**CSCI-GA.3033-012**
# Multicore Processors: Architecture & Programming

# Lecture 5:    Overview of

# Parallel Programming
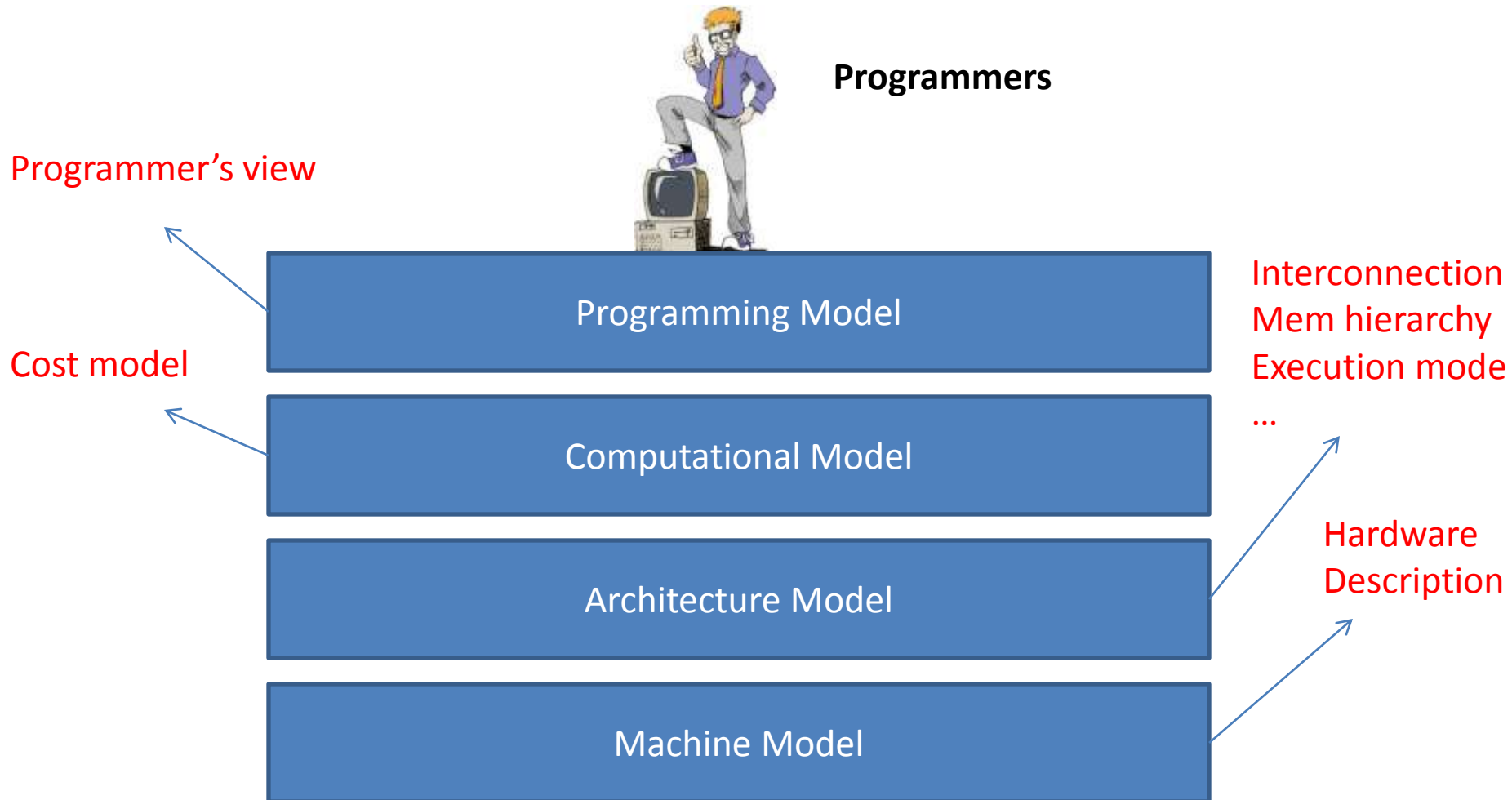
Mohamed Zahran (aka Z)

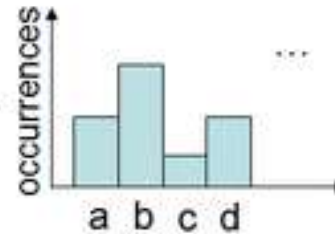mzahran@cs.nyu.edu

http://www.mzahran.com

# Models ... Models

**Programmers**

Programmer's view

Cost model

Interconnection
Mem hierarchy
Execution mode
...

Hardware
Description

| Programming Model |
| Computational Model |
| Architecture Model |
| Machine Model |

# Let's See A Quick Example

- **Problem**: Count the number of times each ASCII character occurs on a page of text.

- **Input**: ASCII text stored as an array of characters.

