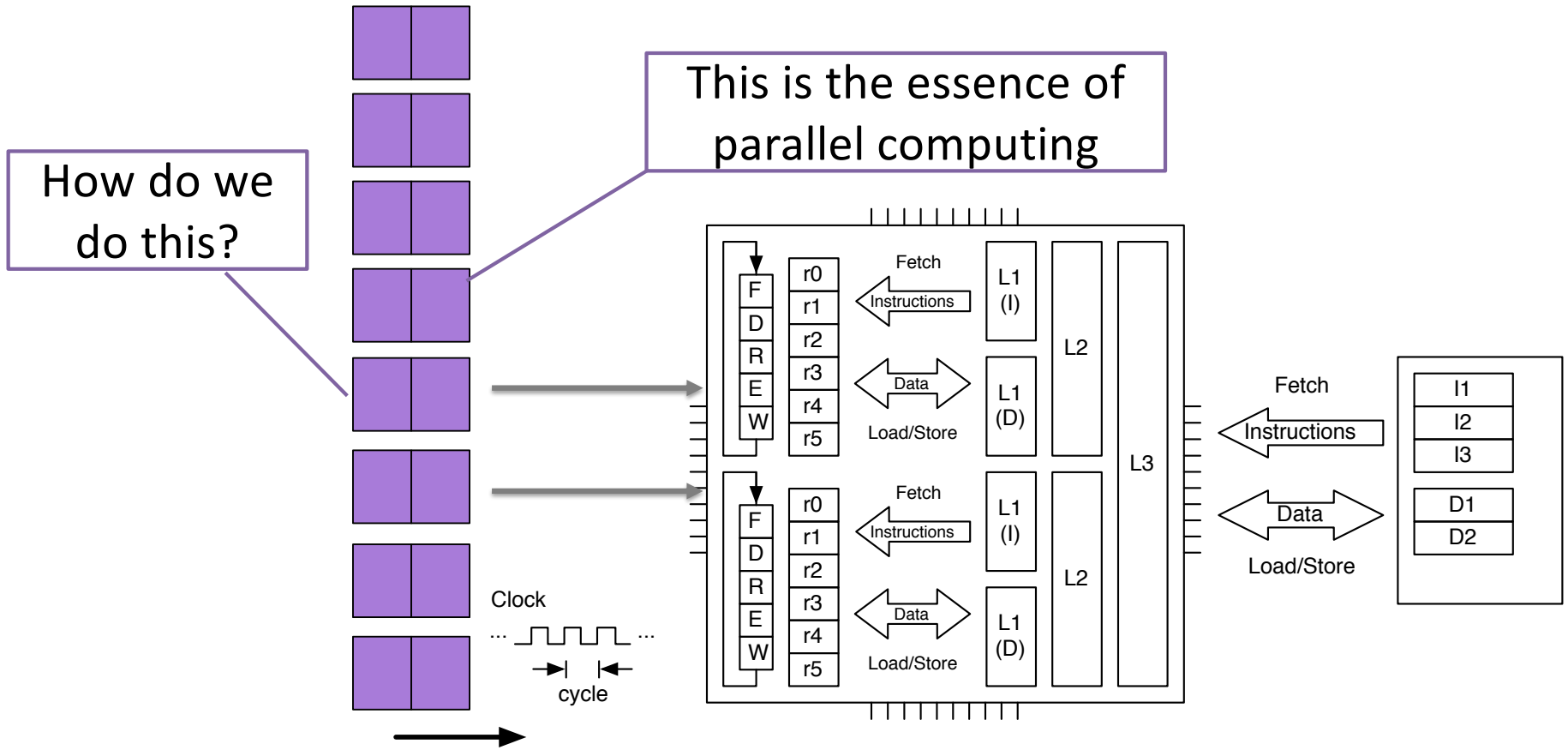
Need enough cores (8)

Work needs to be balanced

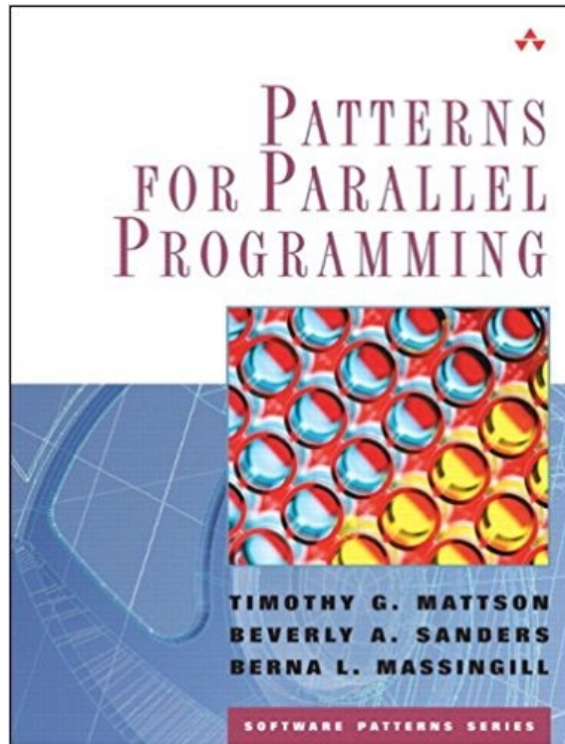
oops



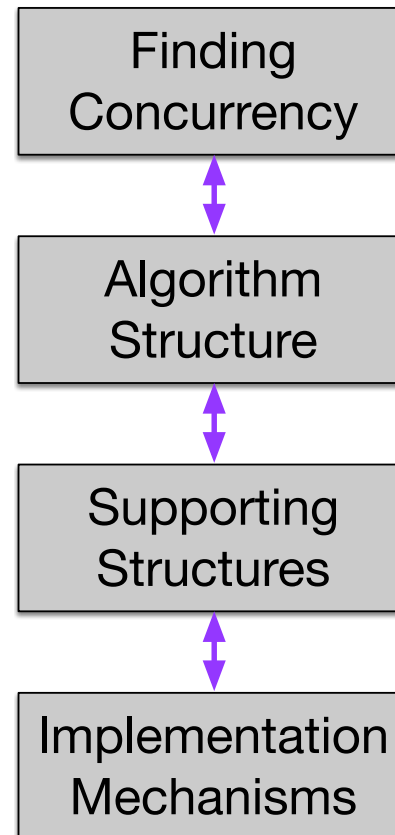
Multitasking on Multicore



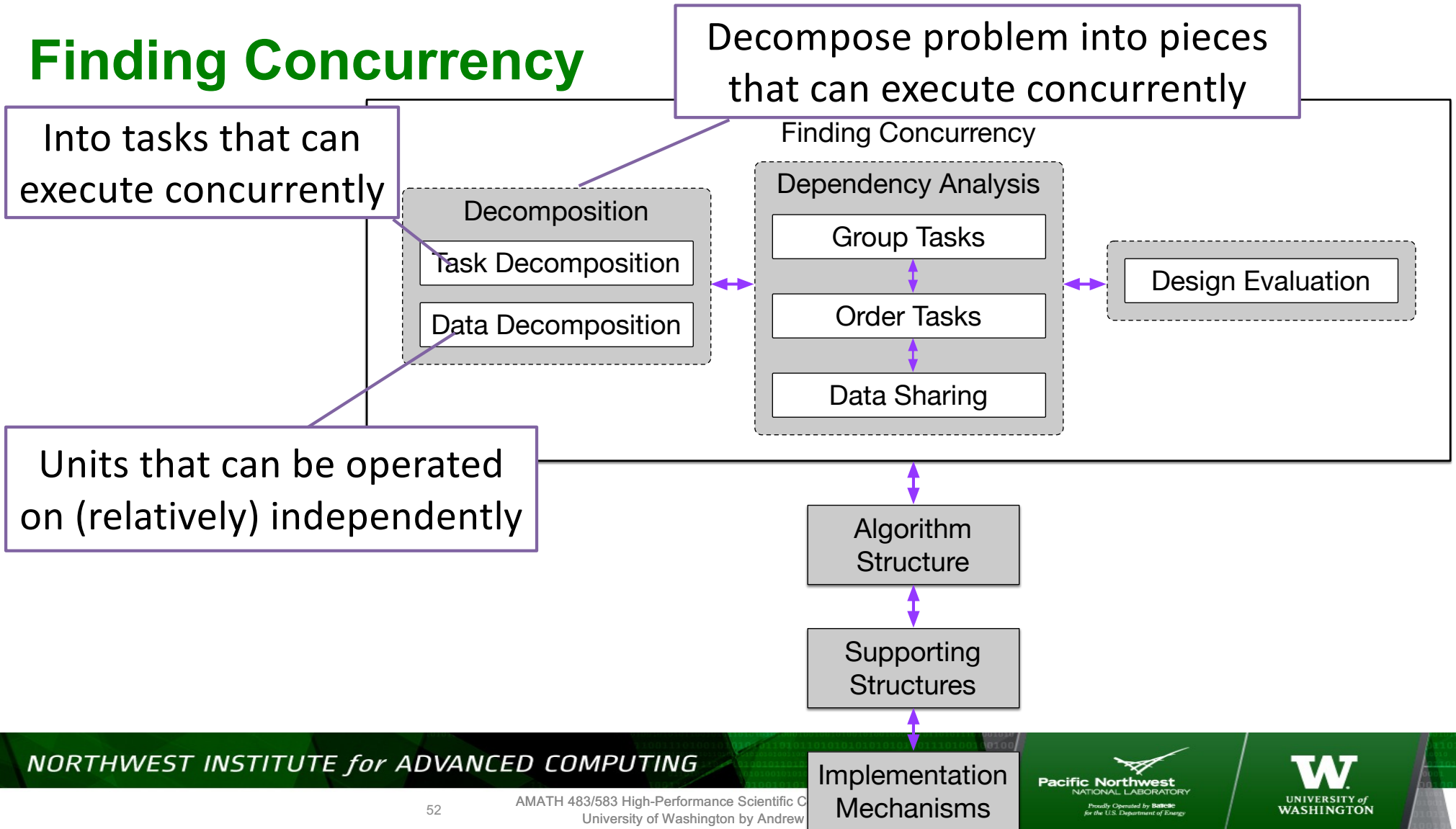
Parallelization Strategy



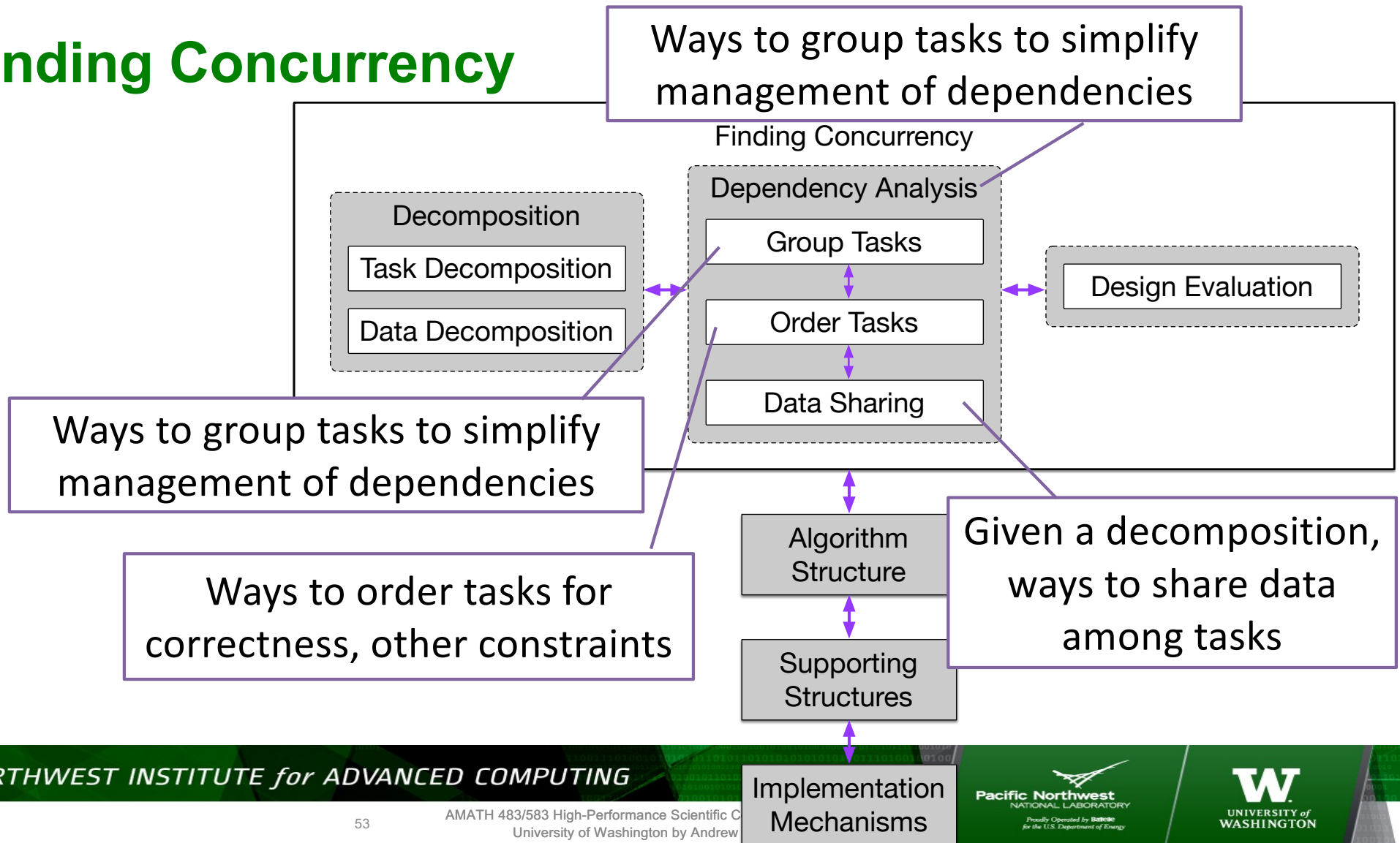
Timothy Mattson, Beverly Sanders, and Berna Massingill.
2004. *Patterns for Parallel Programming*(First ed.). Addison-
Wesley Professional.



