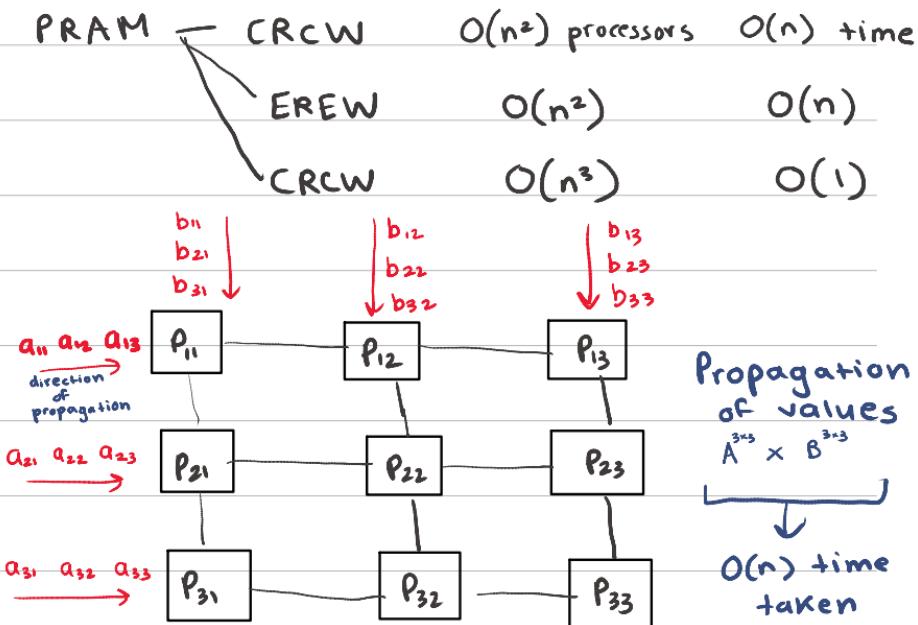


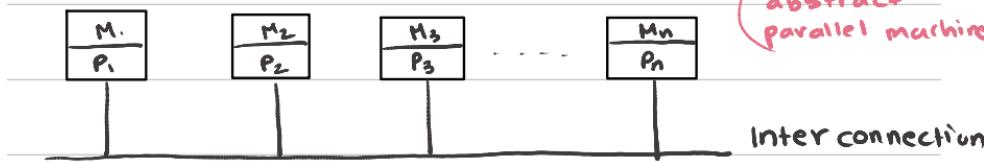


Date: 28/8/23

* Matrix Multiplication



* Bulk Synchronous Parallel Model (BSP)



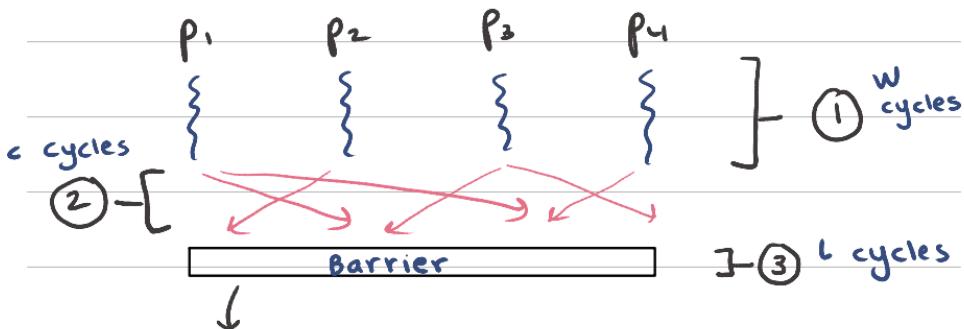
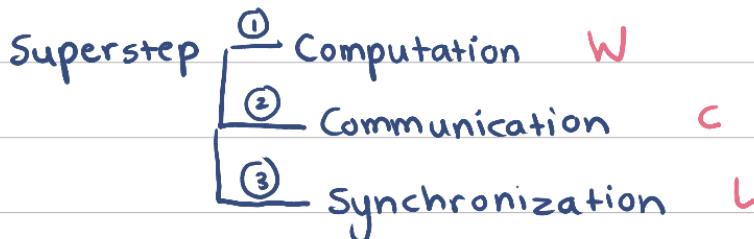
- * Each processor has its own memory (Unlike PRAM → shared memory)

Date: /

Application

* there are n processes; 1 per node

* Each process executes one superstep sequentially



• time taken by $P_i = w_i + c_i + l_i$

$w \approx$ number of operations

$l \approx$ depends on hardware/algo

$$c = gh$$

time/word
(hardware dependant)

max no. of words that can be sent/received by a processor in a superstep

Date: /

Example: $p=4$, Adding n numbers ($n \gg p$)

• As per model, # processes = 4 ($= p$)

(# Steps) → Comp.