- **Output**: A histogram with 128 buckets – one for each ASCII character
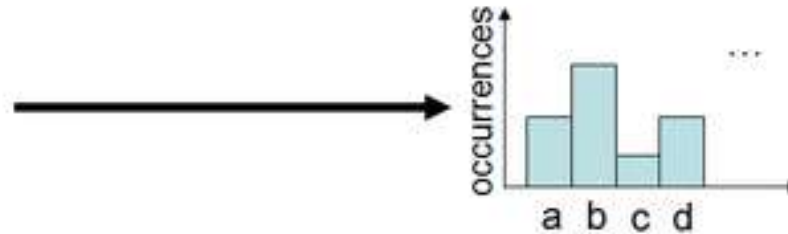
# Let's See A Quick Example
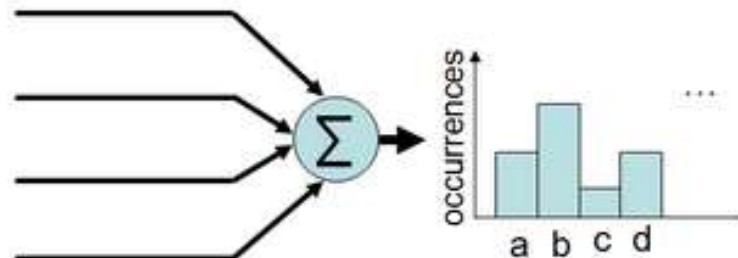
Speed on Quad Core:
10.36 seconds

```
1: void compute_histogram_st(char *page, int page_size, int *histogram){
2: for(int i = 0; i < page_size; i++){
3:     char read_character = page[i];
4:     histogram[read_character]++;
5:   }
6: }
```

**Sequential Version**

# Let's See A Quick Example



**We need to parallelize this.**

# Let's See A Quick Example

```
1: void compute_histogram_st(char *page, int page_size, int *histogram){
2: #pragma omp parallel for
3: for(int i = 0; i < page_size; i++){
4:      char read_character = page[i];
5:       histogram[read_character]++;
6:  }
```

**The above code does not work!!    Why?**

# Let's See A Quick Example

```
1: void compute_histogram_mt2(char *page, int page_size, int *histogram){
2: #pragma omp parallel for
3: for(int i = 0; i < page_size; i++){
4:     char read_character = page[i];
5:     #pragma omp atomic
6:      histogram[read_character]++;
7:    }
8: }
```

Speed on Quad Core:
 114.89 seconds
**> 10x slower than the single thread version!!**

# Let's See A Quick Example

```
1: void compute_histogram_mt3(char *page,
                                int page_size,
                                int *histogram, int num_buckets){
2: #pragma omp parallel
3: {
4:     int local_histogram[111][num_buckets];
5:     int tid = omp_get_thread_num();
6:     #pragma omp for nowait
7:     for(int i = 0; i < page_size; i++){
8:          char read_character = page[i];
9:          local_histogram[tid][read_character]++;
10:    }
11:    for(int i = 0; i < num_buckets; i++){
12:       #pragma omp atomic
13:       histogram[i] += local_histogram[tid][i];
14:    }
15:  }
16: }
```

Runs in 3.8 secs
Why speedup
is not 4 yet?

# Let's See A Quick Example

```
void compute_histogram_mt4(char *page, int page_size,
                                    int *histogram, int num_buckets){
1:          int num_threads = omp_get_max_threads();
2:          #pragma omp parallel
3:          {
4:          __declspec (align(64)) int local_histogram[num_threads+1][num_buckets];
5:          int tid = omp_get_thread_num();
6:          #pragma omp for
7:          for(int i = 0; i < page_size; i++){
8:                  char read_character = page[i];
9:                  local_histogram[tid][read_character]++;
10:         }
11:         #pragma omp barrier
12:         #pragma omp single
13:         for(int t = 0; t < num_threads; t++){
14:                 for(int i = 0; i < num_buckets; i++)
15:                         histogram[i] += local_histogram[t][i];
16:         }
17: }
```

Speed is
4.42 seconds.
Slower than the
previous version.

**source:** http://www.futurechips.org/tips-for-power-coders/writing-optimizing-parallel-programs-complete.html

# Let's See A Quick Example

```
void compute_histogram_mt4(char *page, int page_size,
                                  int *histogram, int num_buckets){
1:          int num_threads = omp_get_max_threads();
2:          #pragma omp parallel
3:          {
4:          __declspec (align(64)) int local_histogram[num_threads+1][num_buckets];
5:          int tid = omp_get_thread_num();
6:          #pragma omp for
7:          for(int i = 0; i < page_size; i++){
8:                      char read_character = page[i];
9:                      local_histogram[tid][read_character]++;
10:         }
11:
12:         #pragma omp for
13:         for(int i = 0; i < num_buckets; i++){
14:                     for(int t = 0; t < num_threads; t++)
15:                                 histogram[i] += local_histogram[t][i];
16:         }
17: }
```

Speed is
3.60 seconds.