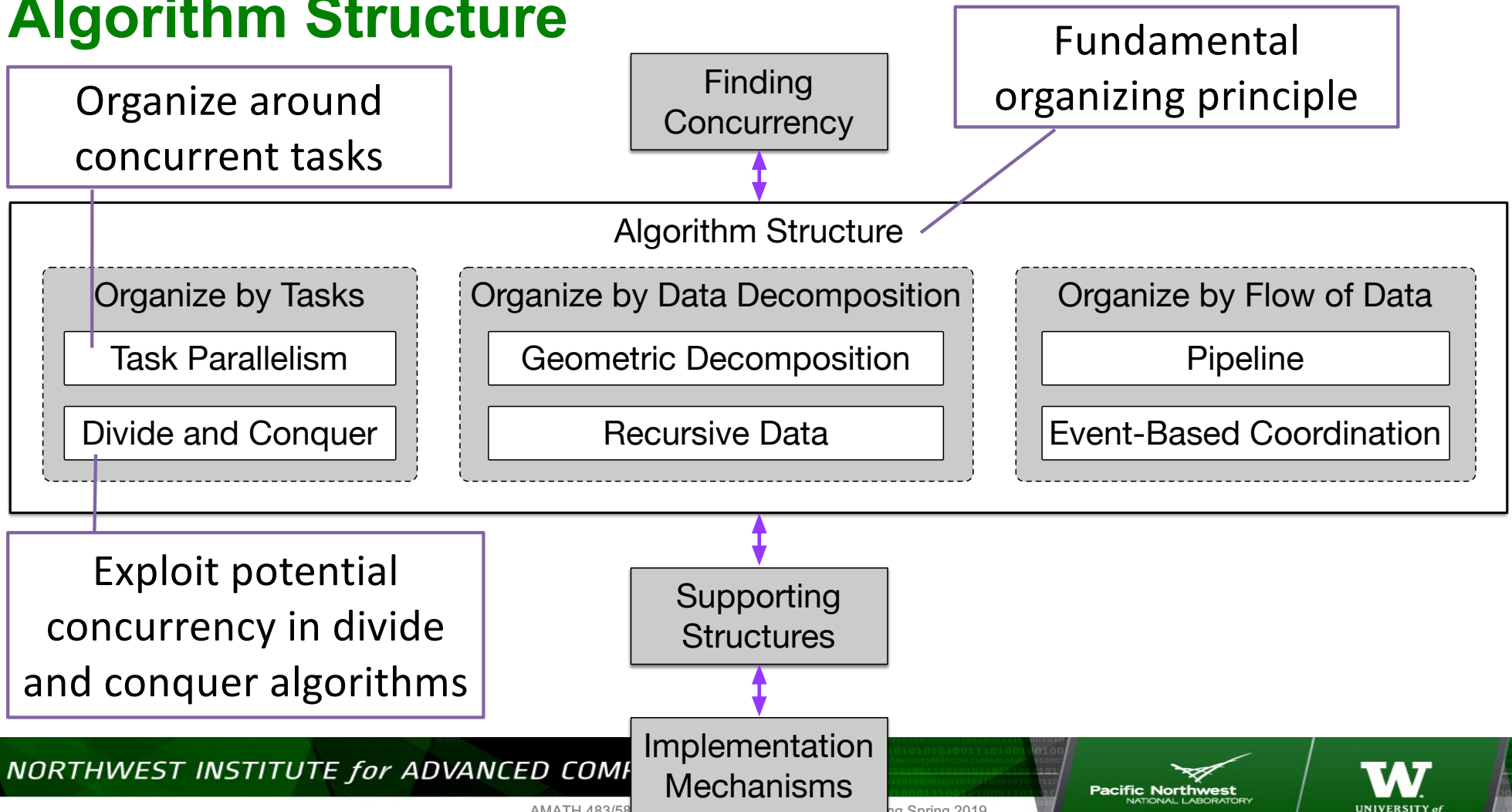
Finding Concurrency



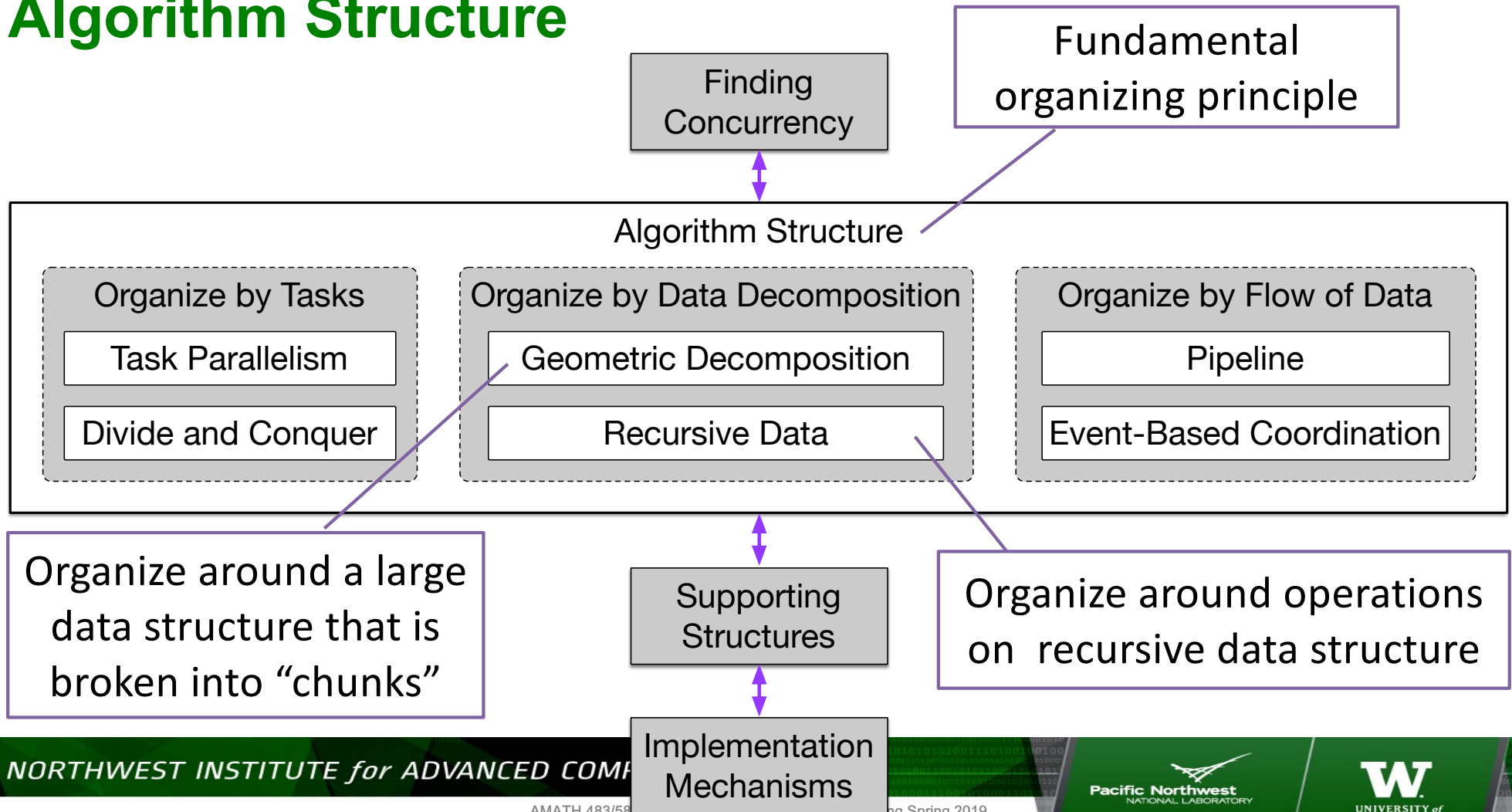
Finding Concurrency



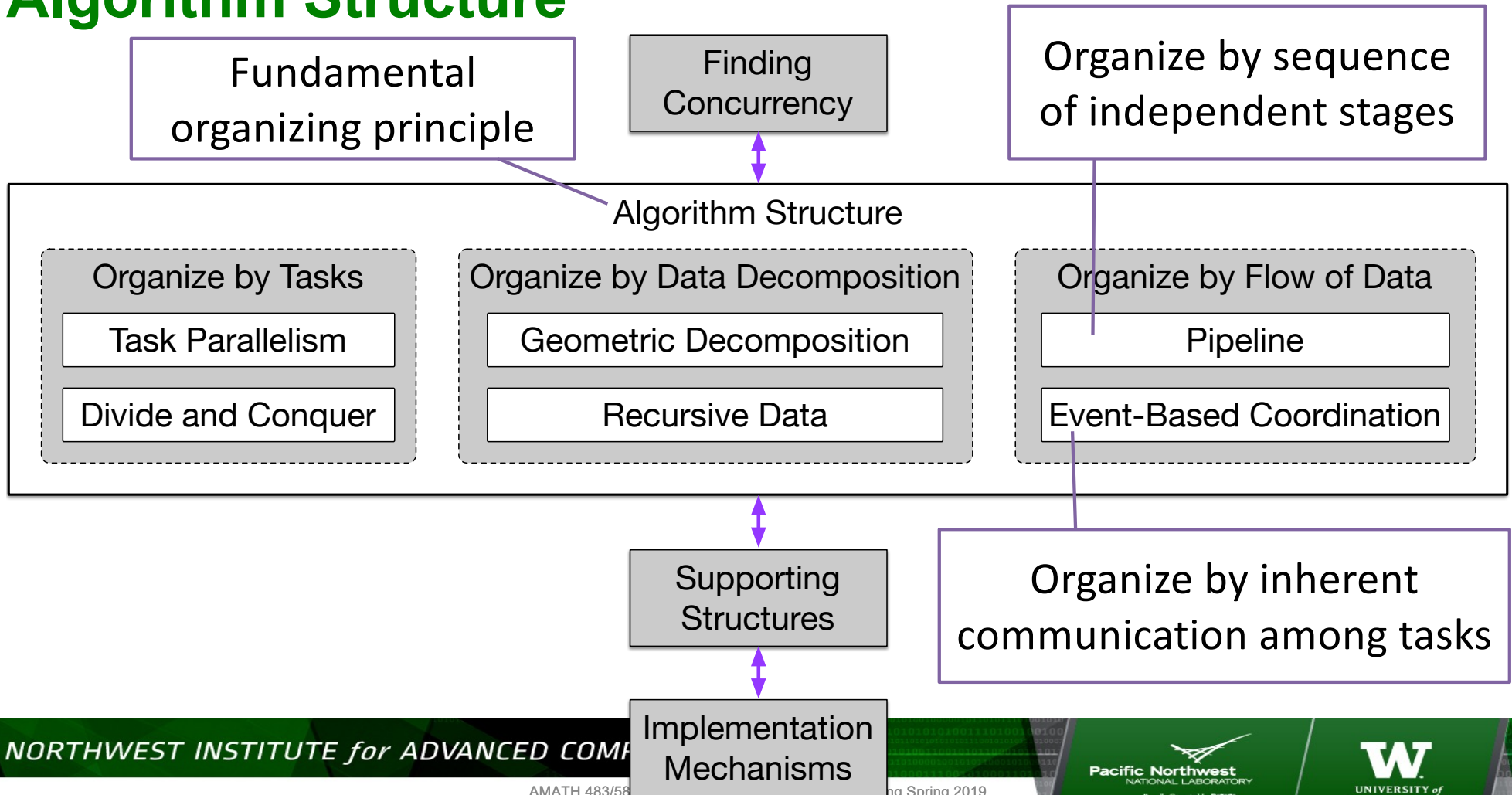
Algorithm Structure



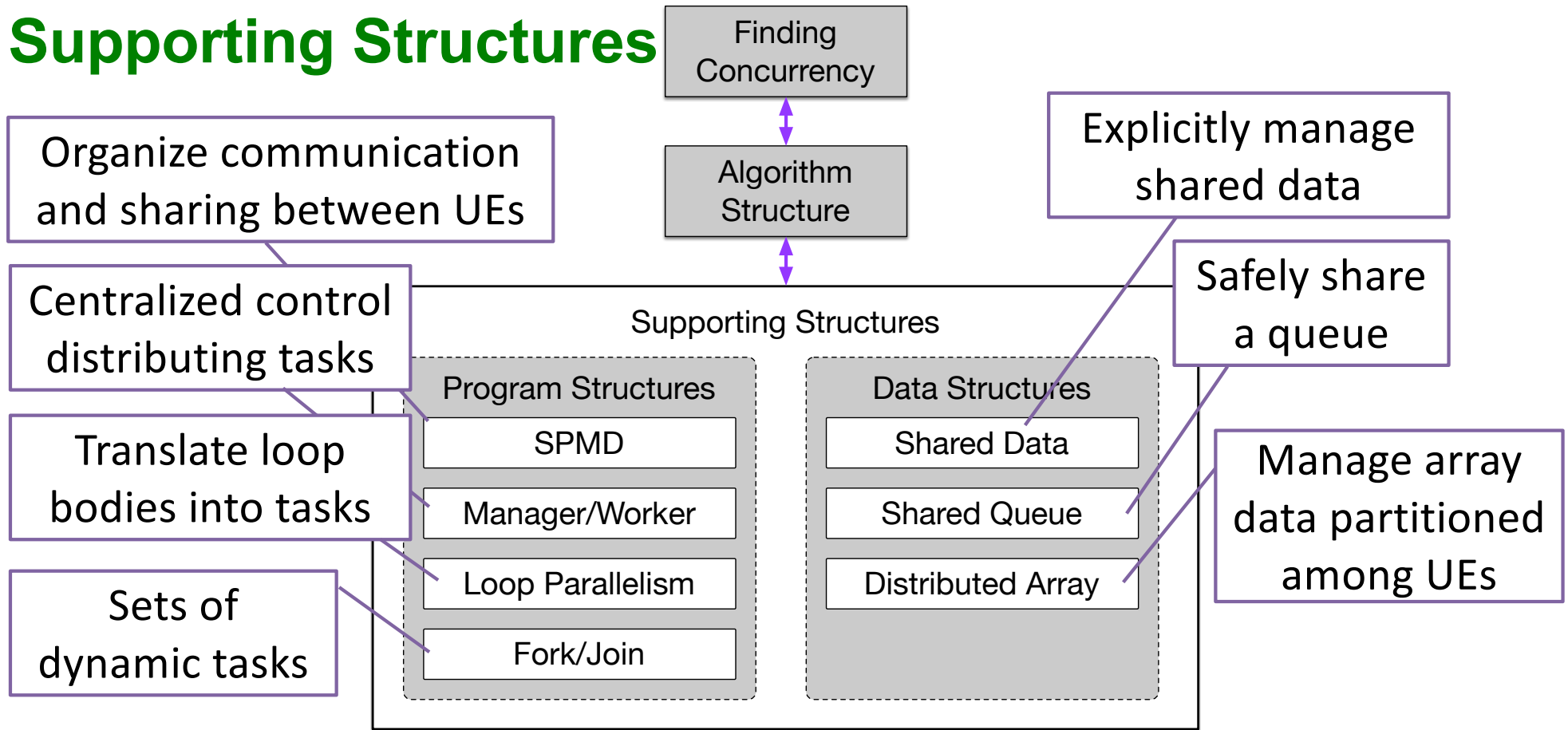
Algorithm Structure



Algorithm Structure



Supporting Structures



Implementation Mechanisms

Implementation Mechanisms

Finding
Concurrency

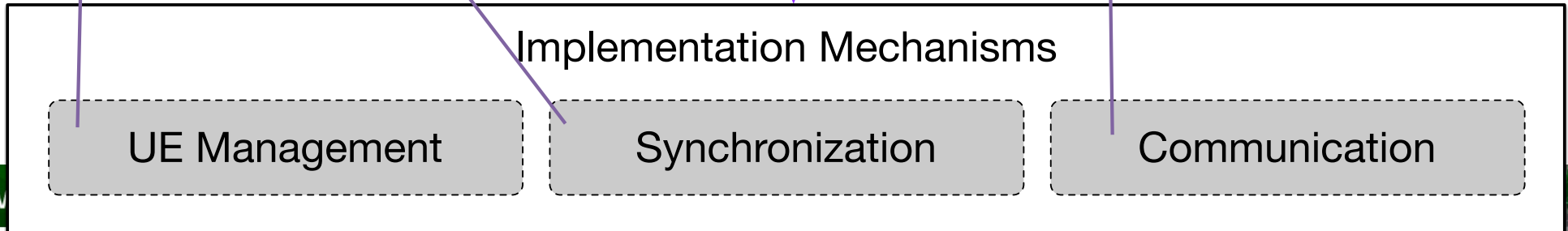
Algorithm
Structure

Supporting
Structures

Manage task
lifetimes

Enforce ordering
constraints

Get data where it
needs to be when UEs
don't share memory



Stay Tuned

- C++ threads
- C++ `async()`
- C++ atomics

Thank you!

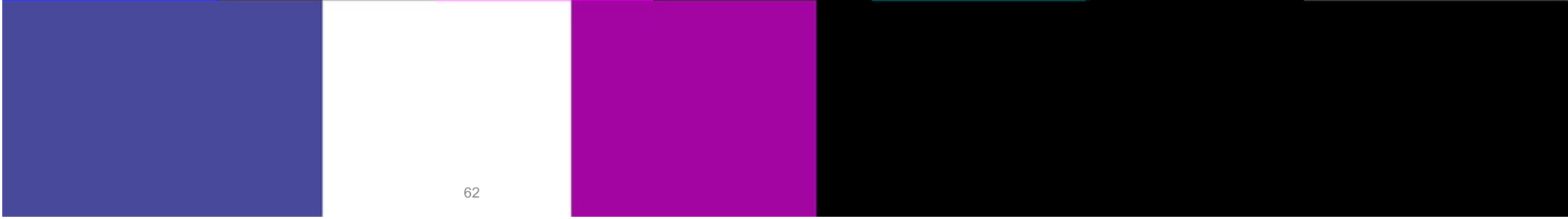
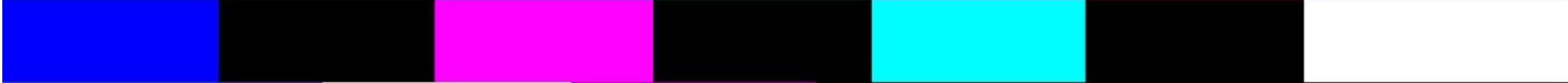
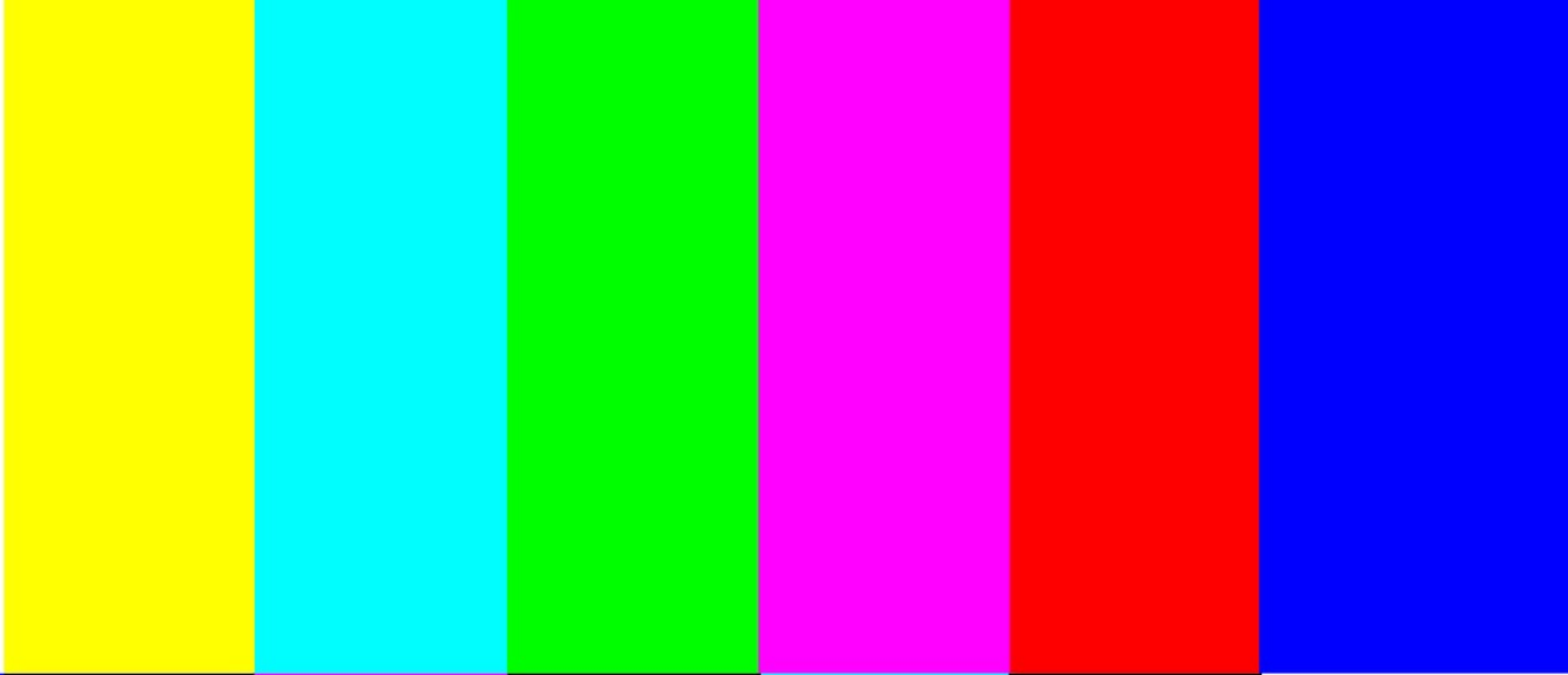
Creative Commons BY-NC-SA 4.0 License



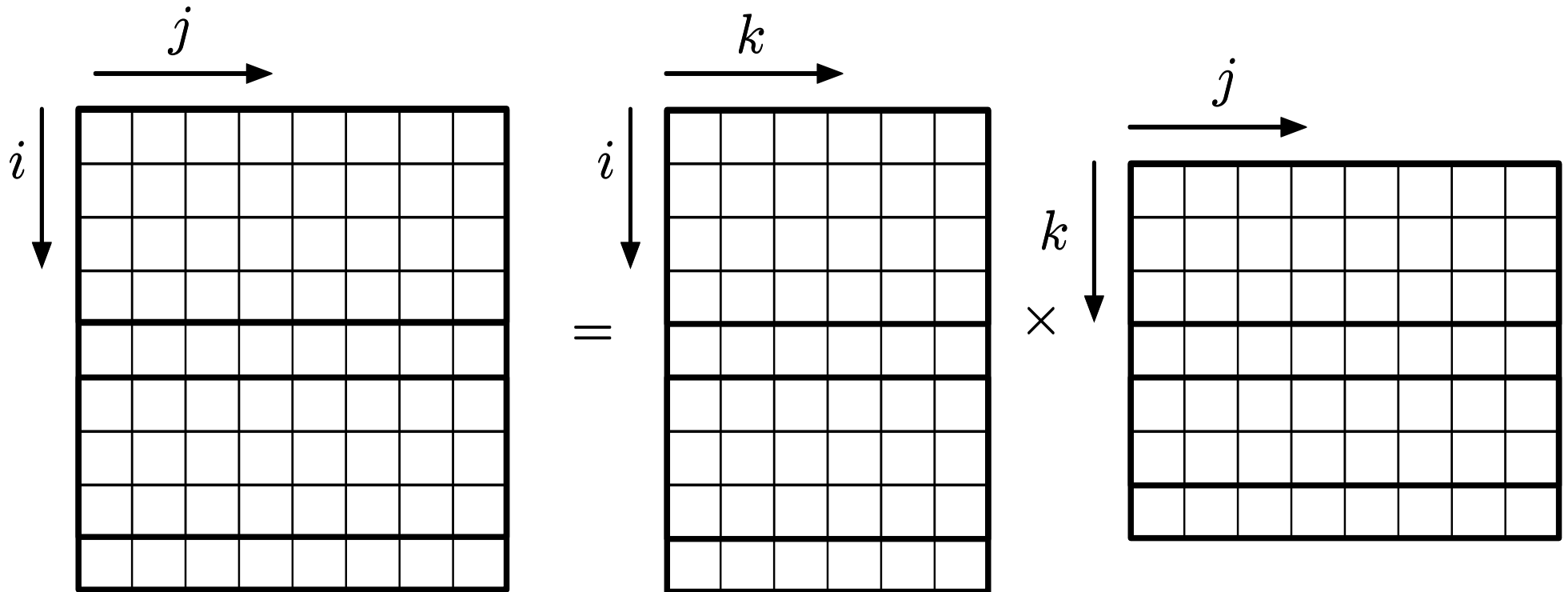
© Andrew Lumsdaine, 2017-2019

Except where otherwise noted, this work is licensed under

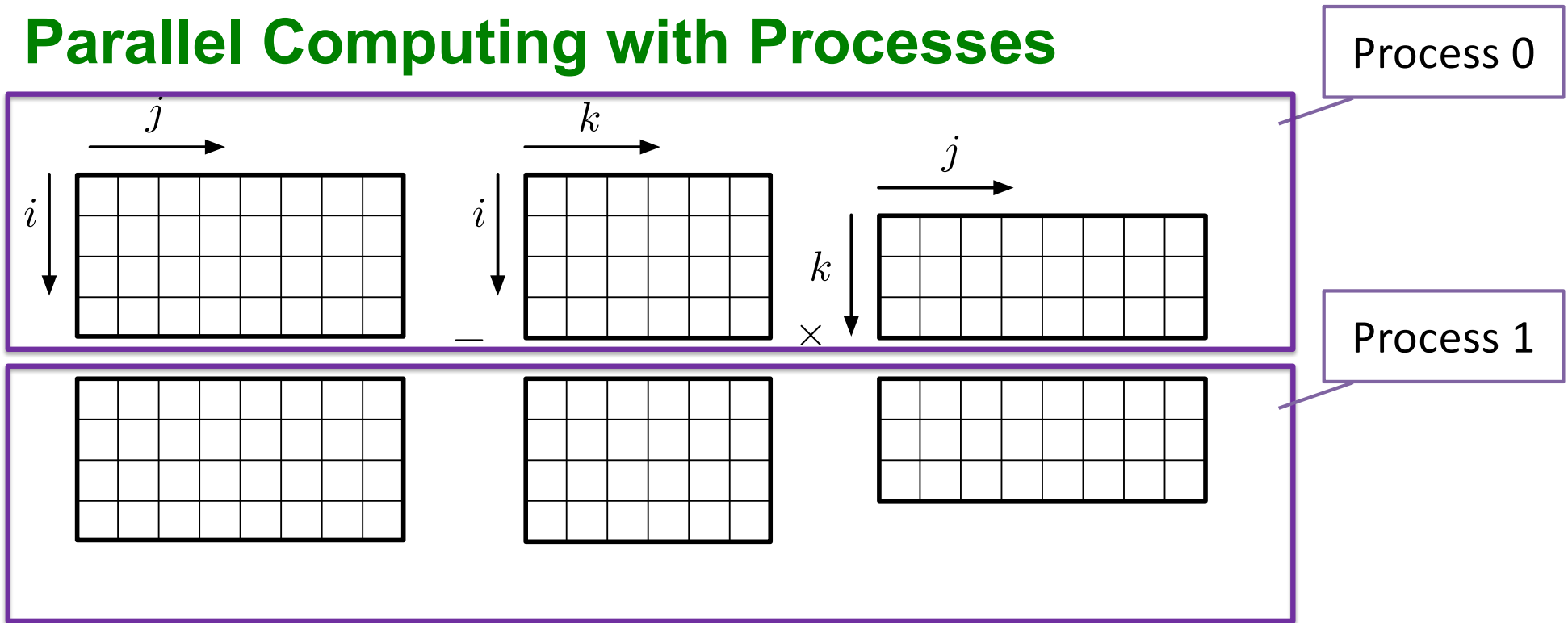
<https://creativecommons.org/licenses/by-nc-sa/4.0/>