Comm.

Sync.

Superstep 1:

$\frac{N}{4}$

Adding the
 $\frac{N}{4}$ nums

1 (g cycles)



1 (L cycles)

Superstep 2:

1

$(S_1 + S_2)$ and
 $(S_3 + S_4)$

1 (g cycles)



1 (L cycles)

Superstep 3:

1

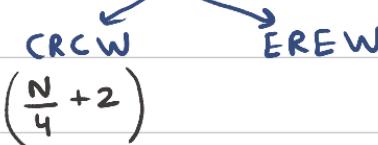
$S_{12} + S_{34}$

—

—

$$\text{Total steps : } \left(\frac{N}{4} + 2\right) + 2g + 2L$$

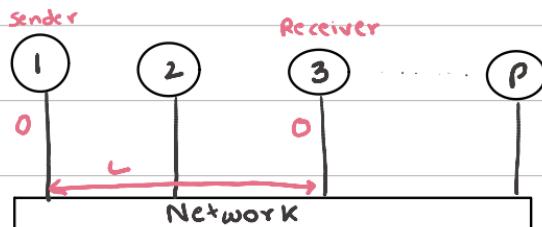
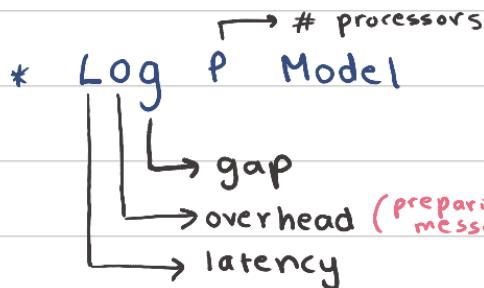
* Same problem on PRAM



- no comm.
- can use sum CW

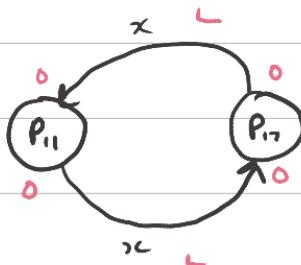
$$\left(\frac{N}{4} + 2\right)$$

Date: 29/8/23



overhead: length of time that a processor is involved in the transmission/reception of msg.

Eg:



$$\text{Total time} = 40 + 2L$$

latency: length of time taken to send a message from one processor to another.

Date:

gap: min time interval between consecutive message termination/reception.

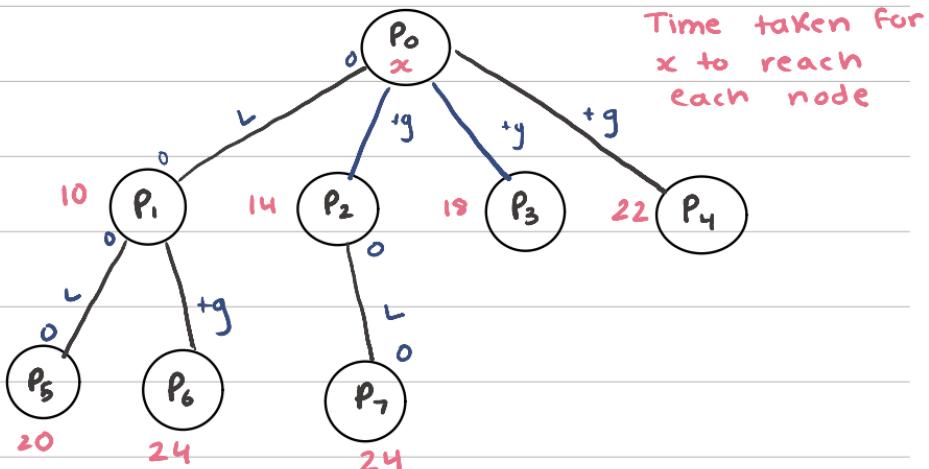
Eg: P_i pumps out messages at intervals

0 g 2g 3g ...

$$\frac{L}{g} = \# \text{ messages in the network}$$

* synchronization happens within sender-receiver pairs (not b/w ALL processors)