# What Can We Learn from the Previous Example?

- Parallel programming is not only about finding a lot of parallelism.
- Critical section and atomic operations
  - Race condition
  - Again: correctness vs performance loss
- Know your tools: language, compiler and hardware

# What Can We Learn from the Previous Example?

- Atomic operations
  - They are expensive
  - Yet, they are fundamental building blocks.
- Synchronization:
  - correctness vs performance loss
  - Rich interaction of hardware-software tradeoffs
  - Must evaluate hardware primitives and software algorithms together

# Sources of Performance Loss in Parallel Programs

- Extra overhead
  - code
  - synchronization
  - communication
- Artificial dependencies
  - Hard to find
  - May introduce more bugs
  - A lot of effort to get rid of
- Contention due to hardware resources
- Coherence
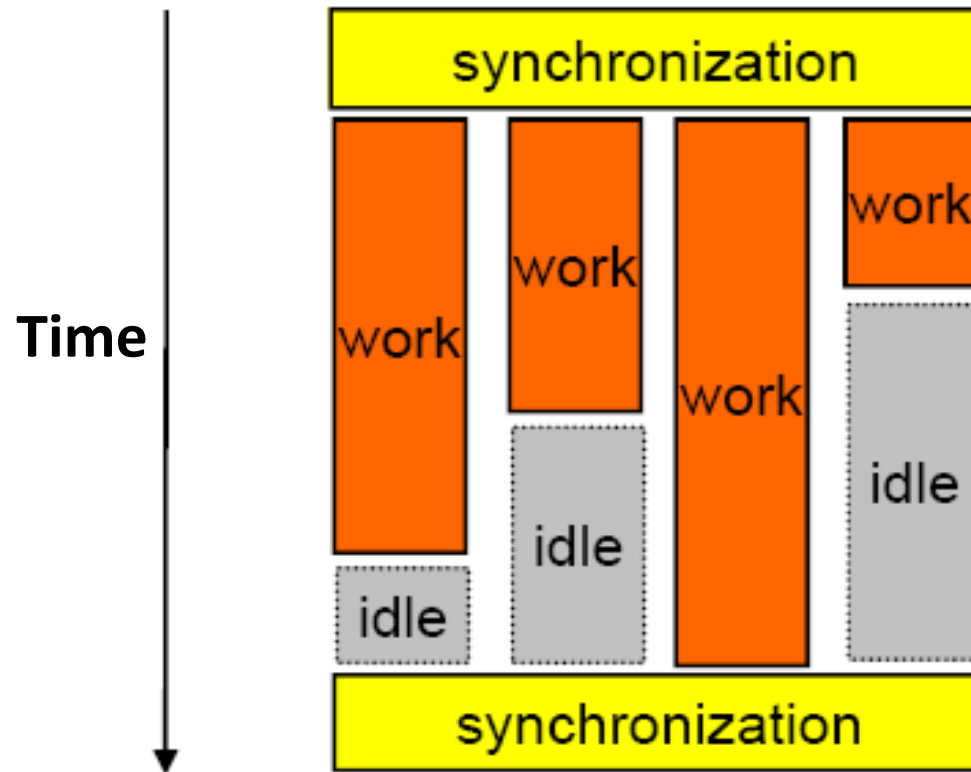- Load imbalance

# Artificial Dependencies

```
int  result;
//Global variable

for (...)  // The OUTER loop
     modify_result(...);
      if(result > threshold)
           break;

void modify_result(...)

   ...
    result = ...
```

What is wrong with that program when we try to paralleize it?