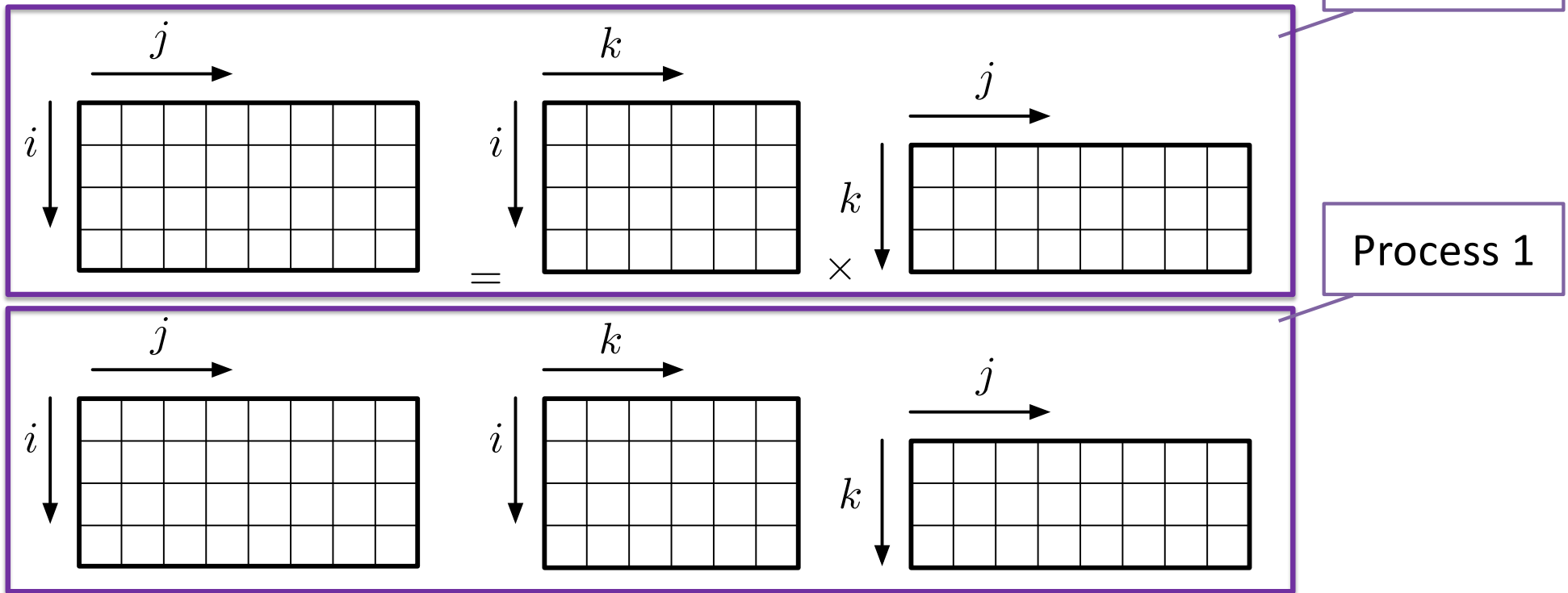
Parallel Computing with Processes



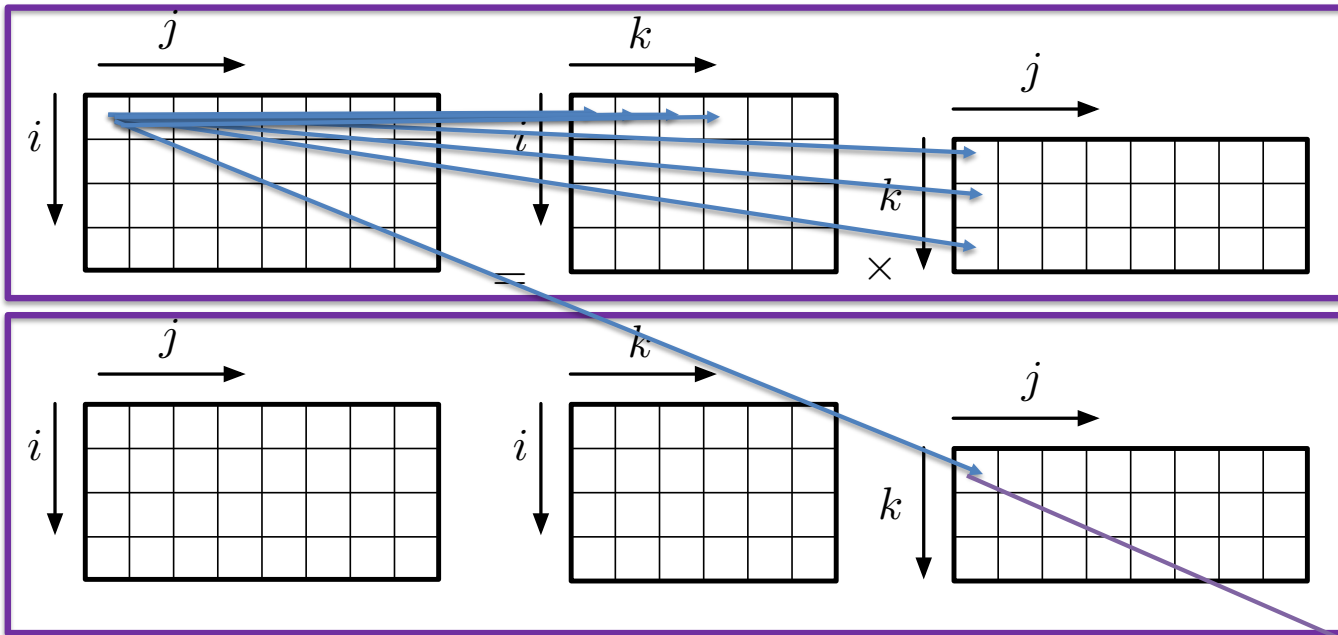
Parallel Computing with Processes



Parallel Computing with Processes



Parallel Computing with Processes



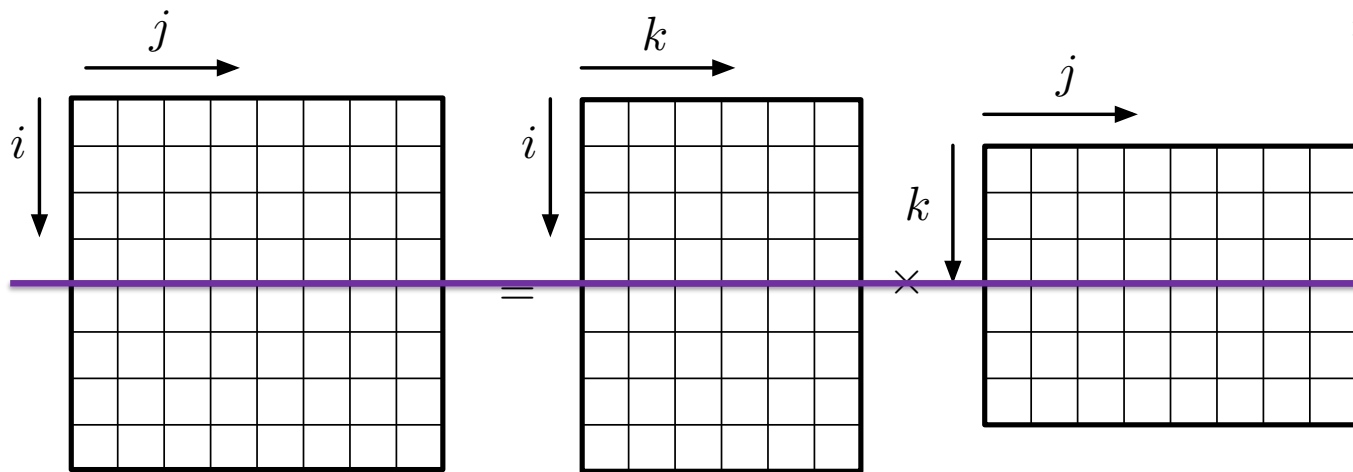
```
for (int i = 0; i < A.numRows(); ++i)
  for (int j = 0; j < B.numCols(); ++j)
    for (int k = 0; k < A.numCols(); ++k)
      C(i,j) += A(i,k) * B(k,j);
}
```

```
for (int i = 0; i < A.numRows(); ++i)
  for (int j = 0; j < B.numCols(); ++j)
    for (int k = 0; k < A.numCols(); ++k)
      C(i,j) += A(i,k) * B(k,j);
}
```

Can't index from different process b/c different address space

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

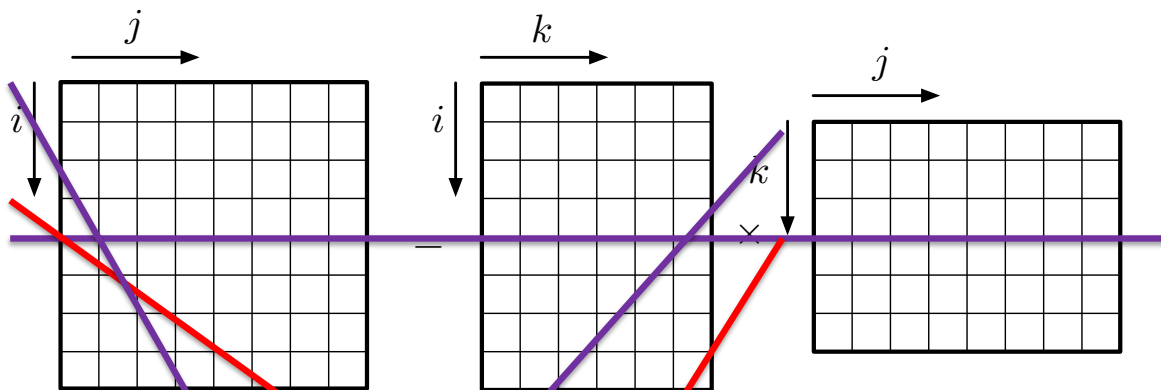
Parallel Computing with One Process



```
for (int i = 0; i < A.numRows(); ++i)
  for (int j = 0; j < B.numCols(); ++j)
    for (int k = 0; k < A.numCols(); ++k)
      C(i,j) += A(i,k) * B(k,j);
}
```

```
for (int i = 0; i < A.numRows(); ++i)
  for (int j = 0; j < B.numCols(); ++j)
    for (int k = 0; k < A.numCols(); ++k)
      C(i,j) += A(i,k) * B(k,j);
}
```

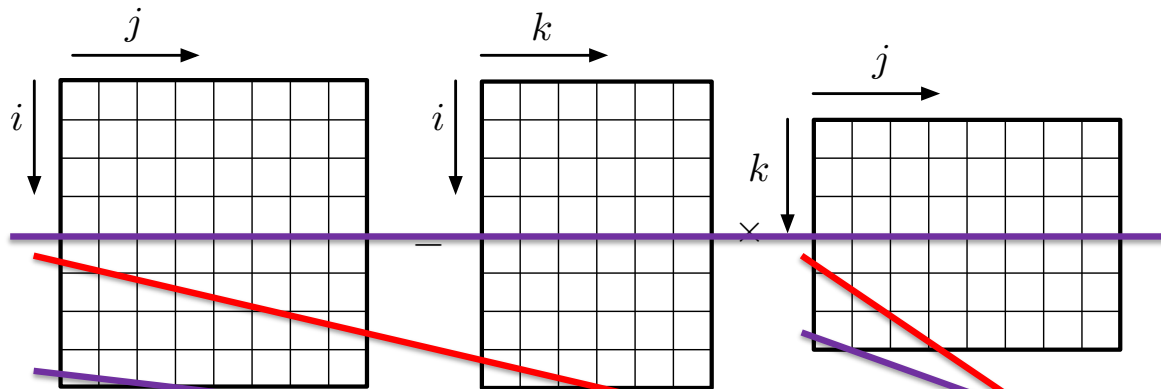
Parallel Computing with One Process



```
for (int i = 0; i < A.numRows(); ++i) {  
    for (int j = 0; j < B.numCols(); ++j) {  
        for (int k = 0; k < A.numCols(); ++k) {  
            C(i,j) += A(i,k) * B(k,j);  
        }  
    }  
}
```

```
for (int i = 0; i < A.numRows(); ++i) {  
    for (int j = 0; j < B.numCols(); ++j) {  
        for (int k = 0; k < A.numCols(); ++k) {  
            C(i,j) += A(i,k) * B(k,j);  
        }  
    }  
}
```

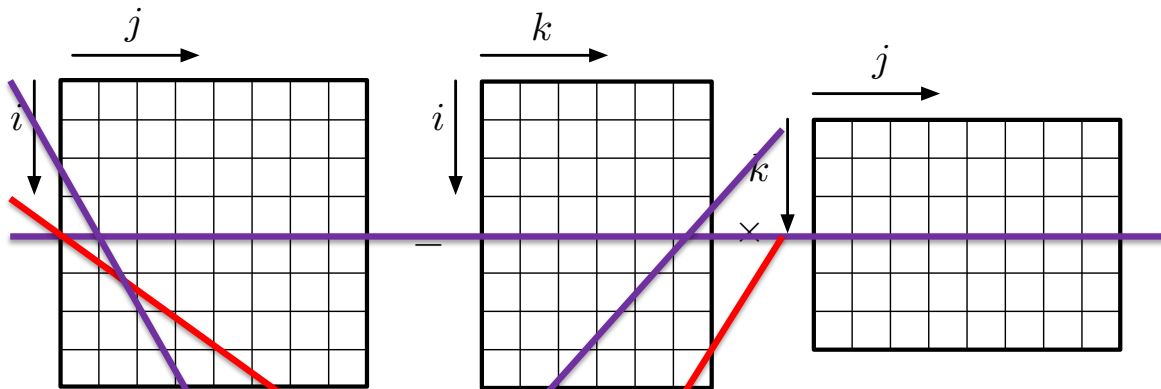
Parallel Computing with One Process



```
for (int i = 0; i < A.numRows(); ++i) {  
  for (int j = 0; j < B.numCols(); ++j) {  
    for (int k = 0; k < A.numCols(); ++k) {  
      C(i,j) += A(i,k) * B(k,j);  
    }  
  }  
}
```

```
for (int i = 0; i < A.numRows(); ++i) {  
  for (int j = 0; j < B.numCols(); ++j) {  
    for (int k = 0; k < A.numCols(); ++k) {  
      C(i,j) += A(i,k) * B(k,j);  
    }  
  }  
}
```

Parallel Computing with One Process



```

for (int i = 0; i < A.numRows(); ++i) {
  for (int j = 0; j < B.numCols(); ++j) {
    for (int k = 0; k < A.numCols(); ++k) {
      C(i,j) += A(i,k) * B(k,j);
    }
  }
}

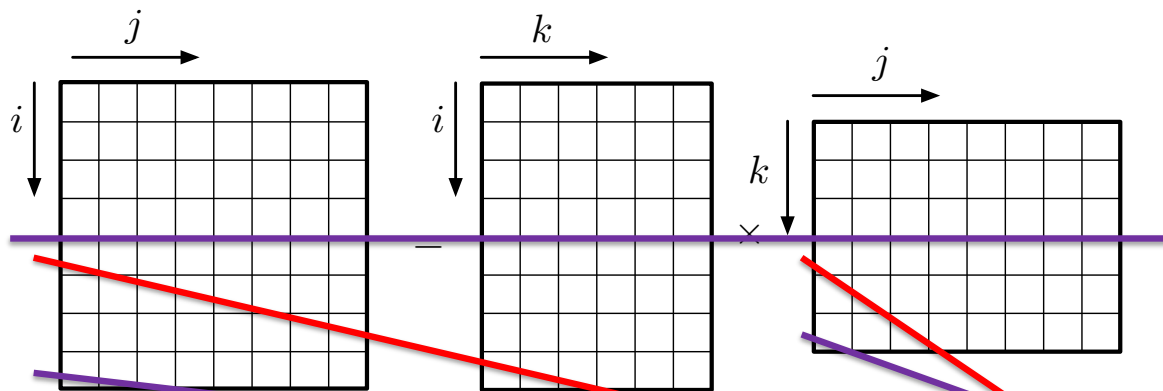
```

```

for (int i = 0; i < A.numRows(); ++i) {
  for (int j = 0; j < B.numCols(); ++j) {
    for (int k = 0; k < A.numCols(); ++k) {
      C(i,j) += A(i,k) * B(k,j);
    }
  }
}

```

Parallel Computing with One Process



```

for (int i = 0; i < A.numRows()/2; ++i) {
  for (int j = 0; j < B.numCols(); ++j) {
    for (int k = 0; k < A.numCols()/2; ++k) {
      C(i,j) += A(i,k) * B(k,j);
    }
  }
}

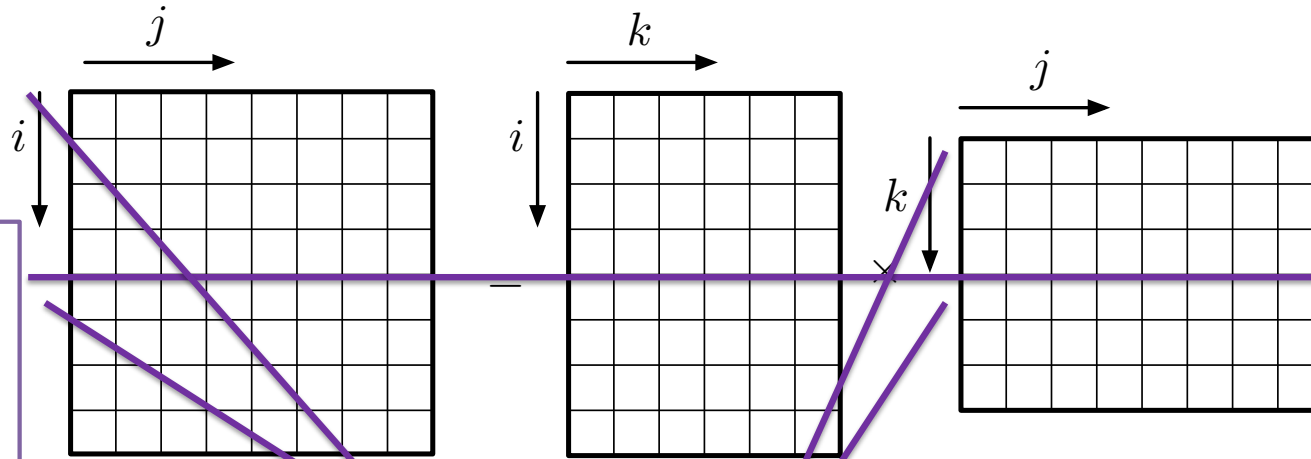
```

```

for (int i = 0; i < A.numRows(); ++i) {
  for (int j = 0; j < B.numCols(); ++j) {
    for (int k = 0; k < A.numCols(); ++k) {
      C(i,j) += A(i,k) * B(k,j);
    }
  }
}

```

Use Same Function in Both Cases

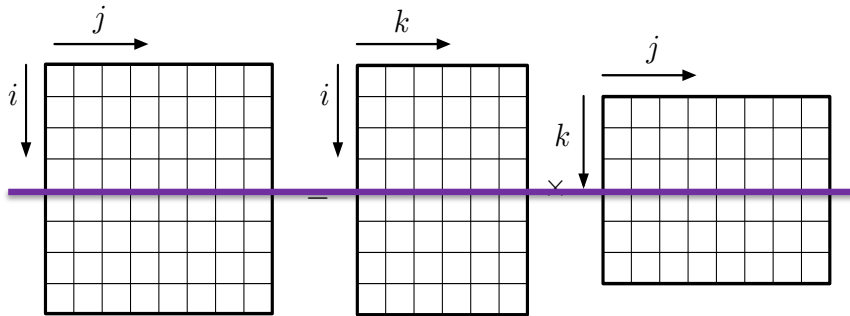


Still need to run two separate instances

Need to run them in **parallel** to get improved performance

```
for (int i = iStart; i < iStart + A.numRows()/2; ++i) {
  for (int j = 0; j < B.numCols(); ++j) {
    for (int k = kStart; k < kStart + A.numCols()/; ++k) {
      C(i,j) += A(i,k) * B(k,j);
    }
  }
}
```


Use Same Function in Both Cases



Run this

```
for (int i = iStart; i < iStart + A.numRows()/2; ++i) {  
  for (int j = 0; j < B.numCols(); ++j) {  
    for (int k = kStart; k < kStart + A.numCols()/2; ++k) {  
      C(i,j) += A(i,k) * B(k,j);  
    }  
  }  
}
```

`int iStart = 0;`
`int kStart = 0;`

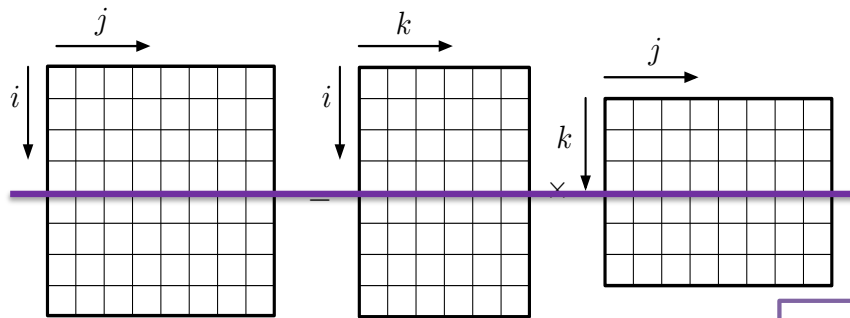
Improved
performance?

```
for (int i = iStart; i < iStart + A.numRows()/2; ++i) {  
  for (int j = 0; j < B.numCols(); ++j) {  
    for (int k = kStart; k < kStart + A.numCols()/2; ++k) {  
      C(i,j) += A(i,k) * B(k,j);  
    }  
  }  
}
```

`int iStart = A.numRows()/2;`
`int kStart = A.numCols()/2;`

Then this

Use Same Function in Both Cases



At the same time

2X faster (?)