Example: $p=8$, $L=6$, $g=4$, $O=2$



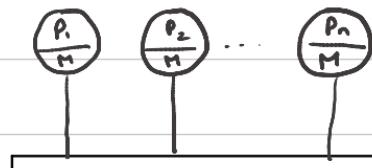
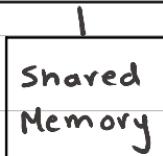
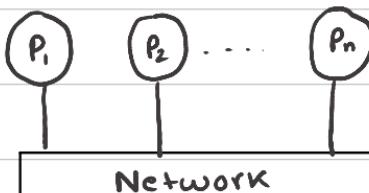
FOR EXAM

Date:

Chapters Completed - 1, 2, 7, 8, 11

* Parallel Programming

(Actual pgm. implementation)



Focus for now

Example pseudo code:

Processor 0

Array A[1:N]

BEGIN

For i=1 to N do

Read A[i]

End For

send (1, A $\left[\frac{N}{2}+1\right]$, $\frac{N}{2}$)

s

Processor 1

Array B[1: $\frac{N}{2}$]

BEGIN

Receive (0, B[i], $\frac{N}{2}$)

set suml = 0

for i = 1 to $\frac{N}{2}$ do

suml =

suml + B[i]

end for

Date:

set sum:=0

for i = 1 to $\frac{N}{2}$ do

sum = sum + A[i]

end for

Receive (1, sum1, 1)

sum = sum + sum1

print (sum)

END

send (0, sum), 1)

END

Note: * for every send, there MUST be a receive
(and vice-versa)

* only one program is written for all processors.

Date: 01/9/23

* Parallel Programming

→ message passing parallel computers

→ send()
receive()

} interprocess
communication

code
portability

→ a single program is written for all processors

→ Message Passing Interface (MPI)

program starts with this

library of standard message passing routines

→ MPI-Init()

} setting up & tearing down

the environment (creation/
deletion)

ends with this → MPI-Finalize()

MPI-Comm-size() → how many

processors

MPI-Comm-rank() → rank of each

processor

MPI-send()

} for the actual message

MPI-receive()

passing

• Single Program Multiple Data (SPMD)

↳ it is a style of programming

• Note: MPI is an interface, the implementation of each routine depends on the vendor

Date: /

Example Program: (running on n processors)

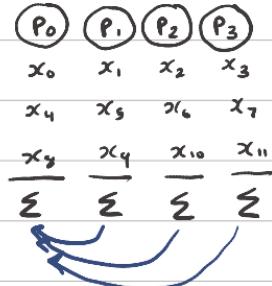
```
#include <stdio.h>
* #include <mpi.h>

int main() {
    int i, x(i), tmp, sum=0,
        group-size, my-rank, N;
```

* MPI_Init();

* MPI_Comm_size (MPI_Comm_World, &group-size); sum of numbers

* MPI_Comm_rank (MPI_Comm_World, &my-rank);



Note: group-size = # processors
my-rank = assigned differently for each processor that runs this code

4/9/23

```
if (my-rank == 0)    → program flow
{                  for P0
    for (i=1; i < group-size; i++)
        MPI_send(&N, 1, MPI_INT, i, i, MPI_Comm_World);
        start   #items  datatype dest tag
        addr
```

① Point-to-point communication

② Collective communication

Date: /

Assuming each processor already has access
to all elements

for (i = my_rank, i < N; i = i + group_size)

 sum = sum + x(i);

for (i = 1; i < group_size; i++)

{ MPI_Recv(&tmp, 1, MPI_INT, i, MPI_COMM_WORLD, &status);
 ^I receiving same as send
 addr

 sum = sum + tmp;

 print(sum);

}

} → END IF

else { ^{program flow for p; i ≠ 0}

 MPI_Recv(&N, 1, MPI_INT, 0, i, MPI_COMM_WORLD, &status);
 ^I src

 for (i = my_rank; i < N; i = i + group_size)

 sum = sum + x(i);

 MPI_Send(&sum, 1, MPI_INT, 0, i, MPI_COMM_WORLD);
 ^I dst

}

 MPI_Finalize()

} → PROGRAM END

Date: /
* For Strings

MPI-Send (string, $\frac{10}{I}$, MPI-CHAR, ...)
length of string

MPI-Recv (string, $\frac{10}{I}$, MPI-CHAR, ...)