# Coherence

- Extra bandwidth (scarce resource)
- Latency due to the protocol
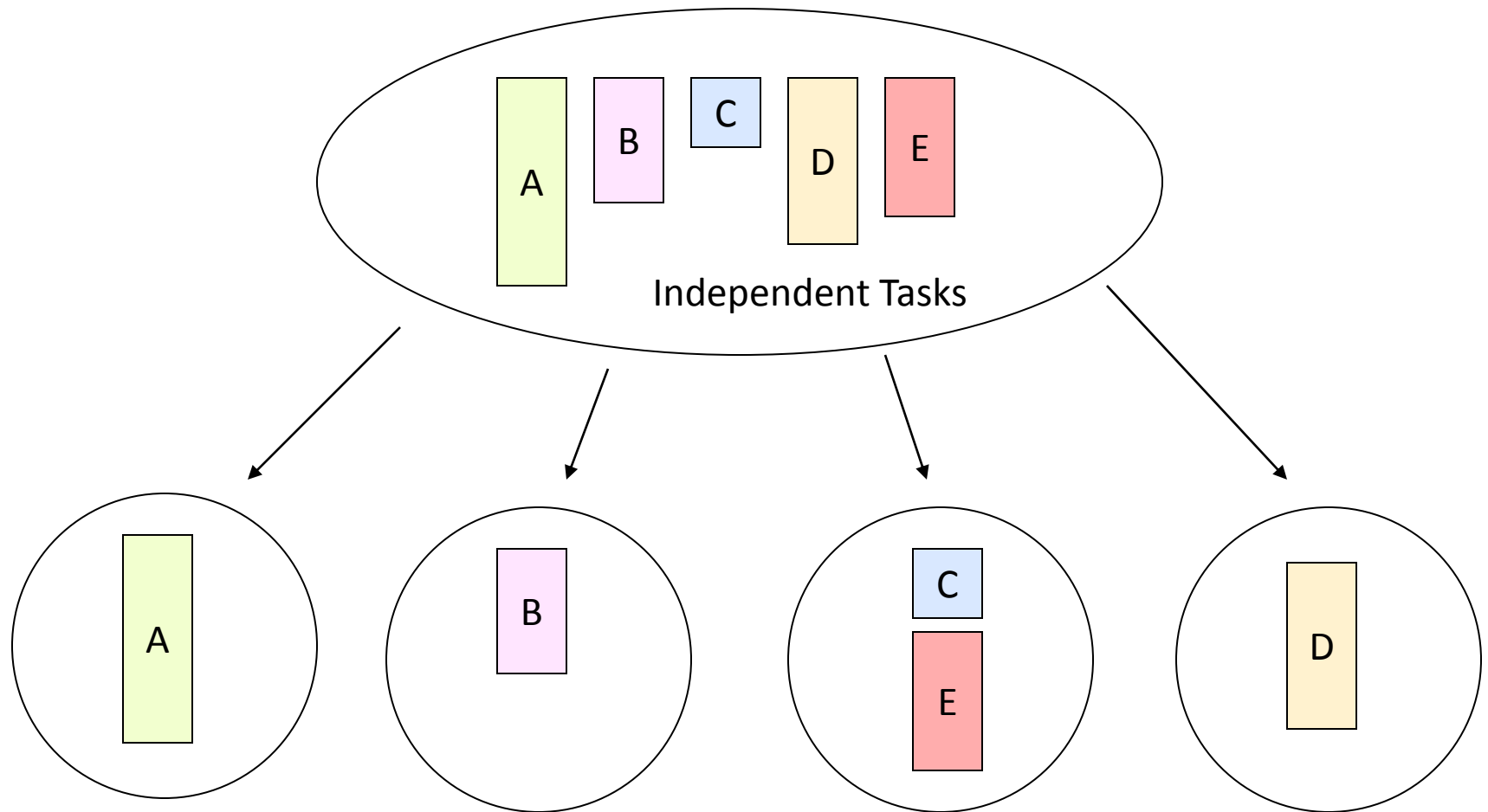- False sharing

# Load Balancing

# Load Balancing

- Assignment of work not data is the key
- If you cannot eliminate it, at least reduce it.
- Static assignment
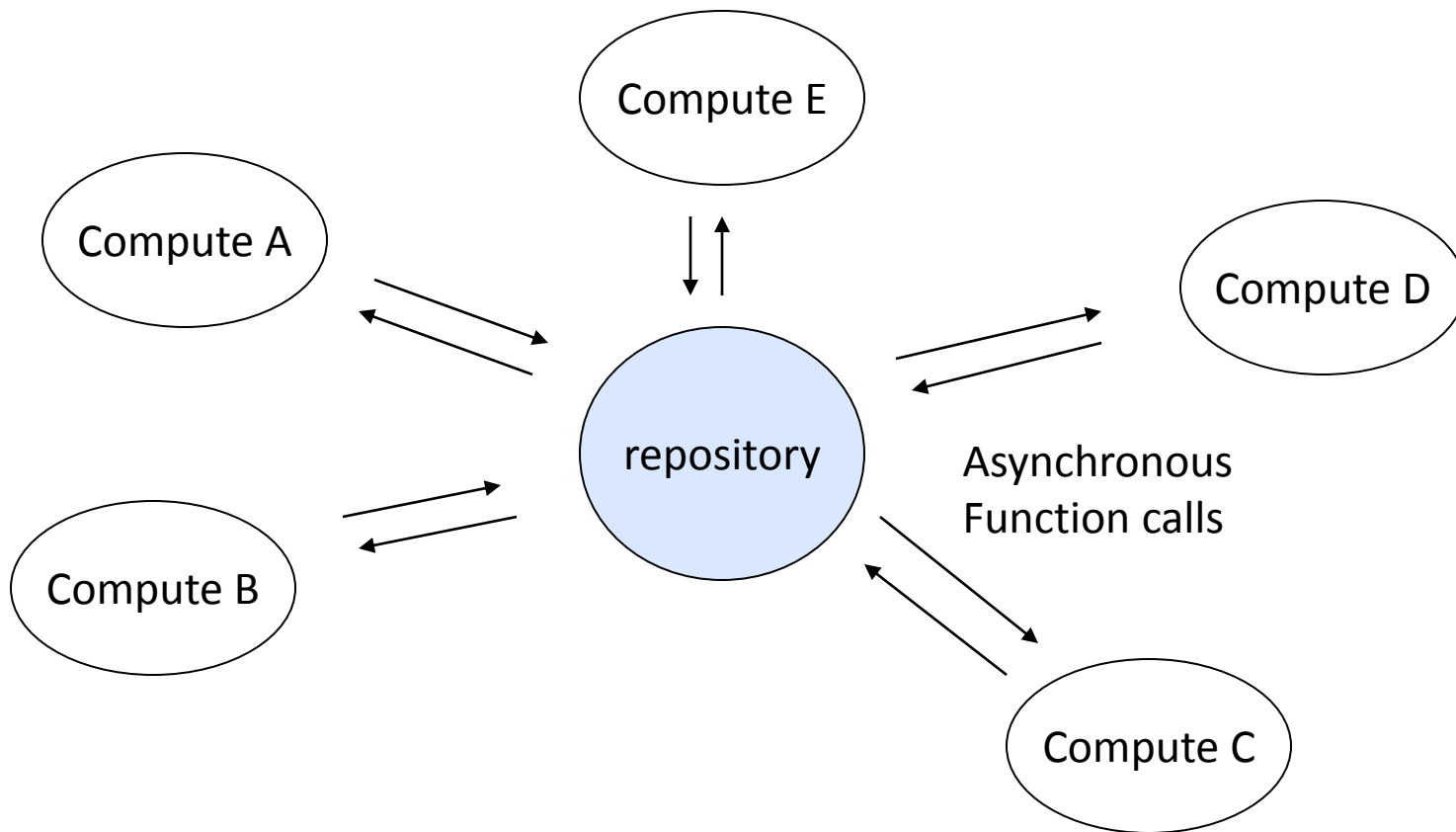- Dynamic assignment
  - Has its overhead