Run this

And this

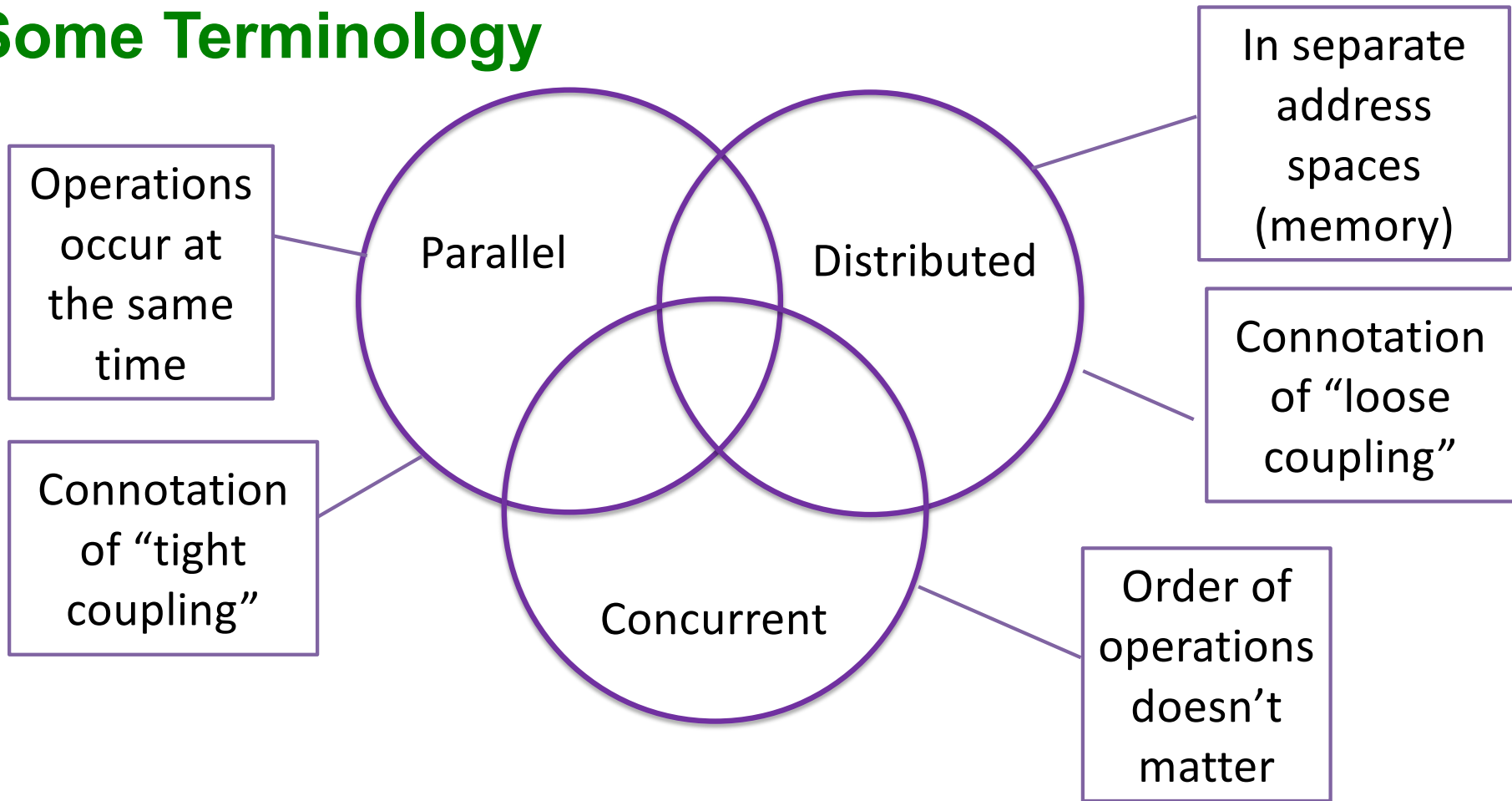
```

for (int i = iStart; i < iStart + A.numRows()/2; ++i) {
    for (int j = 0; j < B.numCols(); ++j) {
        for (int k = kStart; k < kStart + A.numCols()/2; ++k) {
            C(i,j) += A(i,k) * B(k,j);
        }
    }
}
int iStart = 0;
int kStart = 0;
    
```

```

for (int i = iStart; i < iStart + A.numRows()/2; ++i) {
    for (int j = 0; j < B.numCols(); ++j) {
        for (int k = kStart; k < kStart + A.numCols()/2; ++k) {
            C(i,j) += A(i,k) * B(k,j);
        }
    }
}
int iStart = A.numRows()/2;
int kStart = A.numCols()/2;
    
```

Some Terminology



Running Things “At the Same Time”

- Historically, threads evolved as a concurrency mechanism, not parallelism
- Enabled OS and processes to do multiple things “at the same time”
- Can be used for performance if threads are executed in parallel

Running Things “At the Same Time” in C++

```
#include <iostream>
#include <thread>
using namespace std;
```

Pull in thread library

```
void sayHello() {
    cout << "Hello World" << endl;
}
```

Simple function

Create a thread

Join back to main thread

```
int main() {
    thread helloThread(sayHello);
    helloThread.join();

    return 0;
}
```

That runs this function

Multithreading

```
void sayHello(int tnum) {
    cout << "Hello World.  I am thread " << tnum << endl;
}

int main() {
    std::thread tid[16];

    for (int i = 0; i < 16; ++i)
        tid[i] = thread (sayHello, i);

    for (int i = 0; i < 16; ++i)
        tid[i].join();

    return 0;
}
```

Multithreading

```
void sayHello(int tnum) {  
    cout << "Hello World. I am thread " << tnum << endl;  
}  
  
int main() {  
    std::thread tid[16];  
  
    for (int i = 0; i < 16; ++i)  
        tid[i] = thread (sayHello, i);  
  
    for (int i = 0; i < 16; ++i)  
        tid[i].join();  
  
    return 0;  
}
```

Program
output

\$./a.out

Hello World. I am thread Hello World. I am thread Hello
World. I am thread Hello World. I am thread Hello World. I
am thread Hello World. I am thread Hello World. I am
thread Hello World. I am thread Hello World. I am thread
02Hello World. I am thread Hello World. I am thread 13Hello
World. I am thread 5Hello World. I am thread Hello World. I
am thread 6Hello World. I am thread 47Hello World. I am
thread 8

Concurrency?

910

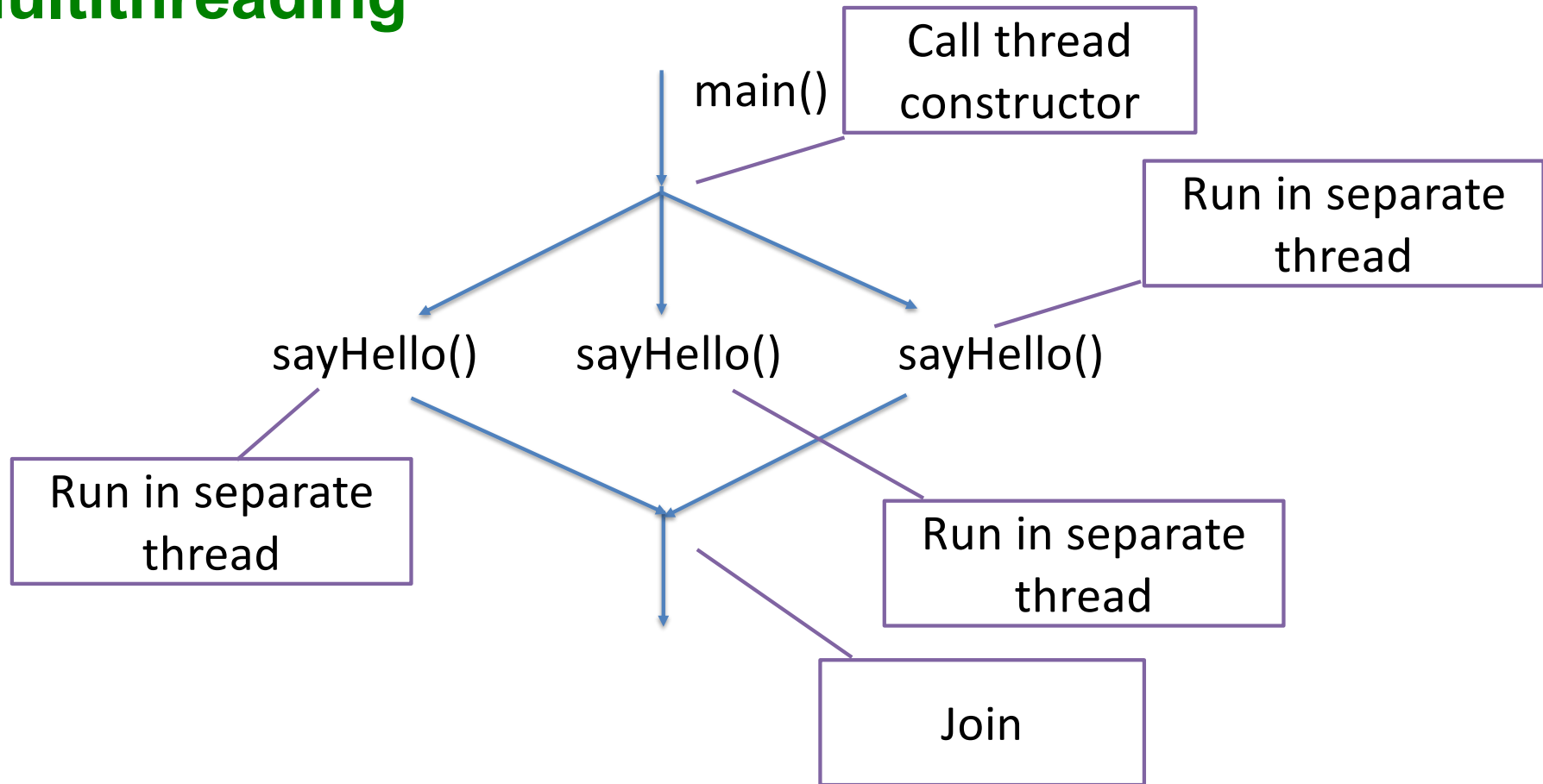
111213

14

Parallelism?

15

Multithreading



Why the Jumbled Output

```
void sayHello(int tnum) {  
    cout << "Hello World. I am thread " << tnum << endl;  
}  
  
int main() {  
    std::thread tid[16];  
  
    for (int i = 0; i < 16; ++i)  
        tid[i] = thread (sayHello, i);  
  
    for (int i = 0; i < 16; ++i)  
        tid[i].join();  
  
    return 0;  
}
```

Concurrency!

\$./a.out

Hello World. I am thread Hello World. I am thread Hello
World. I am thread Hello World. I am thread Hello World. I
am thread Hello World. I am thread Hello World. I am
thread Hello World. I am thread Hello World. I am thread
02Hello World. I am thread Hello World. I am thread 13Hello
World. I am thread 5Hello World. I am thread Hello World. I
am thread 6Hello World. I am thread 47Hello World. I am
thread 8

910

111213

14

15

Another Example

```
int value = 0;

int main() {
    std::thread tid[16];

    for (int i = 0; i < 16; ++i)
        tid[i] = thread (sayHello, i);

    for (int i = 0; i < 16; ++i)
        tid[i].join();

    cout << "Final value is " << value << endl;

    return 0;
}
```

Example

./a.outHello World. I am thread Hello World. I am thread Hello World. I
am thread Hello World. I am thread Hello World. I am thread Hello
World. I am thread Hello World. I am thread Hello World. I am thread
Hello World. I am thread Hello World. I am thread 5302Hello World. I
am thread Hello World. I am thread 64Hello World. I am thread Hello
World. I am thread 1Hello World. I am thread 789Value is Value is Value
is Hello World. I am thread Value is 1011Value is Value is 1213Value is
14Value is Value is Value is 000150Value is Value is 00Value is Value is
0Value is 000Value is 000000

Final value is 1

Not Good!

Race condition

Yet Another Example (Sequential, Synchronous)

```
int bank_balance = 300;

void withdraw(const string& msg, int amount) {
    int bal = bank_balance;
    string out_string = msg + " withdraws " + to_string(amount) + "\n";
    cout << out_string;
    bank_balance = bal - amount;
}

int main() {
    cout << "Starting balance is " << bank_balance << endl;

    withdraw("Bonnie", 100);
    withdraw("Clyde", 100);

    cout << "Final bank balance is " << bank_balance << endl;

    return 0;
}
```

Yet Another Example (Concurrent)

```
int bank_balance = 300;

void withdraw(const string& msg, int amount) {
    int bal = bank_balance;
    string out_string = msg + " withdraws " + to_string(amount) + "\n";
    cout << out_string;
    bank_balance = bal - amount;
}

int main() {
    cout << "Starting balance is " << bank_balance << endl;

    thread bonnie(withdraw, "Bonnie", 100);
    thread clyde(withdraw, "Clyde", 100);

    bonnie.join();
    clyde.join();

    cout << "Final bank balance is " << bank_balance << endl;

    return 0;
}
```

Review

- Process is an abstraction for resource allocation
- Thread is an abstraction for execution
- Concurrency vs Parallelism vs Distributed
- C++ threading library

Thank You!

NORTHWEST INSTITUTE for ADVANCED COMPUTING

90

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

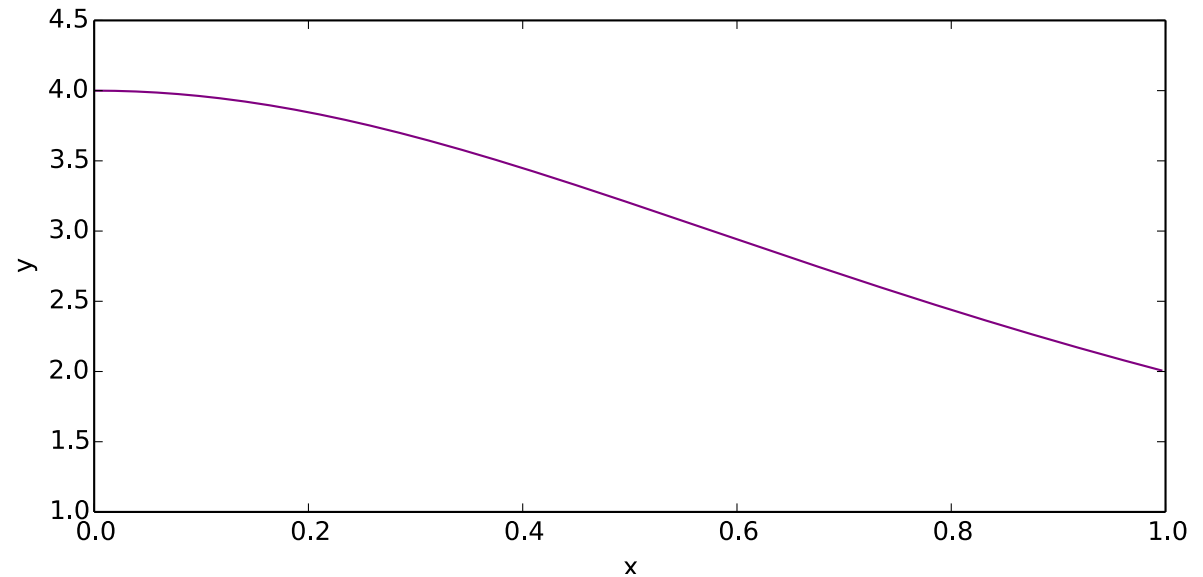

Pacific Northwest
NATIONAL LABORATORY
Proudly Operated by  Battelle
for the U.S. Department of Energy