Important

→ to handle the heterogeneity across diff architectures of storing different datatypes

Eg: $P_0 \rightarrow$ stores int as 2 Bytes

$P_1 \rightarrow$ stores int as 4 Bytes

→ also for underlying network

Processor X

Processor Y

Send (M, 64, Y, tag1) | Recv (P, X, 64, tag1)

Send (M, 32, Y, tag2) | Recv (Q, X, 32, tag2)

* purpose of tags is to handle ambiguity at the receiver's end.

* we cannot guarantee in-order delivery

Date: 5/9/23

- In log P model, there is no explicit synchronization stage

→ But, pairwise synchronization happens during communication

Q. How do you send multiple non-consecutive elements using MPI_send()?

because $\overset{+}{\text{send}}(\text{start}, \# \text{items}, \dots)$

∴ sending $A[0], A[5], A[10] \dots$ together is an issue

Ans: Copy the elements to a new buffer/array

⇒ it becomes consecutive ⇒ use MPI_send()



- 1) Synchronous_send() → waits till the elements are received before moving on.
- 2) Asynchronous_send() → gives the elements to the output buffer & moves on

Sync

→ reliability

→ more wait time

→ doesn't req. output buffer

Asynch

→ no reliability assurance

→ no wait time

→ output buffer is req.

Date: /

* Collective Communications

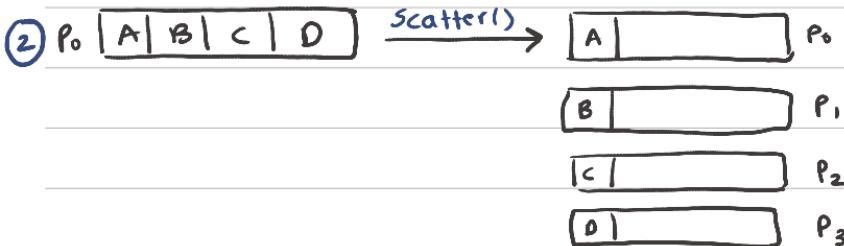


→ instead of using n-1 separate calls, data is sent to all processors in the comm-world.

Adv

→ not just reduced # statements.

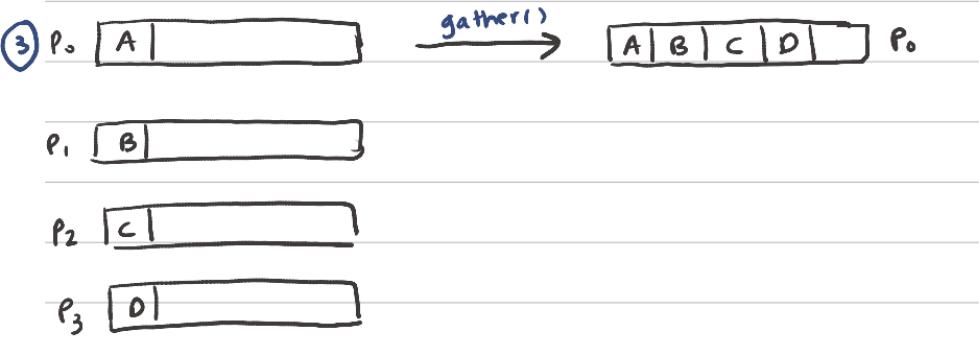
→ the vendor may optimize Broadcast() to $O(\log n)$
∴ It definitely has a speedup over n-1 send() statements.



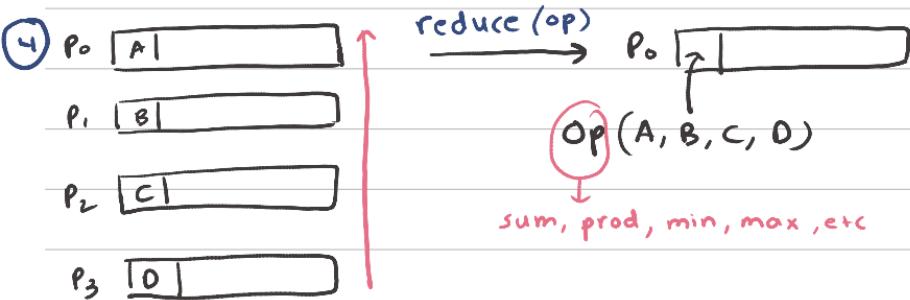
→ "personalized" communication

Date: /

→ the order is which the elements are scattered
is determined by the ranks



→ collecting values from all processes into one.