# Patterns in Parallelism

- Task-level (e.g. Embarrassingly parallel)
- Divide and conquer
- Pipeline
- Iterations (loops)
- Client-server
- Geometric (usually domain dependent)
- Hybrid (different program phases)

# Task Level



Independent Tasks

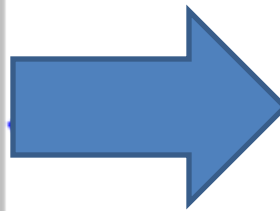# Client-Server/ Repository

# Example

Assume we have a large array and we want to compute its minimum (T1), average (T2), and maximum (T3).

```
#define maxN 1000000000

int m[maxN];
int i;
int min = m[0];
int max = m[0];
double avrg = m[0];

for(i=1; i < maxN; i++) {
    if(m[i] < min)
      min = m[i];
    avrg = avrg + m[i];
    if(m[i] > max)
      max = m[i];
}
avrg = avrg / maxN;
```

```
#define maxN 1000000000
int m[maxN];

int i; int min = m[0];
for(i=1; i < maxN; i++) {
    if(m[i] < min)
      min = m[i];
}                                    T1

int j;
double avrg = m[0];
for(j=1; j < maxN; j++) {
    avrg = avrg + m[j];
}
avrg = avrg / maxN;                  T2

int k; int max = m[0];
for(k=1; k < maxN; k++) {
    if(m[k] > max)
      max = m[k];
}                                    T3
```
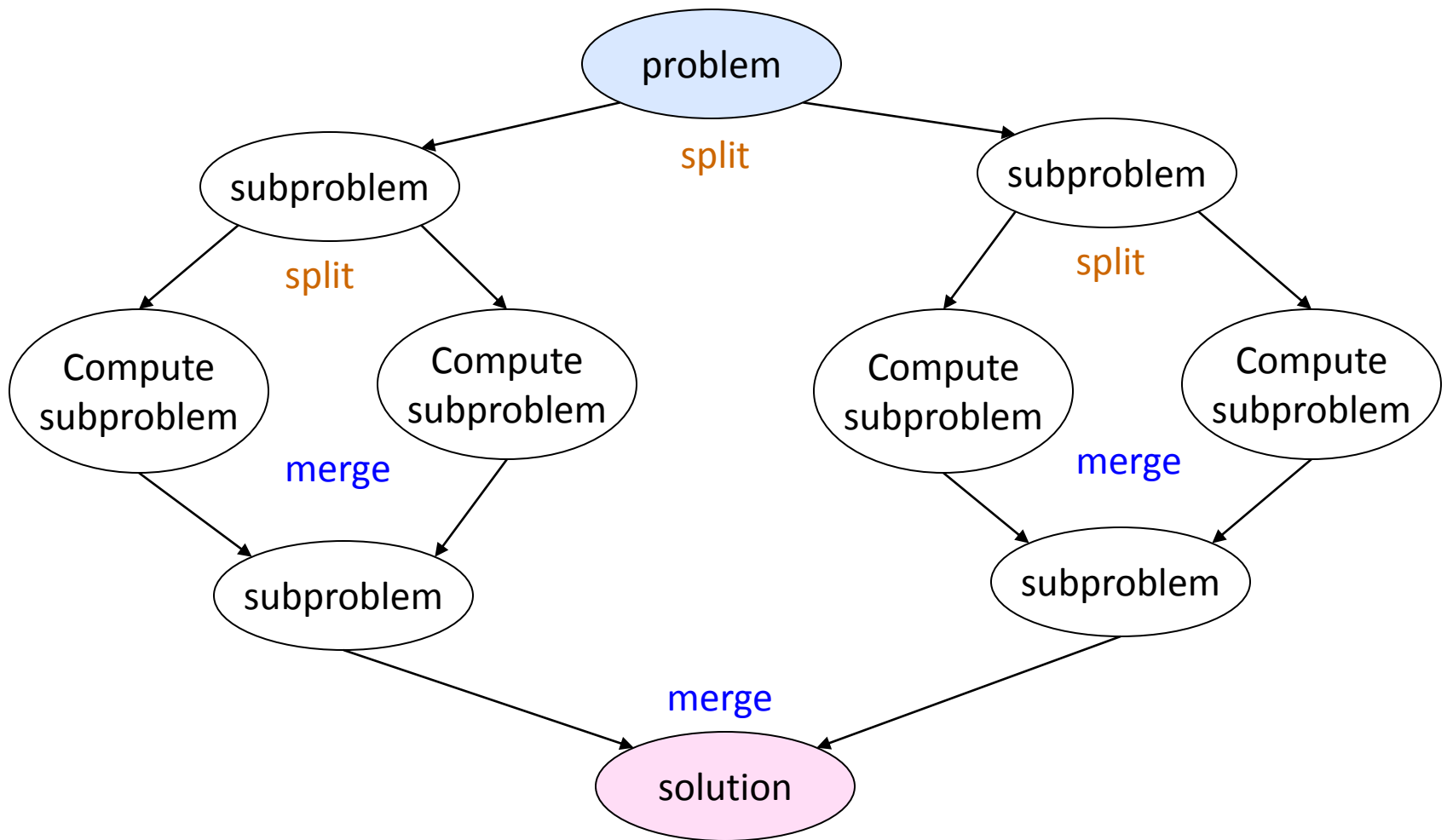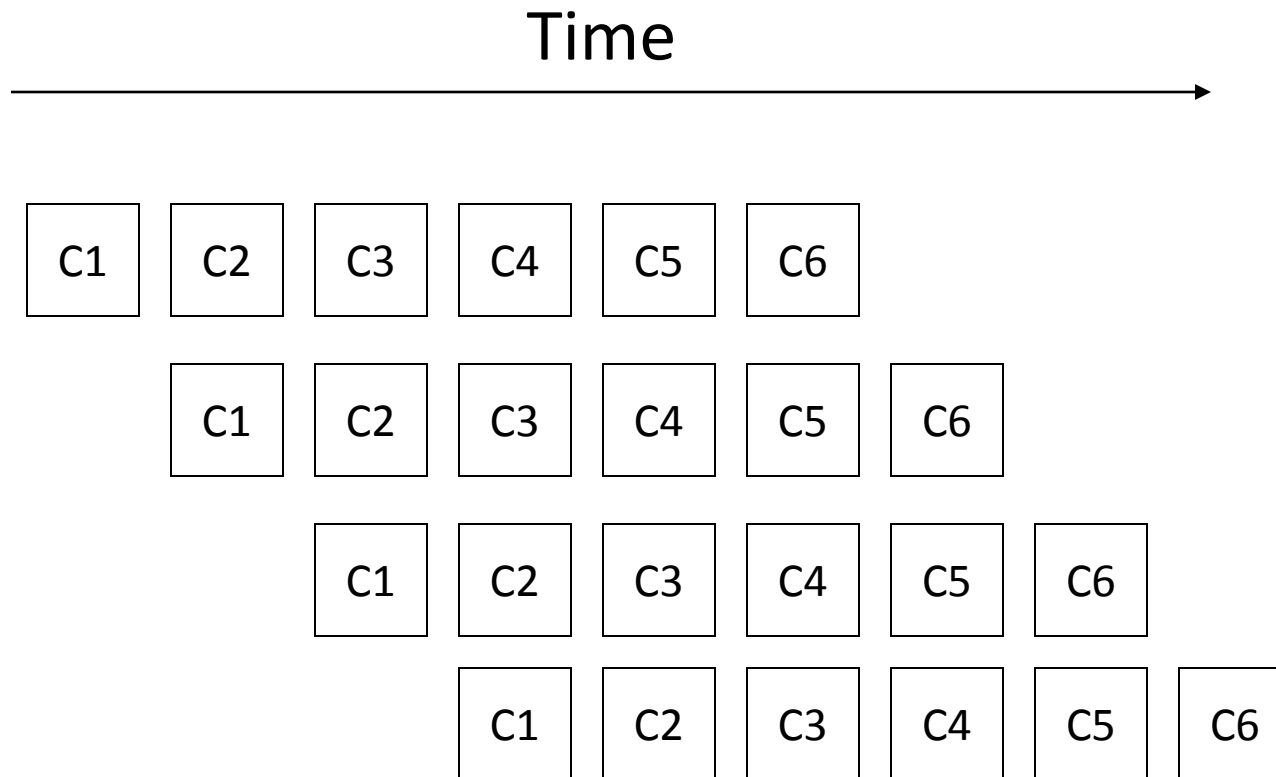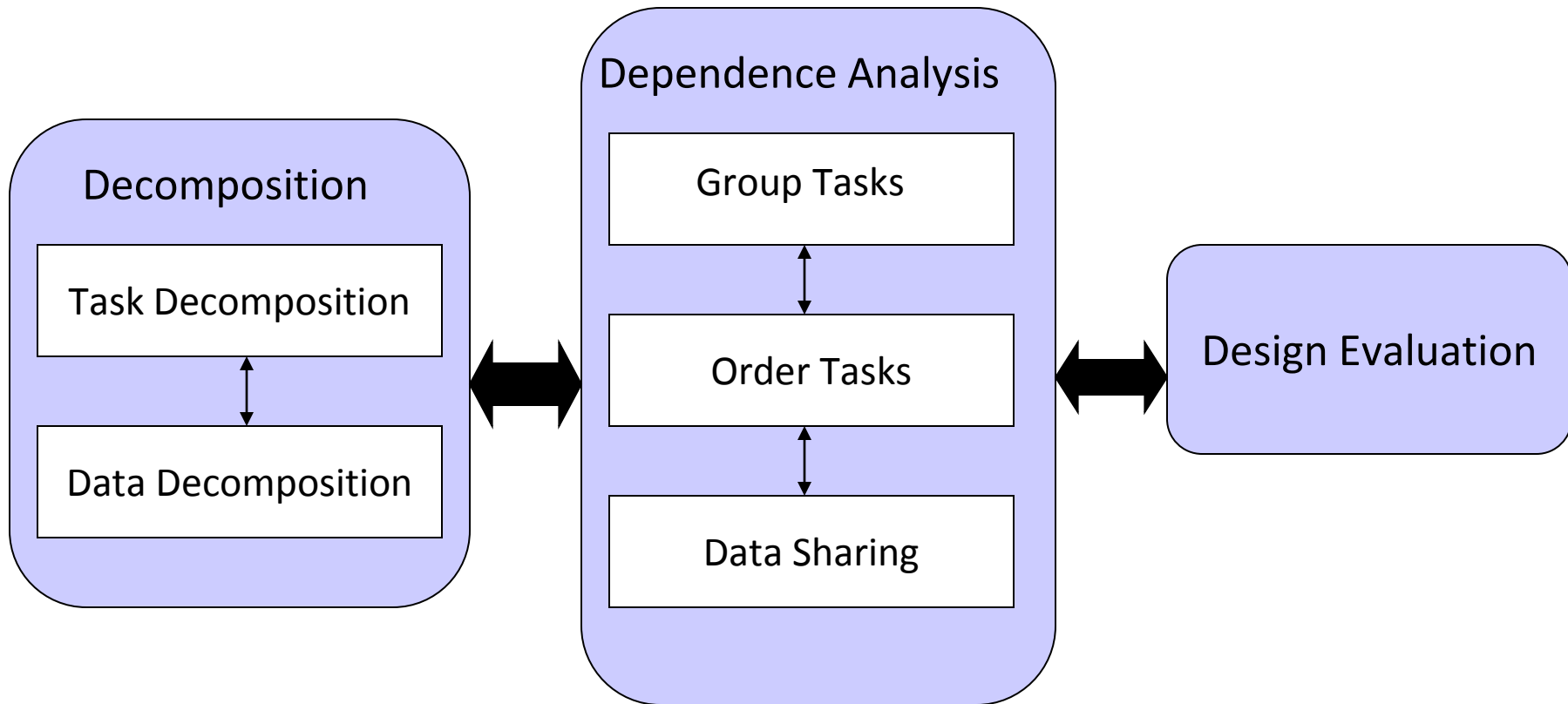
# Divide-And-Conquer