UNIVERSITY of
WASHINGTON

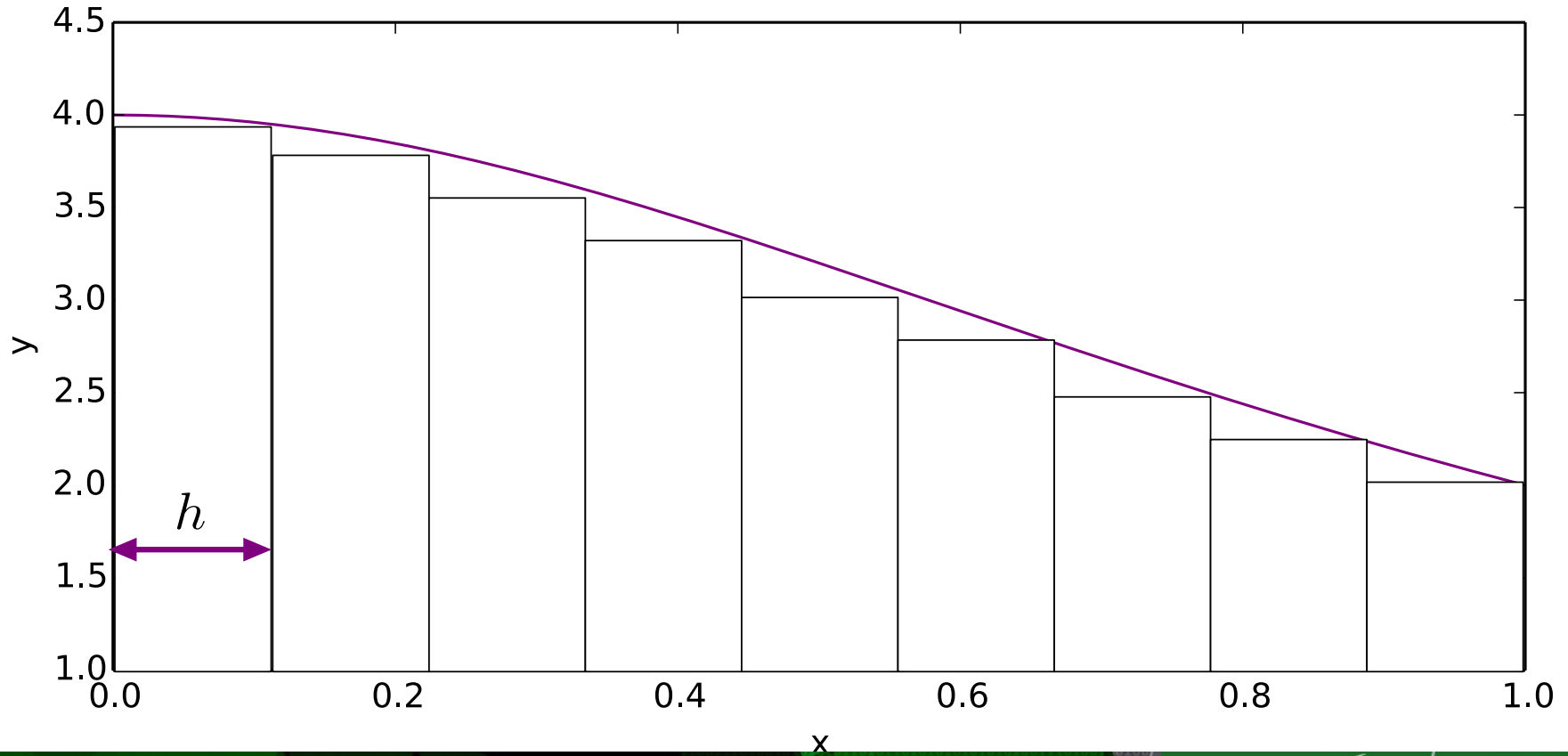
Example

- Find the value of π
- Using formula

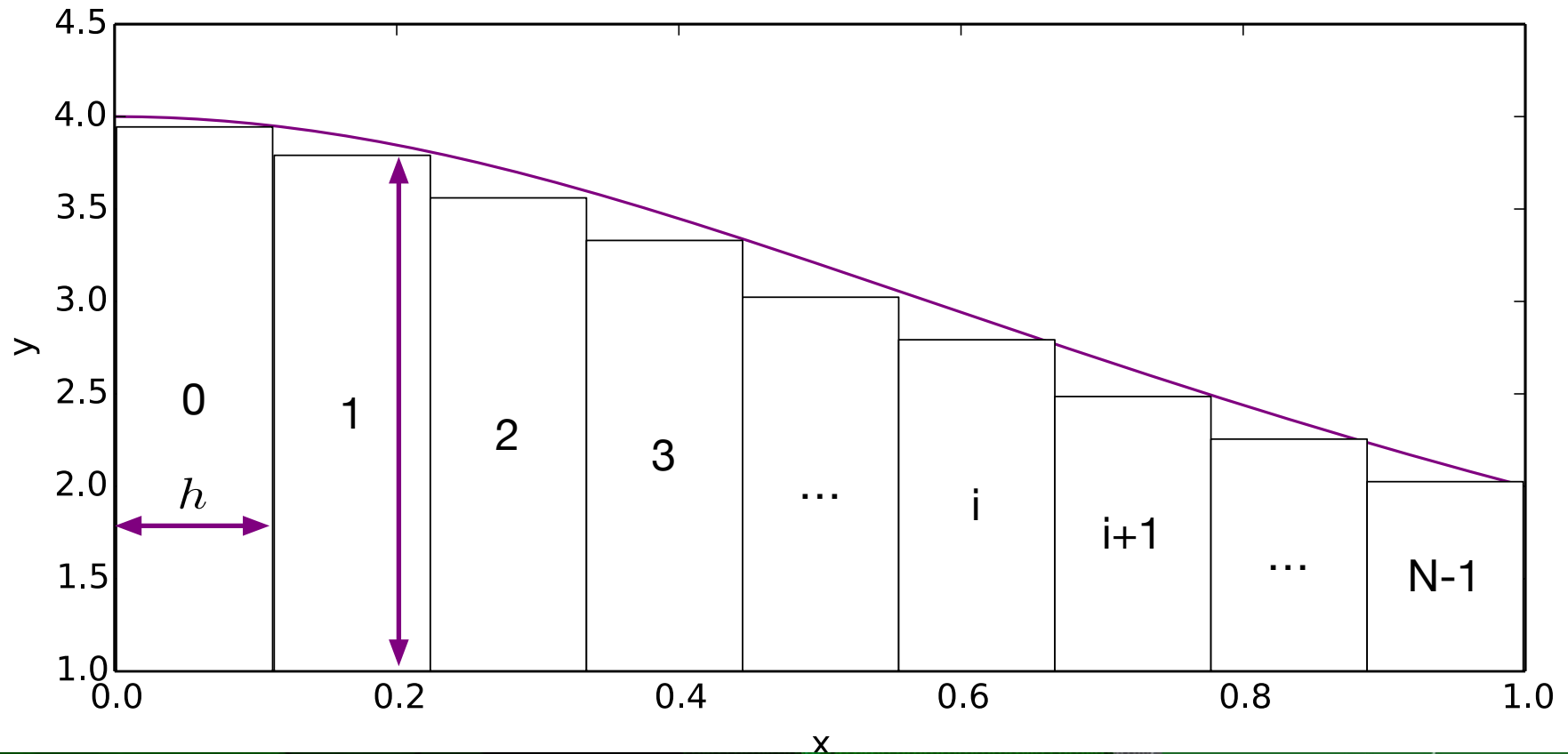
$$\pi = \int_0^1 \frac{4}{1+x^2} dx$$



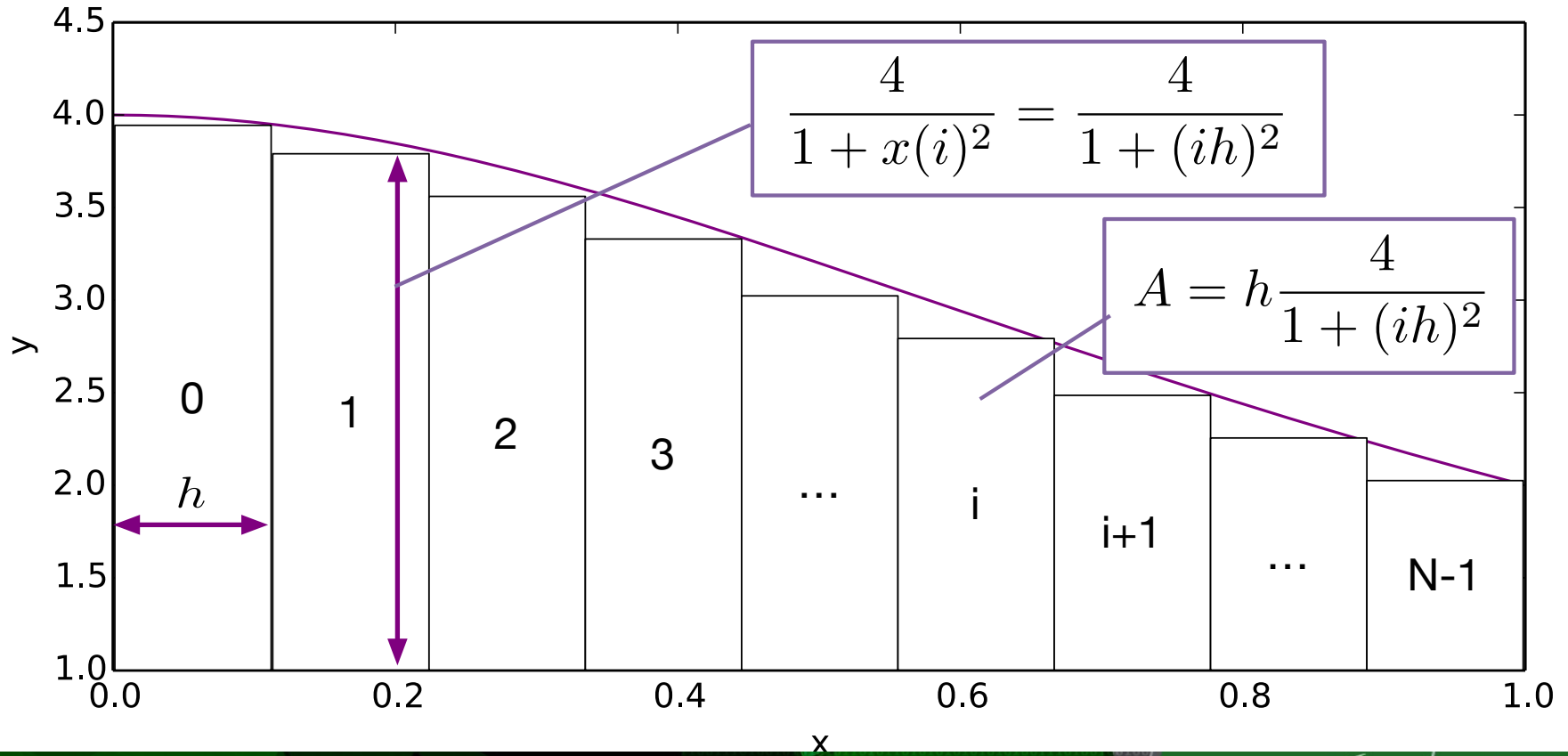
Discretization



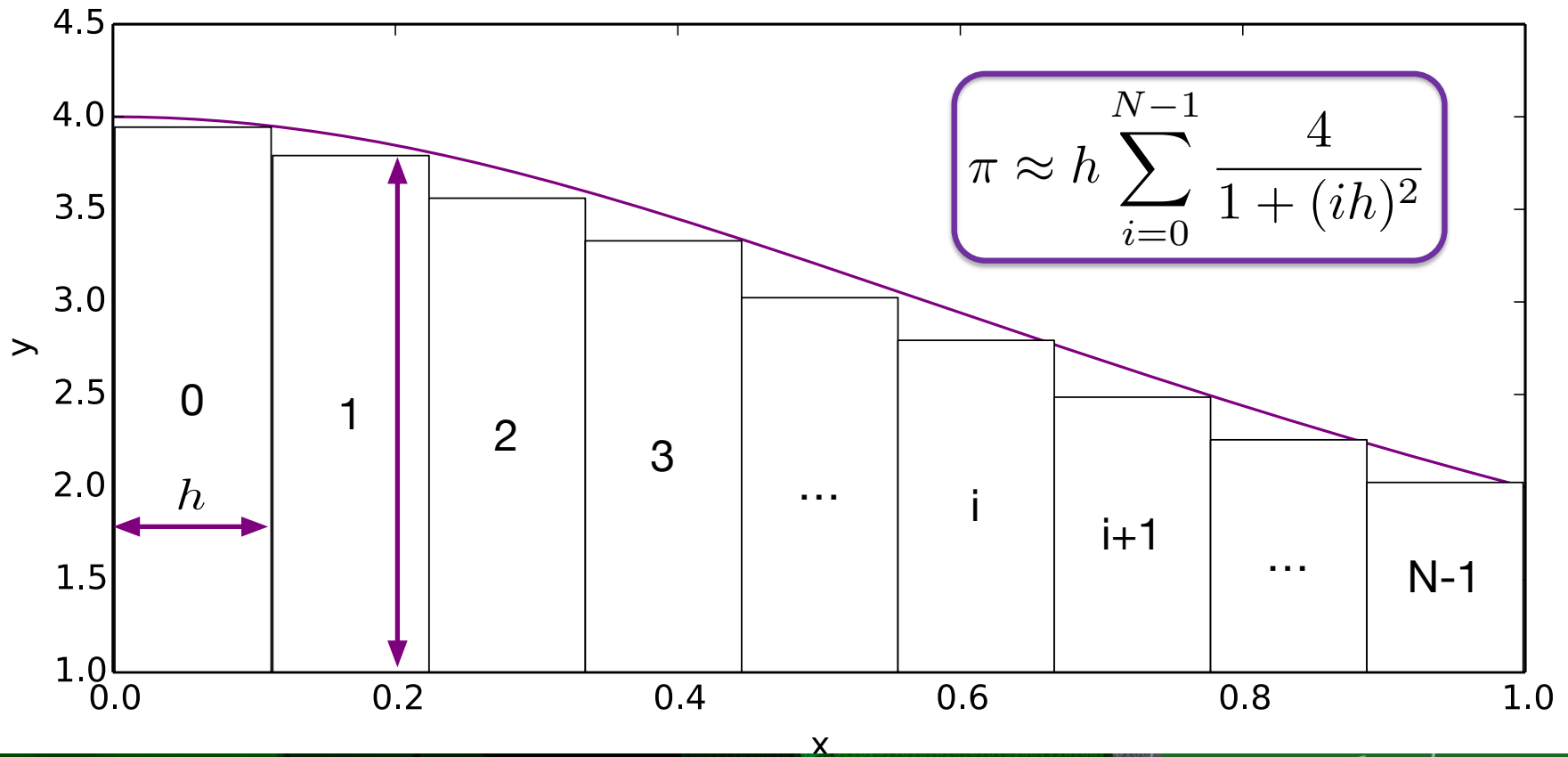
Numerical Quadrature



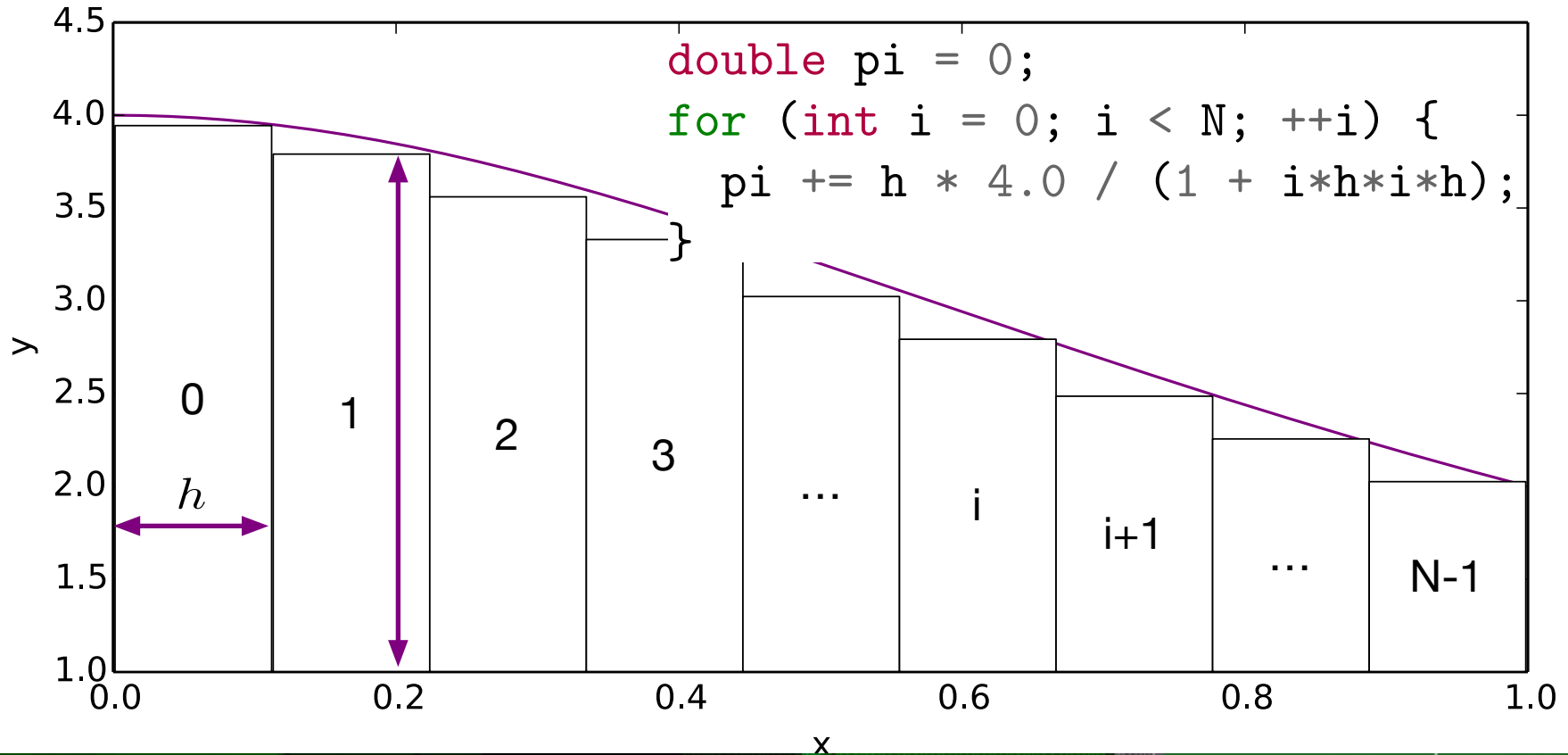
Numerical Quadrature



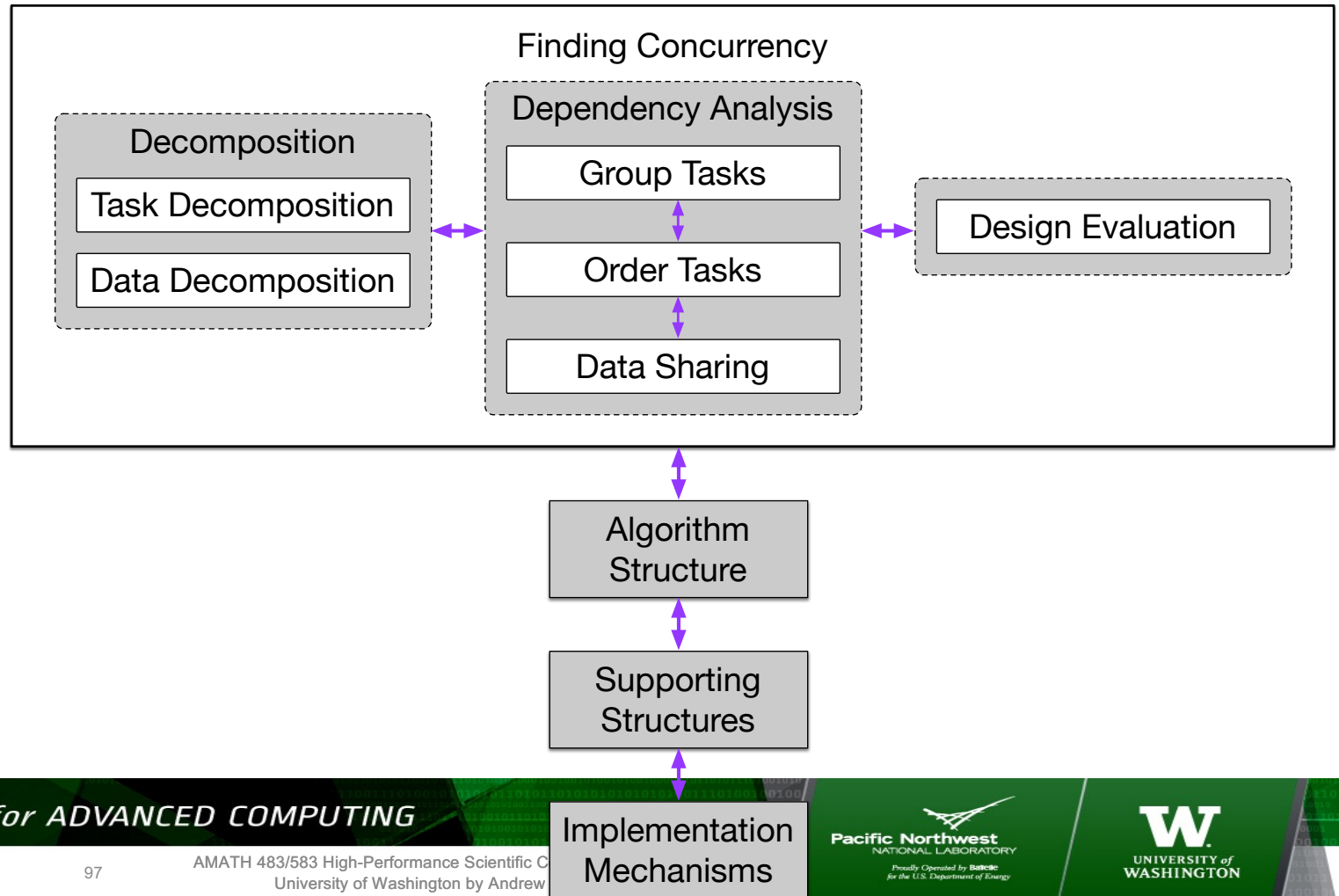
Numerical Quadrature



Numerical Quadrature



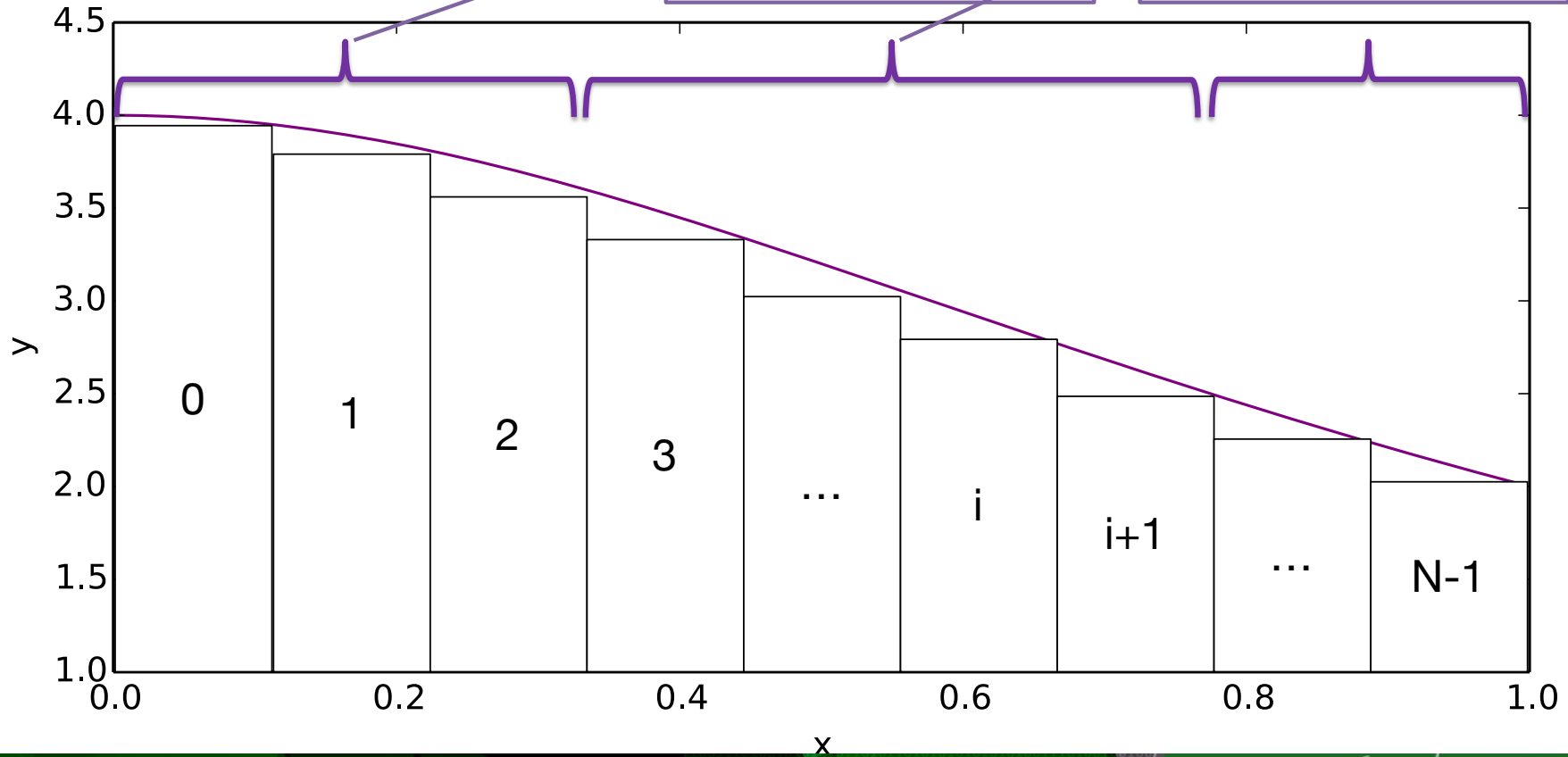
Finding Concurrency



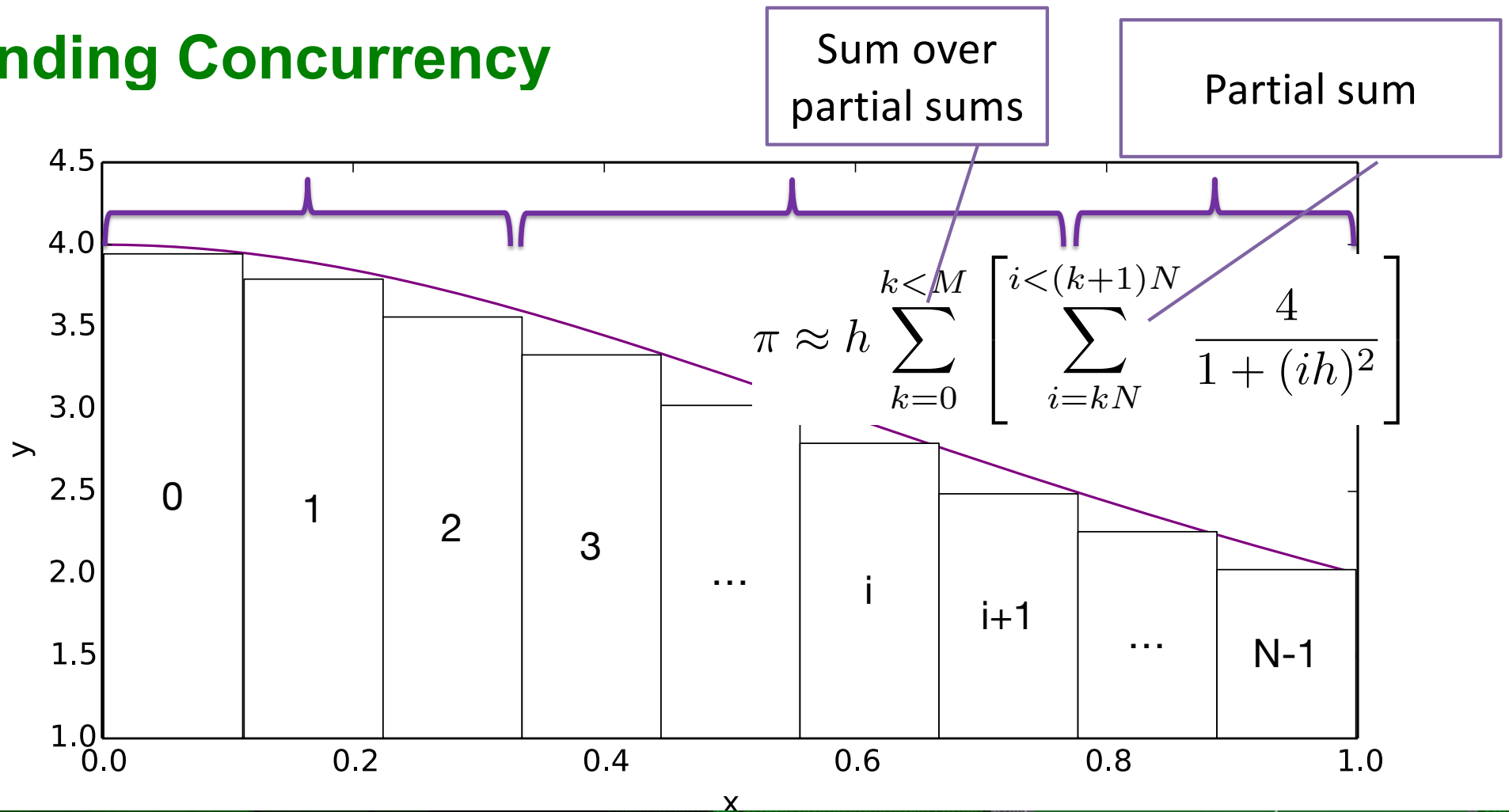
Finding Concurrency

Partial sums are all independent

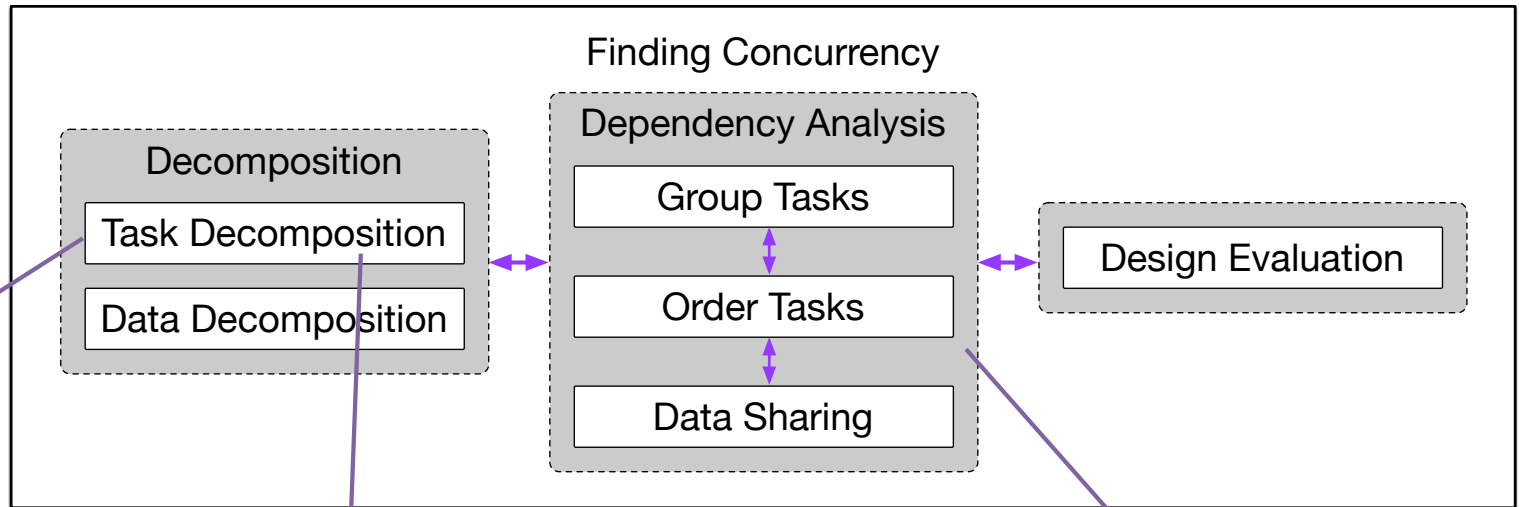
Can be computed concurrently



Finding Concurrency



Finding Concurrency



Decompose total sum into a sum of partial sums

Each task can be computed concurrently

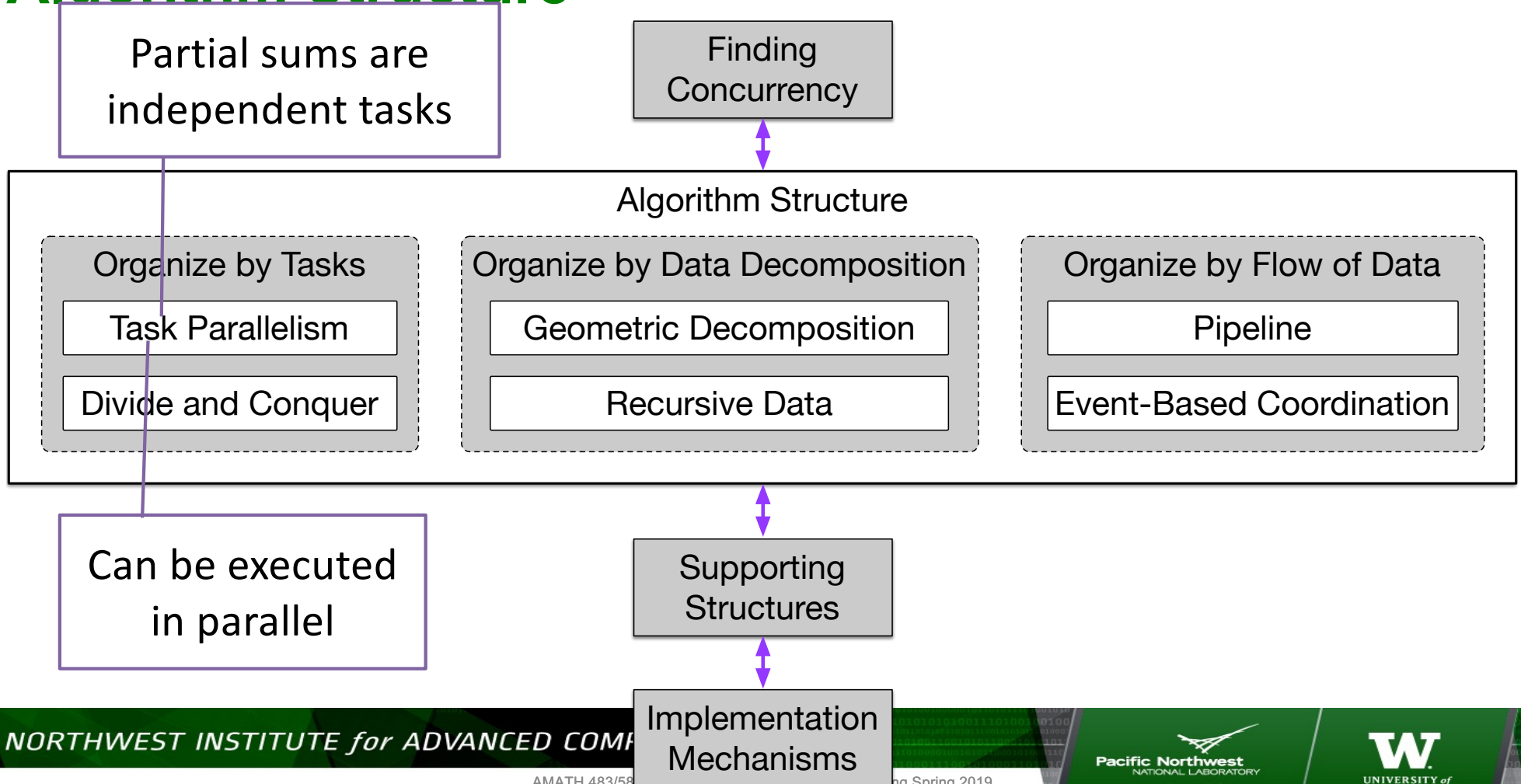
Need to sum up independent partial sums

Algorithm Structure

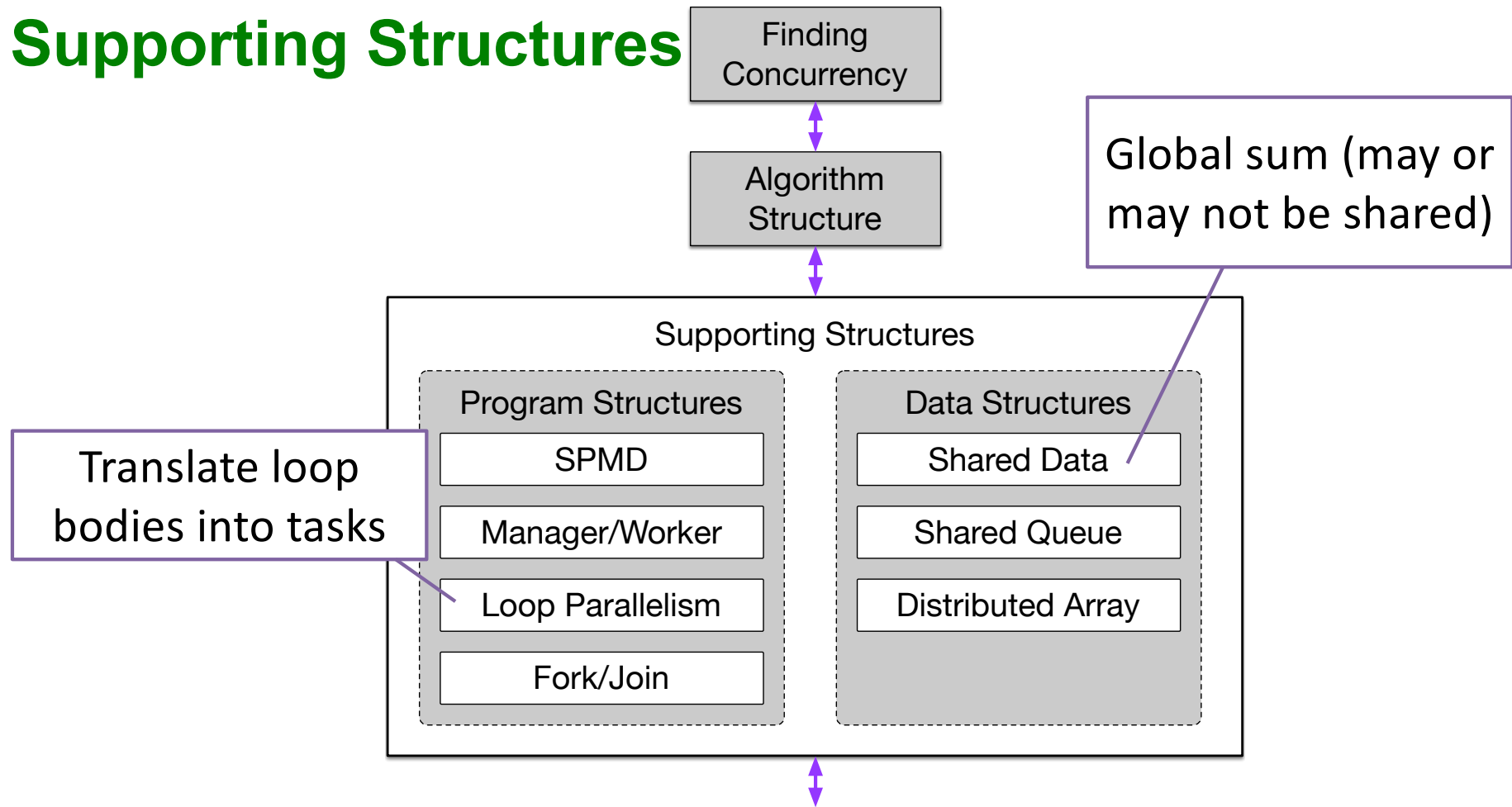
Supporting Structures

Implementation Mechanisms

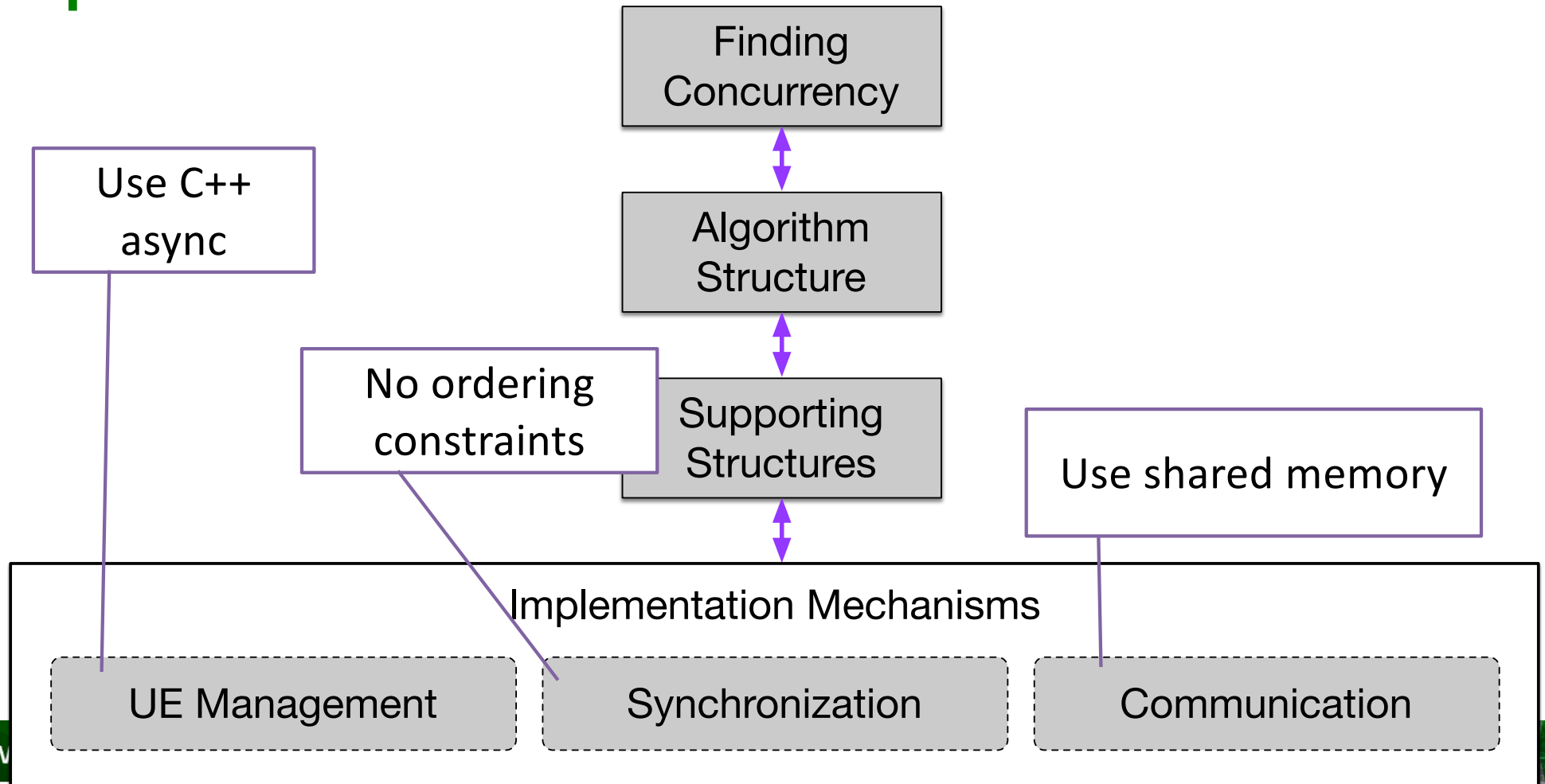
Algorithm Structure



Supporting Structures



Implementation Mechanisms



Sequential Implementation (Two Nested Loops)

```
double h = 1.0 / (double) intervals;
```

Discretization

For each set
of discretized
points

```
double pi = 0.0;
```

```
for (int k = 0; k < intervals; k += blocksize) {
```

```
    double partial_pi = 0.0;
```

```
    for (int i = k; i < (k+blocksize); ++i) {
```

```
        partial_pi += 4.0 / (1.0 + (i*h*i*h));
```

```
    }
```

```
    pi += h * partial_pi;
```

```
}
```

Compute
partial sum

Accumulate
final sum

Threads vs Tasks

```
void sayHello(int tnum) {  
    cout << "Hello World. I am thread " << tnum << endl;  
}  
  
int main() {  
    std::thread tid[16];  
  
    for (int i = 0; i < 16; ++i)  
        tid[i] = thread (sayHello, i);  
  
    for (int i = 0; i < 16; ++i)  
        tid[i].join();  
  
    return 0;  
}
```

Task

Launch
threads

“fork”

Wait for tasks
to finish

“join”

Threads

Thread
returns void

Oops

How do we get
partial sums?

How do we update
global total?

```
void partial_pi(unsigned long begin, unsigned long end) {  
    double partial_pi = 0.0;  
    for (unsigned long i = begin; i < end; ++i) {  
        partial_pi += 4.0 / (1.0 + (i*h*i*h));  
    }  
    return partial_pi;  
}  
  
int  
main(int argc, char *argv[])  
{  
    double h = 1.0 / (double) intervals;  
  
    double pi = 0.0;  
    for (int k = 0; k < intervals; k += blocksize) {  
        pi += h * partial_pi;  
    }  
    std::cout << "pi is approximately " << pi << std::endl;  
  
    return 0;  
}
```

Threads

Task

Assign task
to thread

```
void partial_pi(unsigned long begin, unsigned long end, double h, double& pi) {
    double partial_pi = 0.0;
    for (unsigned long i = begin; i < end; ++i) {
        partial_pi += 4.0 / (1.0 + (i*h*i*h));
    }
    pi += h*partial_pi;
}

int
main(int argc, char *argv[])
{
    std::vector<std::thread> threads;

    double h = 1.0 / (double) intervals;

    double pi = 0.0;
    for (unsigned long k = 0; k < num_blocks; ++k) {
        threads.push_back(std::thread(partial_pi,
            k*blocksize, (k+1)*blocksize, h, std::ref(pi)));
    }

    for (unsigned long k = 0; k < num_blocks; ++k) {
        threads[k].join();
    }
    std::cout << "pi is approximately " << pi << std::endl;

    return 0;
}
```

Threads

Local
variable

Shared
variable

```
void partial_pi(unsigned long begin, unsigned long end, double h, double& pi) {  
    double partial_pi = 0.0;  
    for (unsigned long i = begin; i < end; ++i) {  
        partial_pi += 4.0 / (1.0 + (i*h*i*h));  
    }  
    pi += h*partial_pi;  
}
```

Update shared
variable

Threads

Container for created threads

Thread constructor

Function that will be the task

Arguments to the function

```
int main(int argc, char *argv[]) {
    double h = 1.0 / (double) intervals;

    std::vector<std::thread> threads;

    double pi = 0.0;
    for (unsigned long k = 0; k < num_blocks; ++k) {
        threads.push_back(
            std::thread(
                partial_pi, k*blocksize, (k+1)*blocksize, h, std::ref(pi)));
    }

    for (unsigned long k = 0; k < num_blocks; ++k) {
        threads[k].join();
    }

    std::cout << "pi is approximately " << pi << std::endl;

    return 0;
}
```

Have to explicitly tag this as a reference

We are invoking `std::thread`, not `partial_pi`

Results

```
$ ./thrpi  
pi is approximately 3.14159
```

```
$ ./thrpi  
pi is approximately 3.14159
```

Correct

Correct

Incorrect!

Exactly same
program!

What
happened?

Name This Famous Couple

Bonnie Parker

Clyde Barrow



Bonnie and Clyde Use ATMs



```
int bank_balance = 300;

void withdraw(const string& msg, int amount) {
    int bal = bank_balance;
    string out_string = msg + " withdraws " + to_string(amount) + "\n";
    cout << out_string;
    bank_balance = bal - amount;
}

int main() {
    cout << "Starting balance is " << bank_balance << endl;

    thread bonnie(withdraw, "Bonnie", 100);
    thread clyde(withdraw, "Clyde", 100);

    bonnie.join();
    clyde.join();

    cout << "Final bank balance is " << bank_balance << endl;

    return 0;
}
```


Withdraw Function

```
int bank_balance = 300;

void withdraw(const string& msg, int amount) {
    int bal = bank_balance;