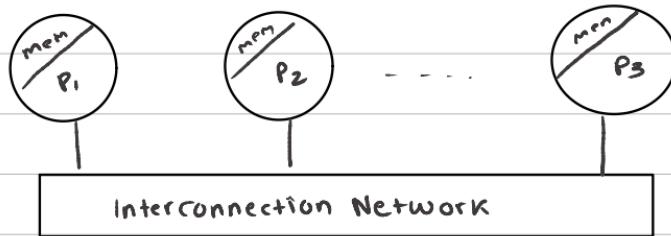


→ performs some sort of aggregation across all processors

⑤ MPI_Barrier() → to make all processes in the MPI_World hit the synchronization barrier

Date: 11 / 9 / 23

* Message Passing Parallel Computer



MPI

- >Create Terminate
- Who am I? rank
- Function
 - point to point communication
 - collective communication

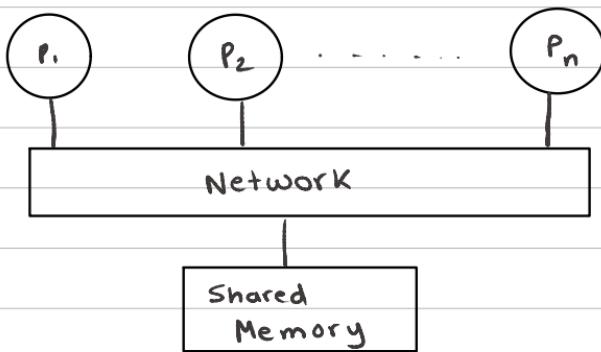
{ { } } → MPI-Barrier required for synchronization

Barrier

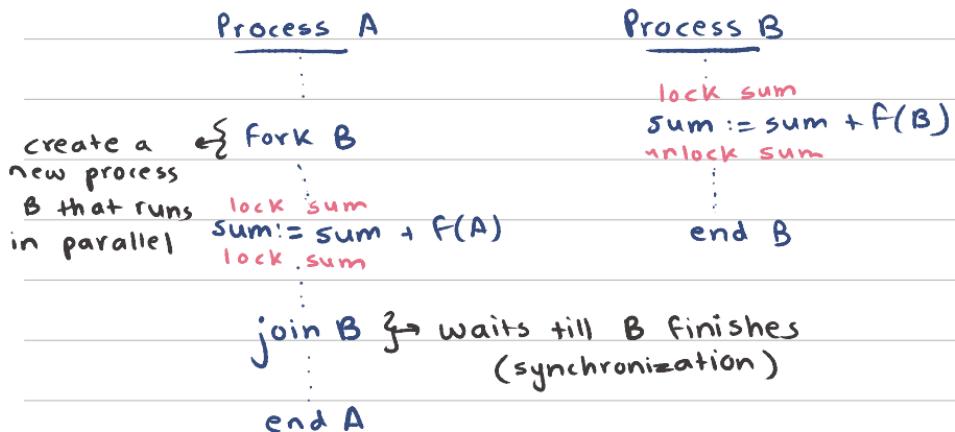
- * In collective communication, the barrier is implied.
- * In MPI, all functions are bi-directional

Date: /

* Shared Memory Parallel Computer



- * each processor has its own local memory for doing computations (not shared)