# Pipeline

A series of ordered but independent computation stages need to be applied on data.

Time

# Pipeline

- Useful for
  - streaming workloads
  - Loops that are hard to parallelize
    - due inter-loop dependence
- Usage for loops: split each loop into stages so that multiple iterations run in parallel.
- Advantages
  - Expose intra-loop parallelism
  - Locality increases for variables uses across stages
- How shall we divide an iteration into stages?
  - number of stages
  - inter-loop vs intra-loop dependence
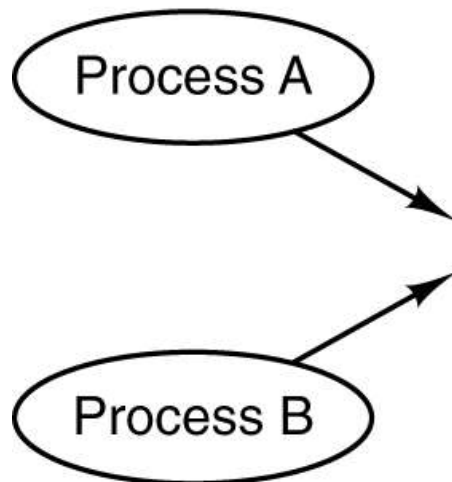
# The Big Picture of Parallel Programming

**Decomposition**

- Task Decomposition
- Data Decomposition

**Dependence Analysis**

- Group Tasks
- Order Tasks
- Data Sharing

**Design Evaluation**

# BUGS

- Sequential programming bugs + more
- Hard to find
- Even harder to resolve ☹
- Due to many reasons:
  - example: race condition

# Example of Race Condition

1. Process A reads **in**
2. Process B reads **in**
3. Process B writes file name in slot 7
4. Process A writes file name in slot 7
5. Process A makes **in = 8**

RACE CONDITION!!