    string out_string = msg + " withdraws " + to_string(amount) + "\n";
    cout << out_string;
    bank_balance = bal - amount;
}
```

Get balance

Save new
balance

Compute the
new balance

Making Concurrent Withdrawals

```
int main() {  
    cout << "Starting balance is " << bank_balance << endl;  
    thread bonnie(withdraw, "Bonnie", 100);  
    thread clyde(withdraw, "Clyde", 100);  
  
    bonnie.join();  
    clyde.join();  
  
    cout << "Final bank balance is " << bank_balance << endl;  
  
    return 0;  
}
```

Launch
threads

Run withdraw
function

Constructor

Wait for
completion

Bonnie and Clyde Use ATMs



\$./a.out

Starting balance is 300

Bonnie withdraws 100

Clyde withdraws 100

Is this
correct?

What Happened?

Bonnie's thread,
bal = 300

Clyde's thread,
bal = 300

```
void withdraw(const string& msg, int amount) {  
    int bal = bank_balance;  
    string out_s = msg + " withdraws " + to_string(amt) + "\n";
```

```
void withdraw(const string& msg, int amount) {  
    int bal = bank_balance;  
    string out_s = msg + " withdraws " + to_string(amt) + "\n";
```

Context switch

Context switch

```
    cout << out_s;  
    bank_balance = bal - amount;  
}
```

```
    cout << out_s;  
    bank_balance = bal - amount;  
}
```

bal is still 300

bal is still 300

bank_balance
gets 200

bank_balance
gets 200

Profit!

What Happened: Race Condition

- Final answer depends on instructions from different threads are interleaved with each other
- Often occurs with shared writing of shared data
- Often due to read then update shared data
- What was true at the read is not true at the update

Critical Section Problem

```
int bank_balance = 300;

void withdraw(const string& msg, int amount) {
    int bal = bank_balance;
    string out_string = msg + " withdraws " + to_string(amount) + "\n";
    cout << out_string;
    bank_balance = bal - amount;
}
```

We want to tell
operating system not to
run anything else here

When some thread is executing
this *critical section*, no other
thread may execute it

The Critical-Section Problem

- n processes all competing to use some shared data
- Each process has a code segment, called critical section, in which the shared data is accessed.
- Problem – ensure that when one process is executing in its critical section, no other process is allowed to execute in its critical section.
- What do we mean by “execute in its critical section”?

Solution to Critical-Section Problem

- **Mutual Exclusion** - If process P_i is executing in its critical section, then no other processes can be executing in their critical sections
- **Progress** - If no process is executing in its critical section and there exist some processes that wish to enter their critical section, then the selection of the processes that will enter the critical section next cannot be postponed indefinitely
- **Bounded Waiting** - A bound must exist on the number of times that other processes are allowed to enter their critical sections after a process has made a request to enter its critical section and before that request is granted
 - Assume that each process executes at a nonzero speed
 - No assumption concerning relative speed of the N processes

Critical Section Problem

```
int bank_balance = 300;

void withdraw(const string& msg, int amount) {
    int bal = bank_balance;
    string out_string = msg + " withdraws " + to_string(amount) + "\n";
    cout << out_string;
    bank_balance = bal - amount;
}
```

This is a *critical section*

Let's just think about mutual exclusion for now

Critical Section Problem

```
bool lock = false;

int bank_balance = 300;

void withdraw(const string& msg, int amount) {
    while (lock == true)
        ;
    lock = true;

    int bal = bank_balance;
    string out_string = msg + " withdraws " + to_string(amount) + "\n";
    cout << out_string;
    bank_balance = bal - amount;

    lock = false;
}
```

Take the lock

Execute
critical
section

Release lock

Test if another
thread is holding
the lock

Spin if it is

Fall through when lock == false

Aside

```
bool lock = false;

int bank_balance = 300;

void withdraw(const string& msg, int amount) {

    string out_string = msg + " withdraws " + to_string(amount) + "\n";
    cout << out_string;
    bank_balance -= amount;

}
```

Still a race

Aside

Then write

```
bool lock = false;

int bank_balance = 300;

void withdraw(const string& msg, int amount) {

    string out_string = msg + " withdraws " + to_string(amount) + "\n";
    cout << out_string;
    bank_balance = bank_balance - amount;
}
```

Still a race

Read

Compute

Critical Section Problem

Critical
section

```
bool lock = false;

int bank_balance = 300;

void withdraw(const string& msg, int amount) {

    string out_string = msg + " withdraws " + to_string(amount) + "\n";
    cout << out_string;

    bank_balance = bank_balance - amount;

}
```

Solution (?)

```
bool lock = false;

int bank_balance = 300;

void withdraw(const string& msg, int amount) {

    string out_string = msg + " withdraws " + to_string(amount) + "\n";
    cout << out_string;

    while (lock == true)
        ;

    lock = true;

    bank_balance = bank_balance - amount;

    lock = false;
}
```

Test if another thread is holding the lock

Take the lock

Execute critical section

Release lock

Spin if it is

Fall through when lock == false

Solution (?)