Spooler directory

| | |
|---|---|
| 4 | abc |
| 5 | prog.c |
| 6 | prog.n |
| 7 | |
| | |

out = 4

in = 7

Process A

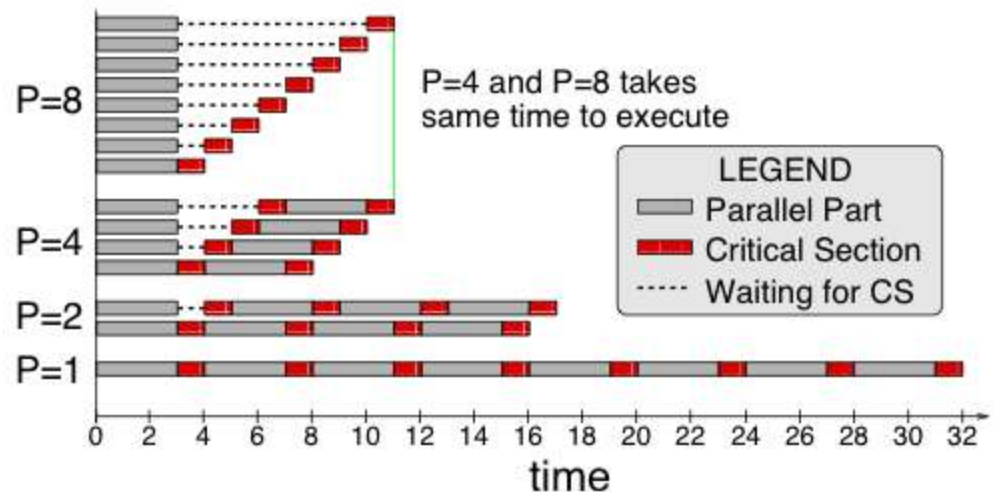Process B

# How to Avoid Race Condition?

- Prohibit more than one process from reading and writing the shared data at the same time -> <span style="color:red">mutual exclusion</span>

- The part of the program where the shared memory is accessed is called the <span style="color:red">critical region</span>



P=4 and P=8 takes
same time to execute

LEGEND
- Parallel Part
- Critical Section
- Waiting for CS
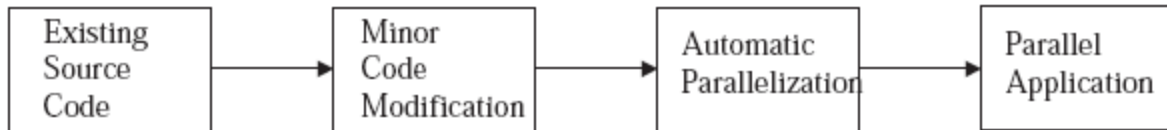
# Conditions of Good Solutions to Race Condition

1. No two processes may be simultaneously inside their critical region
2. No assumptions may be made about speeds or the number of CPUs/Cores
3. No process running outside its critical region may block other processes
4. No process has to wait forever to enter its critical region

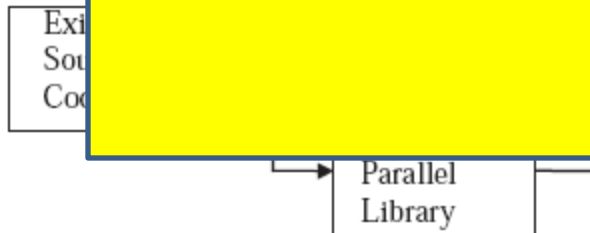# Importance Characteristic of Critical Sections

- How severe a critical section on performance depends on:
  - The position of the critical section (in the middle or at the end)
  - Kernel executed on the same or different core(s)
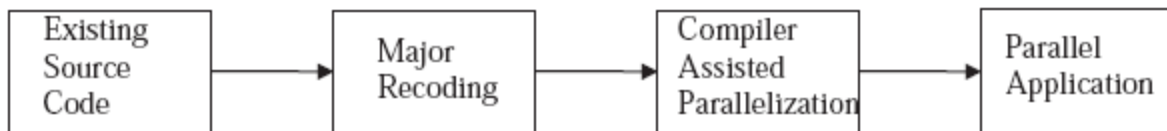
# Traditional Way of Parallel

**Strategy 1: Automatic Parallelization**

| Existing Source Code | → | Minor Code Modification | → | Automatic Parallelization | → | Parallel Application |

**Do We Have To Start With Sequential Code?**

Existing
Source
Code

Parallel
Library

**Strategy 3: Major Recoding**

| Existing Source Code | → | Major Recoding | → | Compiler Assisted Parallelization | → | Parallel Application |

# Conclusions

- Pick your programming model
- Task decomposition
- Data decomposition
- Refine based on:
  - What compiler can do
  - What runtime can do
  - What the hardware provides