Common pattern (when correct)

Take the lock

Lock might be taken between the test and the set

```
bool lock = false;

int bank_balance = 300;

void withdraw(const string& msg, int amount) {

    string out_string = msg + " withdraws " + to_string(amount) + "\n";
    cout << out_string;

    while (lock == true)
        ;

    lock = true;

    bank_balance = bank_balance - amount;

    lock = false;
}
```

Test if another thread is holding the lock

Spin if it is

Fall through when lock == false

We've traded one critical section problem for another

Synchronization Hardware

- Many systems provide hardware support for critical section code
- Uniprocessors – could disable interrupts
 - Currently running code would execute without preemption
 - Generally too inefficient on multiprocessor systems
 - Operating systems using this not broadly scalable
- Modern machines provide special **atomic** hardware instructions
 - Atomic = non-interruptable
 - Either test memory word and set value
 - Or swap contents of two memory words

Test and Set

```
bool TestAndSet (bool& target)
{
    bool rv = target;
    target = TRUE;
    return rv;
}
```

```
bool TestAndSet (bool *target)
{
    bool rv = *target;
    *target = TRUE;
    return rv;
}
```

These are the semantics, not the implementation

Implemented in hardware as an invisible instruction

Compare And Swap

```
void CompareAndSwap (bool& a, bool& b)
{
    bool temp = a;
    a = b;
    b = temp;
}
```

These are the semantics, not the implementation

```
void CompareAndSwap (bool *a, bool *b)
{
    bool temp = *a;
    *a = *b;
    *b = temp;
}
```

Implemented in hardware as an invisible instruction

Correct Withdraw

```
int bank_balance = 300;
bool lock = false;

void withdraw(const string& msg, int amount) {
    string out_s = msg + " withdraws " + to_string(amt) + "\n";
    cout << out_s;

    while (TestAndSet(lock) == true)
        ;

    bank_balance -= amount;

    lock = false;
}
```

Under what condition will we fall through?

Spin while the value is true (another thread holds the lock)

What is the state of the lock?

Release the lock

Correct Withdraw

```
int bank_balance = 300;
bool lock = false;

void withdraw(const string& msg, int amount) {
    string out_s = msg + " withdraws " + to_string(amt) + "\n";
    cout << out_s;
```

```
    while (TestAndSet(lock) == true)
```

```
    ;
```

```
    bank_balance -= amount;
```

```
    lock = false;
```

What is the CPU doing?

How is it affecting other threads?

"Spin lock"
(common pattern)

Is this a good programming abstraction?

Parallel Speedup, Parallel Efficiency

Speedup on p processing units

Time to run problem size n on one PU

Time to run problem size n on p PUs

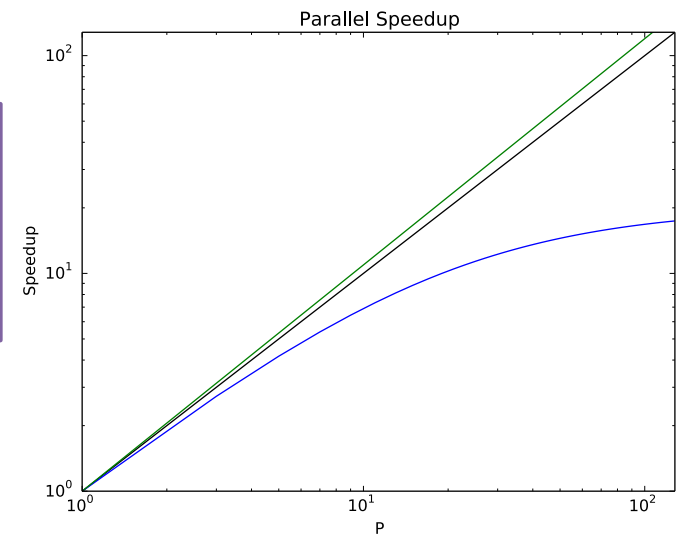
$$S(p) = \frac{T(n, 1)}{T(n, p)}$$

Efficiency on p processing units

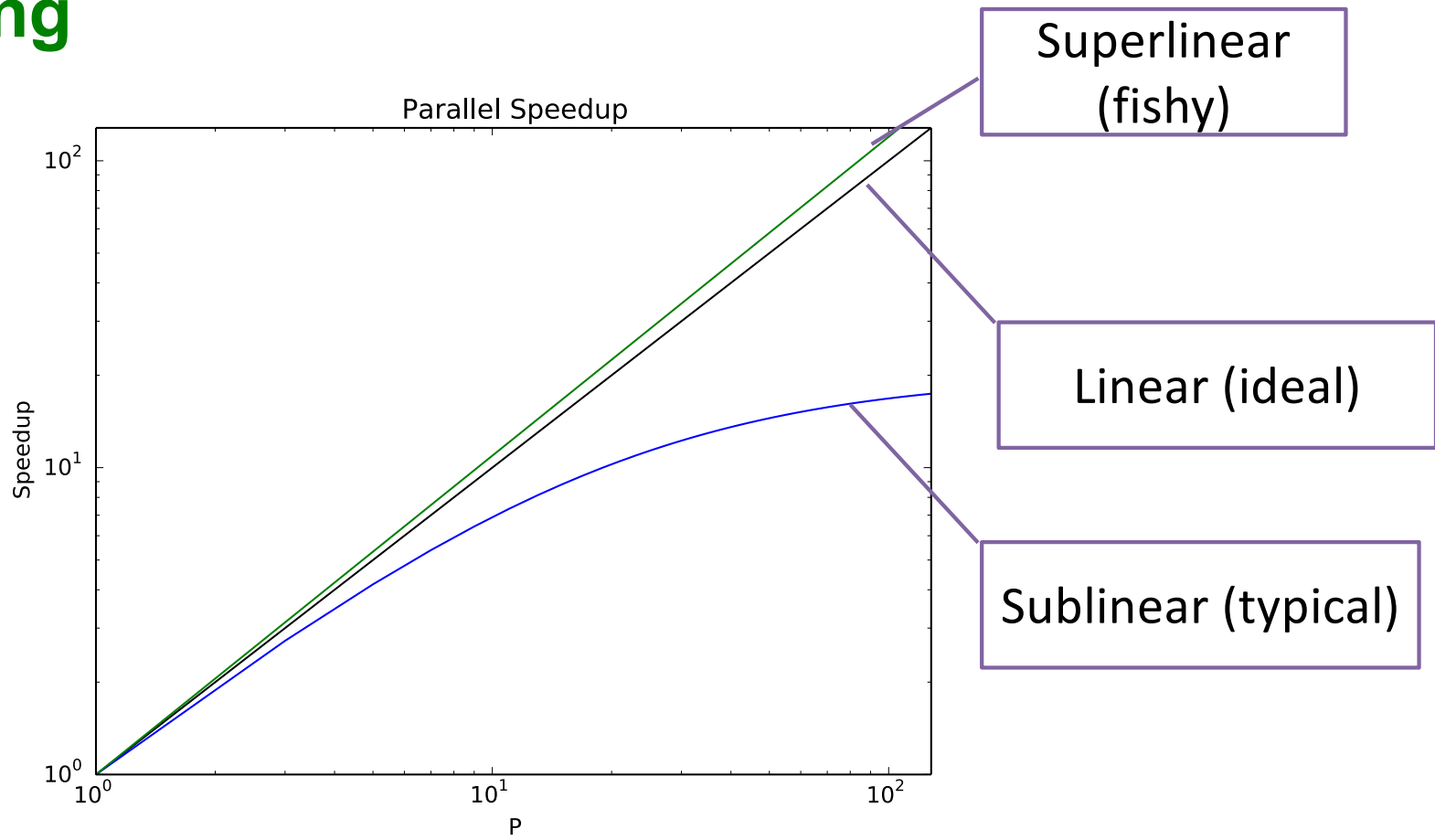
Ideal parallel execution time

Divided by actual parallel execution time

$$E(p) = \frac{T(n, 1)/p}{T(n, p)} = \frac{T(n, 1)/T(n, p)}{p} = \frac{S(p)}{p}$$



Scaling



Name This Famous Person

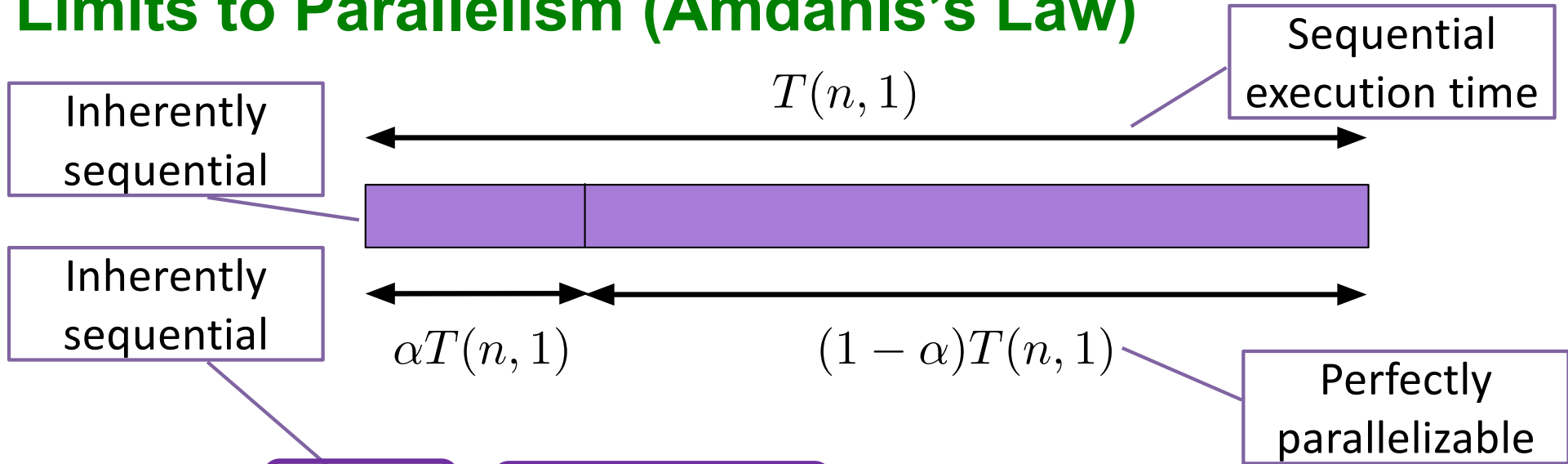


"Validity of the single processor approach to achieving large-scale computing capabilities," AFIPS Conference Proceedings (30): 483–485, 1967.

Gene Amdahl (1922-2015)

Amdahl's Law

Limits to Parallelism (Amdahl's Law)



$$T(n, 1) = \alpha T(n, 1) + (1 - \alpha)T(n, 1)$$

Perfectly parallelizable

$$T(n, p) = \alpha T(n, 1) + \frac{1}{p}(1 - \alpha)T(n, 1)$$

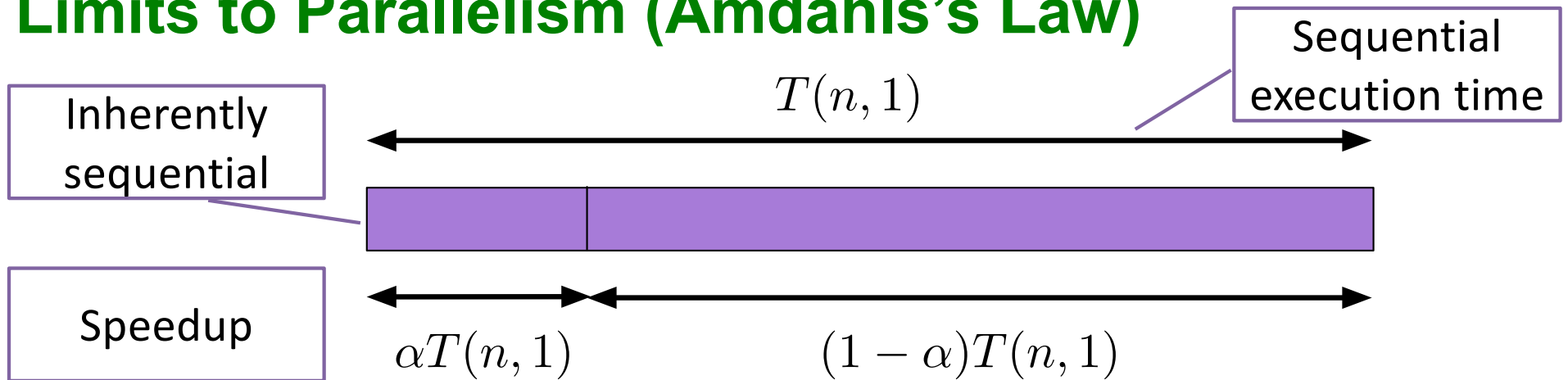
$$= T(n, 1)\left(\alpha + \frac{1}{p}(1 - \alpha)\right)$$

Ideal speedup (for parallelizable portion)

Sequential portion

INSTITUTE for ADVANCED COMPUTING

Limits to Parallelism (Amdahl's Law)



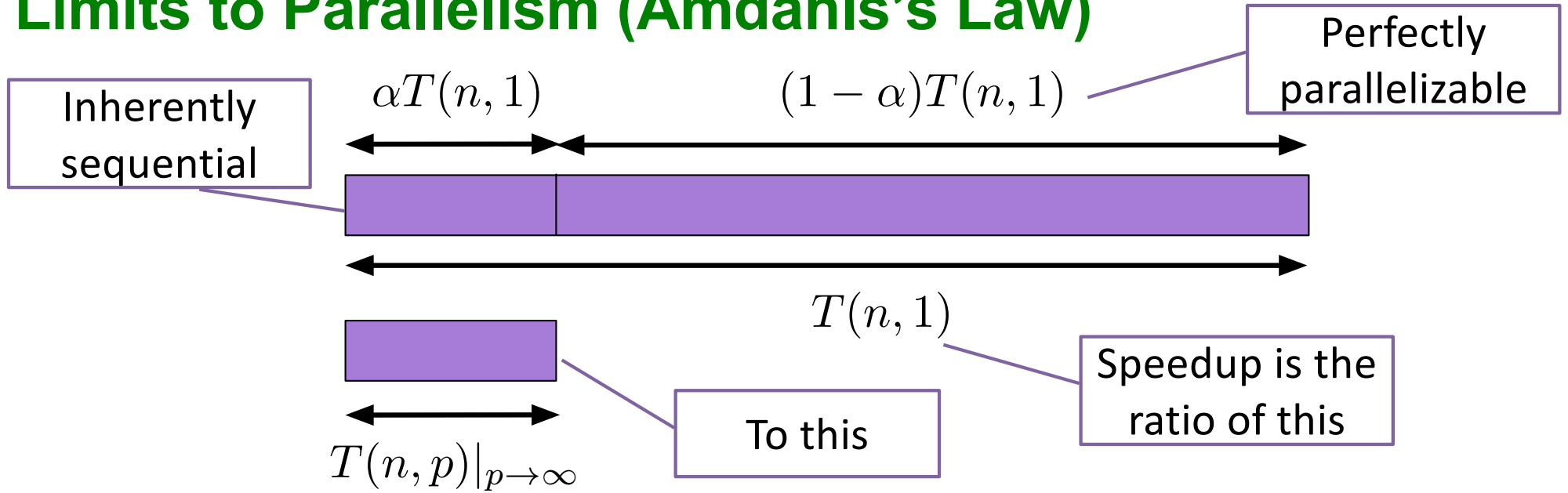
$$S(p) = \frac{T(n, 1)}{T(n, p)} = \frac{T(n, 1)}{T(n, 1) \left[\alpha + \frac{1}{p}(1 - \alpha) \right]}$$

$$= \frac{1}{\alpha + \frac{1}{p}(1 - \alpha)} \leq \frac{1}{\alpha}$$



$$\lim_{p \rightarrow \infty} S(p) = \frac{1}{\alpha}$$

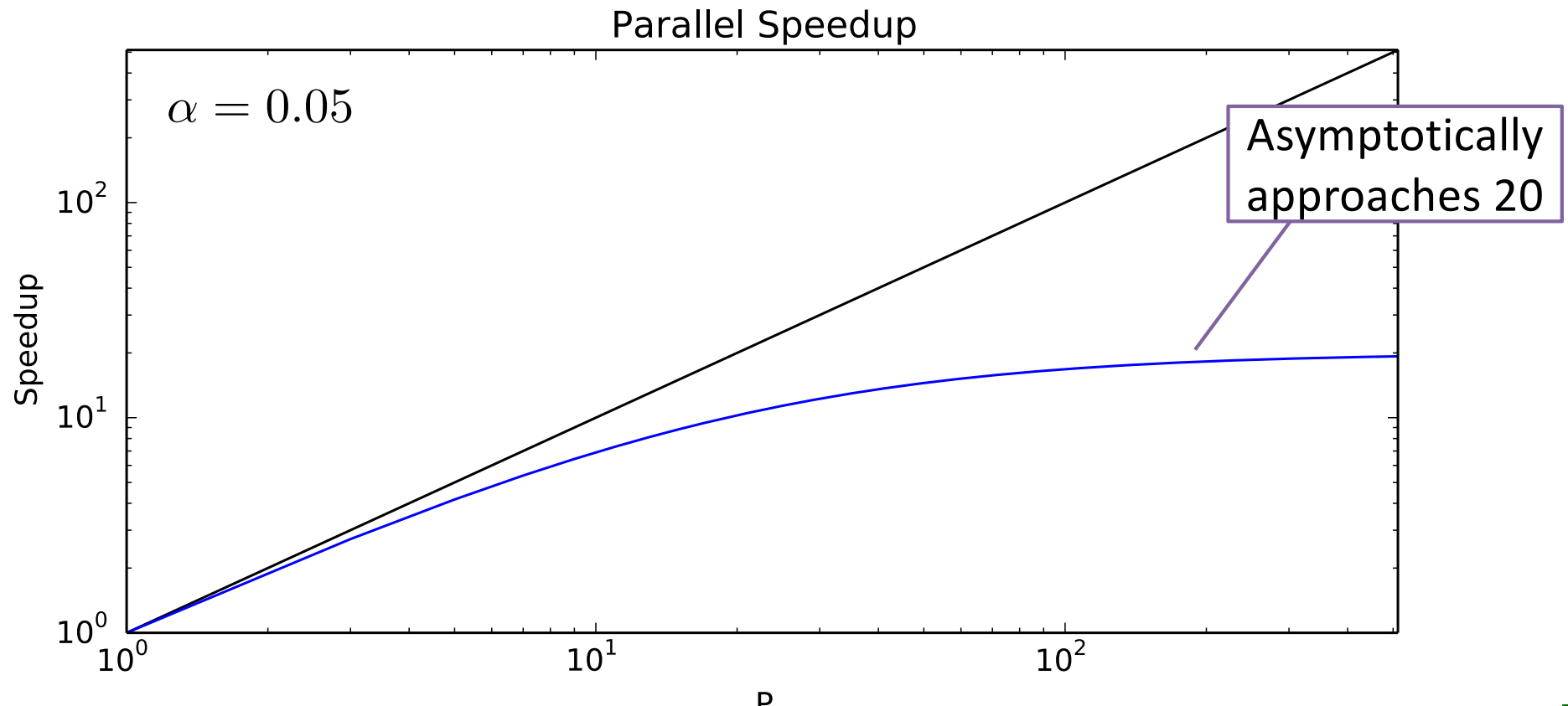
Limits to Parallelism (Amdahl's Law)



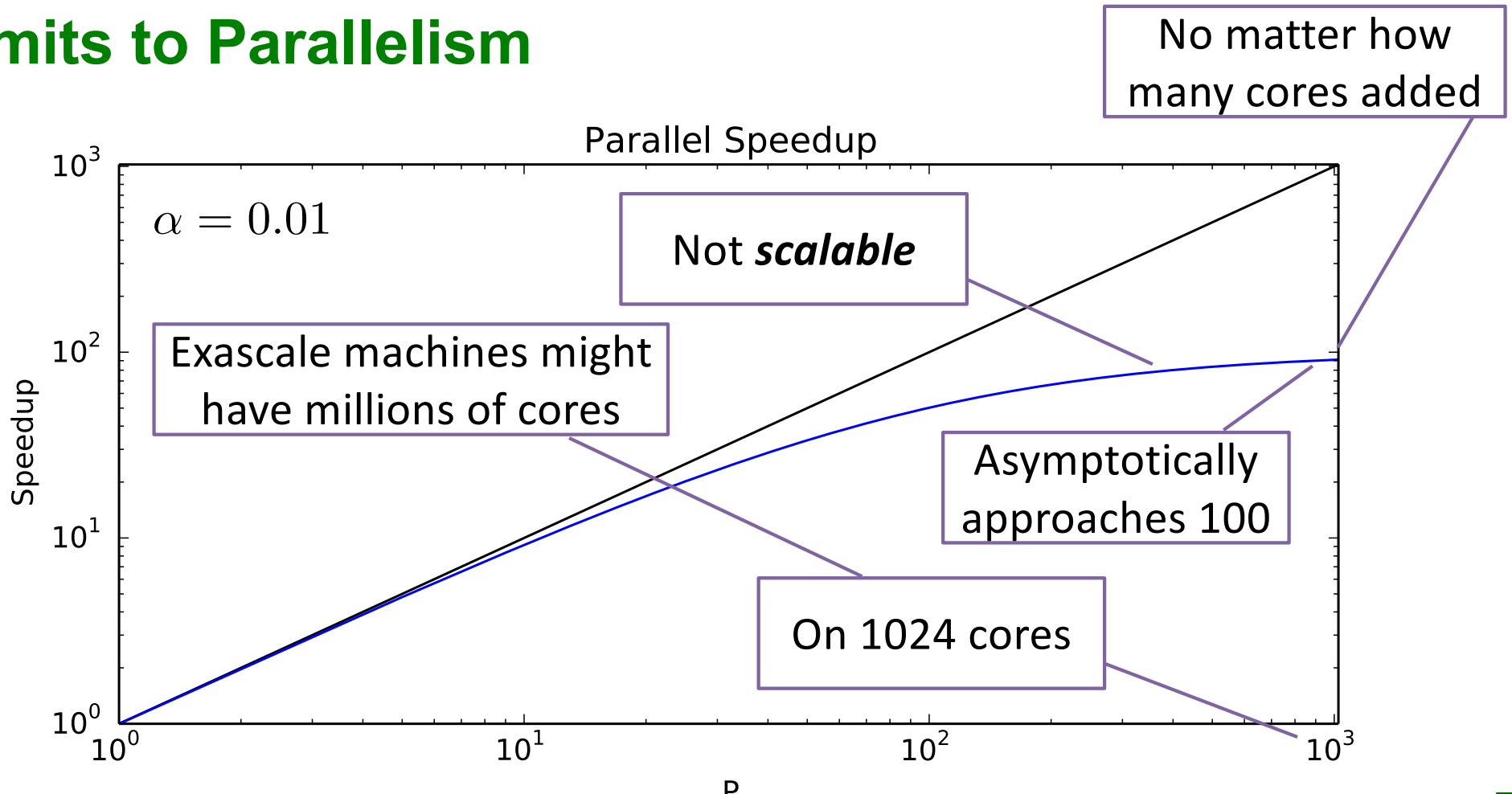
$$\lim_{p \rightarrow \infty} S(p) = \frac{1}{\alpha}$$

$$S(p) = \frac{T(n, 1)}{T(n, p)}$$

Limits to Parallelism (Amdahl's Law)

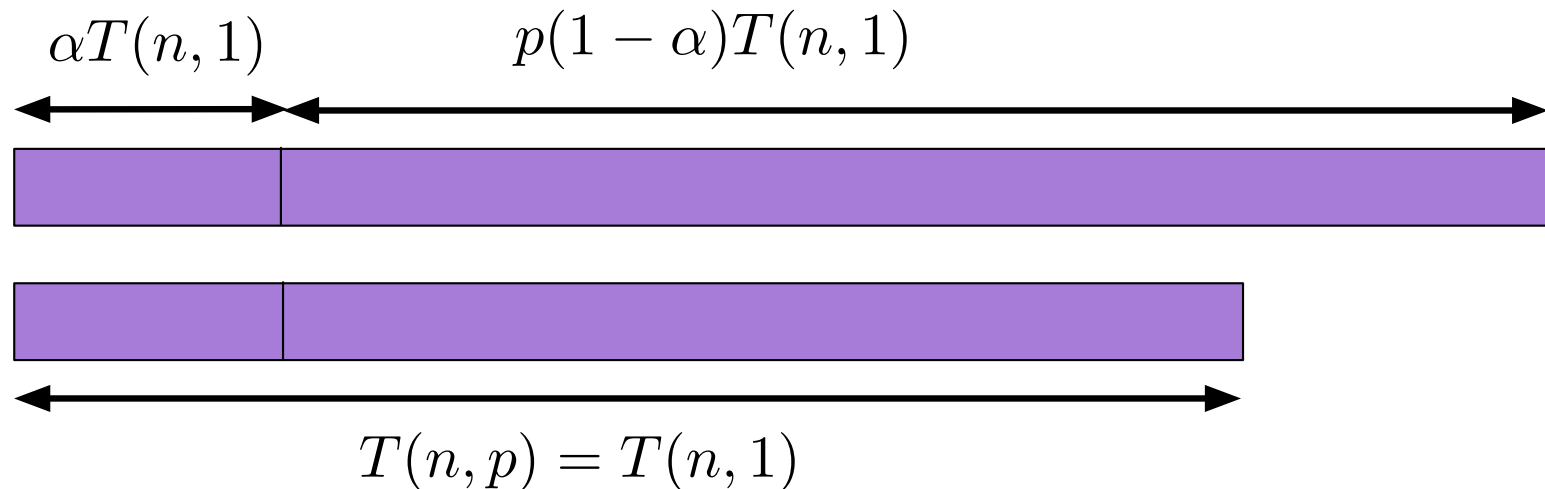


Limits to Parallelism

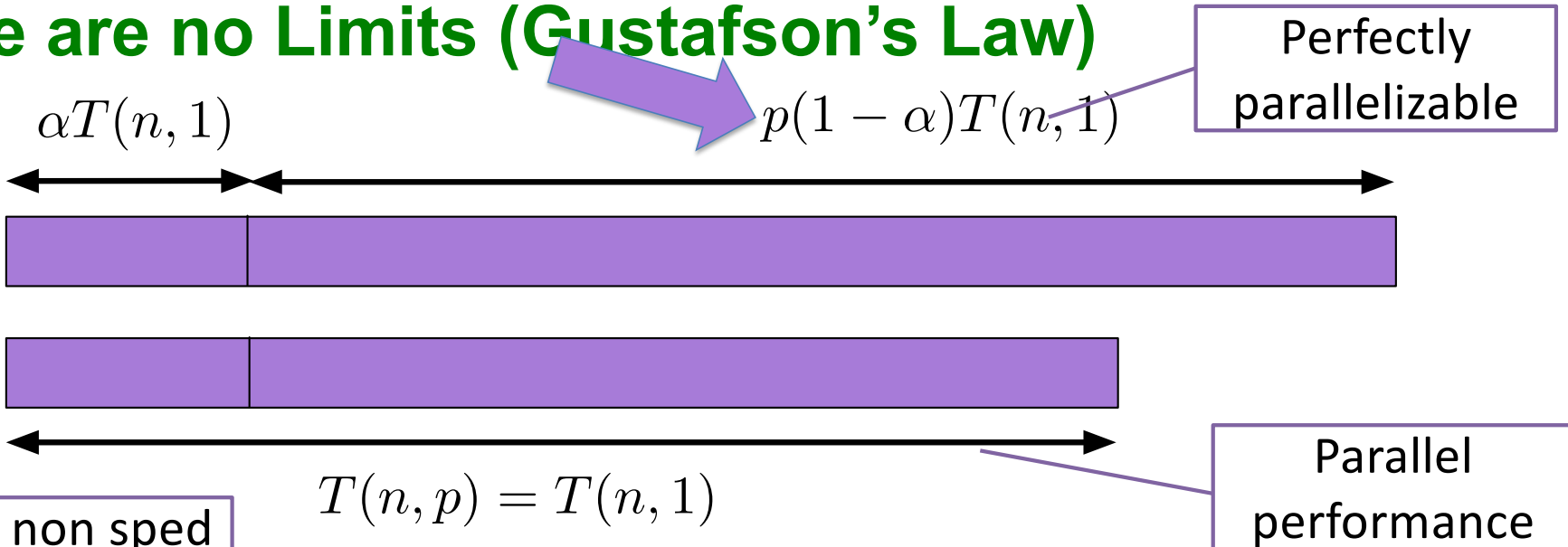


There are no Limits (Gustafson's Law)

- Doing the same problem faster and faster is not how we use parallel computers
- Rather, we solve bigger and more difficult problems
- I.e., the amount of parallelizable work grows



There are no Limits (Gustafson's Law)

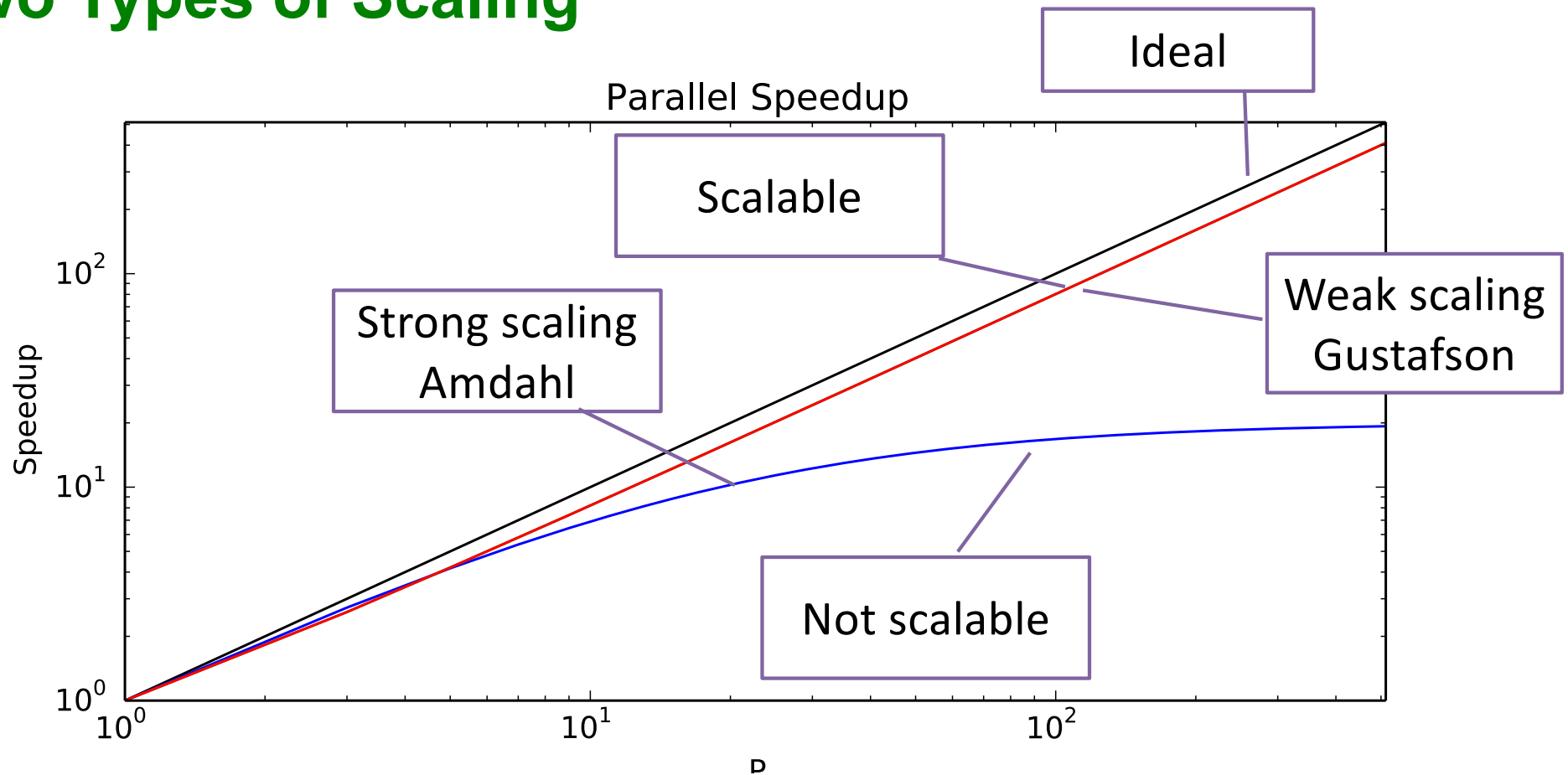


Ratio of non sped up to sped up

$$S(p) = \frac{\alpha T(n,1) + p(1-\alpha)T(n,1)}{T(n,p)} = \frac{\alpha T(n,1) + p(1-\alpha)T(n,1)}{T(n,1)} = \alpha + p(1 - \alpha)$$

$$E(p) = \frac{S(p)}{p} \quad \longrightarrow \quad \lim_{p \rightarrow \infty} E(p) = (1 - \alpha)$$

Two Types of Scaling



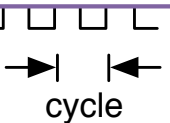
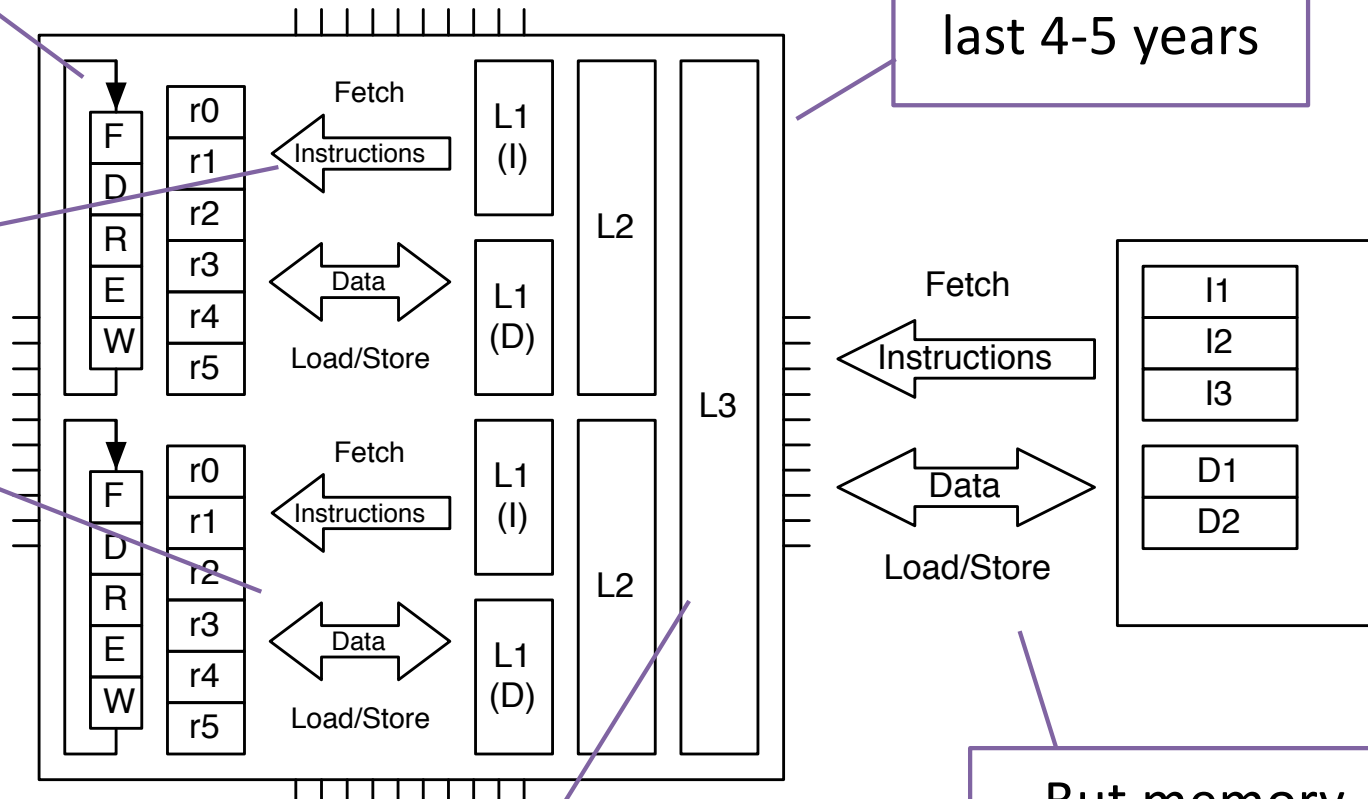
Multicore Architecture

Core is a
FDREW + regs

Each runs its
own sequence
of instructions

Each can access
its own data

Any CPU in the
last 4-5 years



But memory
might be shared

Each has memory
hierarchy

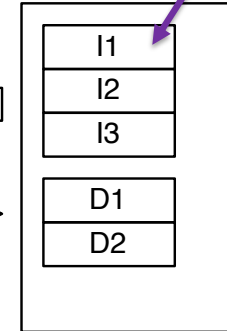
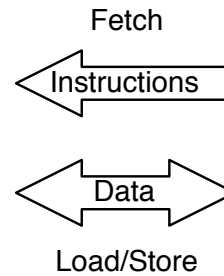
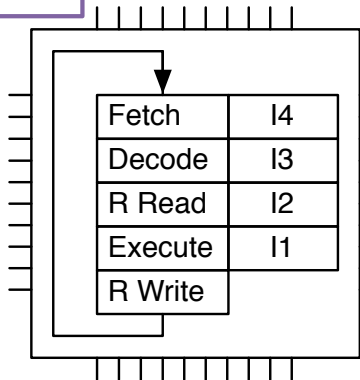
Running a Program

When a CPU is executing bytes from one program

It isn't executing bytes from another

Including from the OS (just another program)

Bytes from program stored in memory

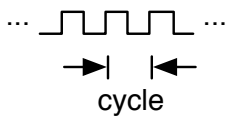


```

.global      __Z15hoistedMultiplyRKMatrixS1_RS_
.p2align    4, 0x0
__Z15hoistedMultiplyRKMatrixS1_RS_:
.cfi_startproc
## BB#0:
pushq      %rbp
Ltmp16:
.cfi_def_cfa_offset 16
Ltmp17:
.cfi_offset 4,%rbp,-16
movq       %rsp,%rbp
Ltmp18:
.cfi_def_cfa_register %rbp
pushq      %r15
pushq      %r14
pushq      %r13
pushq      %r12
pushq      %rbx
Ltmp19:
.cfi_offset 4,%rbx,-56
Ltmp20:
.cfi_offset 4,%r12,-48
Ltmp21:
.cfi_offset 4,%r13,-40
Ltmp22:
.cfi_offset 4,%r14,-32
Ltmp23:
.cfi_offset 4,%r15,-24
movq       8(%rsi),%rcx
testq     %rcx,%rcx
je         LBB2_9
## BB#1:
movq       16(%rsi),%r12
movq       8(%rdx),%rax
movq       %rax,-104(%rbp)
movq       16(%rdx),%rdx
movq       8(%rsi),%rax
movq       16(%rdi),%r13
leaq      -1(%rcx),%rsi
movq       %rsi,-88(%rbp)
movl      %ecx,%esi
    
```

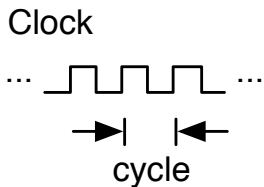
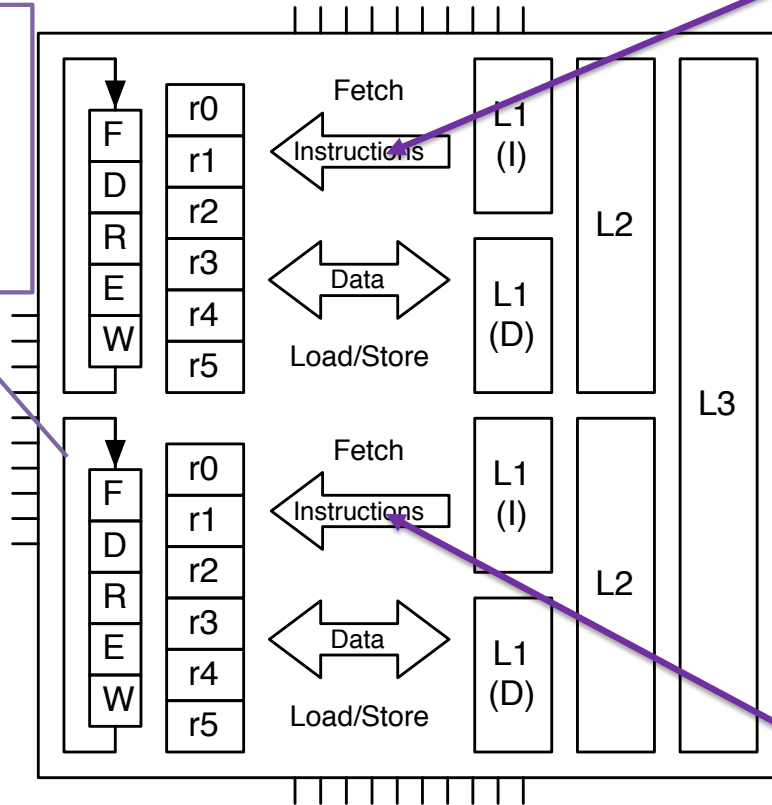
How does another program run?

How did the bytes get here?



Multicore Architecture

Cores are a resource managed by the OS



```

        .globl  __Z15hoistedMultiply80MatrixS1_RS_
        .p2align 4, 0x90
        __Z15hoistedMultiply80MatrixS1_RS_:
        .cfi_startproc
        ## BB#0:
        pushq   %rbp
        ltmp16: .cfi_def_cfa_offset 16
        ltmp17: .cfi_offset %rbp, -16
        movq    %rsp, %rbp
        ltmp18: .cfi_def_cfa_register %rbp
        pushq   %r15
        pushq   %r14
        pushq   %r13
        pushq   %r12
        pushq   %rbx
        ltmp19: .cfi_offset %rbx, -56
        ltmp20: .cfi_offset %r12, -48
        ltmp21: .cfi_offset %r13, -40
        ltmp22: .cfi_offset %r14, -32
        ltmp23: .cfi_offset %r15, -24
        movq   (%rdi), %rax
        movq   %rax, -120(%rbp)
        testq  %rax, %rax
        je     LBB2_9
        ## BB#1:
        movq   8(%rsi), %rcx
        testq  %rcx, %rcx
        je     LBB2_9
        ## BB#2:
        movq   16(%rsi), %r12
        movq   8(%r12), %rax
        movq   %rax, -104(%rbp)
        Fetch
        Instructions
        Data
        Load/Store
        L1 (I)
        L1 (D)
        L2
        L3
        Fetch
        Instructions
        Data
        Load/Store
        L1 (I)
        L1 (D)
        L2
        L3
        Fetch
        Instructions
        Data
        Load/Store
        L1 (I)
        L1 (D)
        L2
        L3
        I1
        I2
        I3
        D1
        D2
        .globl  __Z15hoistedMultiply80MatrixS1_RS_
        .p2align 4, 0x90
        __Z15hoistedMultiply80MatrixS1_RS_:
        .cfi_startproc
        ## BB#0:
        pushq   %rbp
        ltmp16: .cfi_def_cfa_offset 16
        ltmp17: .cfi_offset %rbp, -16
        movq    %rsp, %rbp
        ltmp18: .cfi_def_cfa_register %rbp
        pushq   %r15
        pushq   %r14
        pushq   %r13
        pushq   %r12
        pushq   %rbx
        ltmp19: .cfi_offset %rbx, -56
        ltmp20: .cfi_offset %r12, -48
        ltmp21: .cfi_offset %r13, -40
        ltmp22: .cfi_offset %r14, -32
        ltmp23: .cfi_offset %r15, -24
        movq   (%rdi), %rax
        movq   %rax, -120(%rbp)
        testq  %rax, %rax
        je     LBB2_9
        ## BB#1:
        movq   8(%rsi), %rcx
        testq  %rcx, %rcx
        je     LBB2_9
        ## BB#2:
        movq   16(%rsi), %r12
        movq   8(%r12), %rax
        movq   %rax, -104(%rbp)
    
```

A process may use one or more cores

Memory is not shared between processes