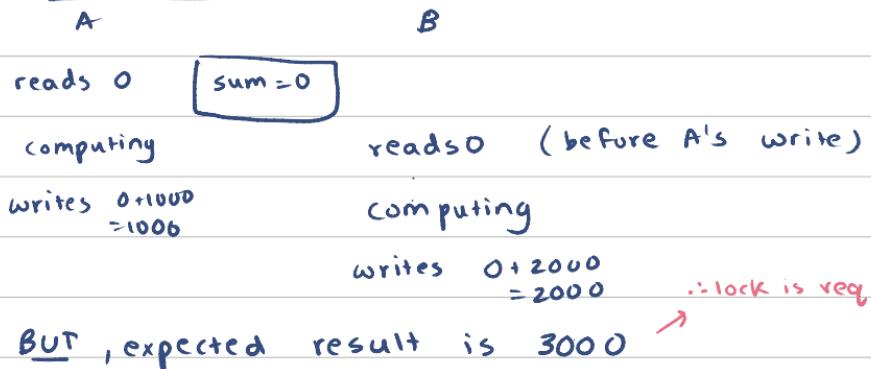


- * fork \rightarrow create new parallel process

- * join \rightarrow synchronize with another parallel process

Date: /

without lock

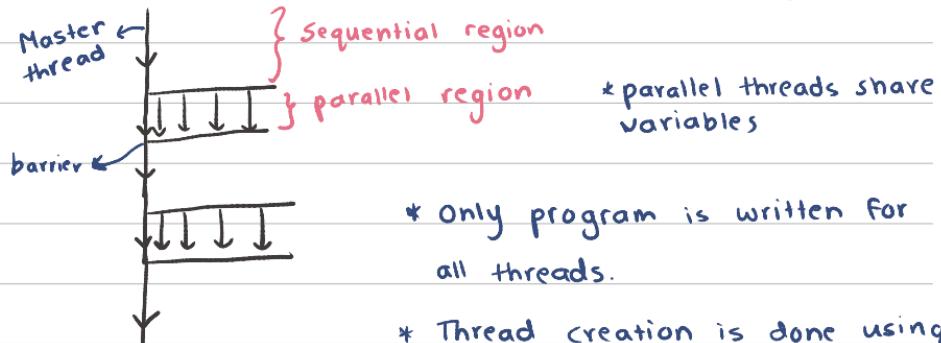


* Open MP - Standard → APIs for parallel machines with shared memory

↓
Open specification for Multi-Processing

- ① Compiler directives
- ② Run-time library function
- ③ Environment variables

* Open MP talks in terms of threads, not processes.



* Only program is written for all threads.

* Thread creation is done using an OpenMP function call.

Date: /

* thread id = processor rank

* Code using OpenMP

```
#include <stdio.h>
```

```
#include <omp.h>
```

```
int main ()
```

```
{
```

shared variable [int noOfThreads, threadID, THREAD-COUNT=4;]

fork \equiv #pragma omp parallel private (threadID)

num-threads (THREAD-COUNT)

* converting threadID to a ptv variable
(specific to each thread)

run by all of the threads [threadID = omp_get_thread_num();
printf ("I am thread %d\n", threadID);]

destroys all threads [#pragma omp barrier
if (threadID == 0)
{ }

noOfThreads = omp_get_num_threads()

printf ("Master: My team has %d threads\n", noOfThreads);

}

return 0;

3

Date: /

* Output sample

I am thread 0
I am thread 1
I am thread 2
I am thread 3

can be printed
in any order
• 4 threads are competing
for std::io

Master: My team has 4 threads

* Master thread is included in THREAD-COUNT

12/9/23

omp-get-thread-num() → threadID

omp-get-num-threads() → no. of threads

#pragma omp parallel → creates 11el threads
compiler directive

#pragma omp barrier → synchronization of threads

* Parallel Loop (code)

#include <stdio.h>

* #include <omp.h>

#define N 7

int main (void)

{

Date: /

parallel:
for loop int i, THREAD-COUNT = 3;

in
OpenMP \triangleright #pragma omp parallel for num-thread (THREAD-COUNT)
{
 for (i=0; i<N; i++) \rightarrow execute the iterations
}
of this for loop in
parallel

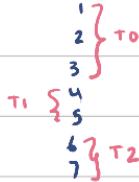
return 0;

}

Note:

- * you don't have to explicitly create new threads when you use omp parallel for
- * the for loop is unrolled & different iterations are assigned to THREAD-COUNT # threads.

Eg: Iteration



- * the #iterations MUST BE KNOWN at compile time, so that it can be split across threads by the compiler.

Eg: for (i=2; i<7; i++)
{
 a[i] = a[i] + c;
 b[i] = a[i-1] + b[i];
}

\Rightarrow iterations are not independent.

\therefore can't use parallel for

Date: /

- * Can't use parallel for if there are dependencies across iterations.

dependency within iteration ✓

across iterations ✗

- * you can provide additional information to

the compiler via clauses ↴ w/o specification = 1

Eg: i) Schedule (static, CHUNK-SIZE)

{ Suppose chunk-size = 8 , iterations 1 to 8
go to thread 1 , 9 to 16 → thread 2 and
so on (round-robin till all iterations are
allotted to a thread)}

used when
iters take
the same
time

ii) Schedule (Dynamic, CHUNK-SIZE)

- if chunk-size = 8
 - 1) 8 iters are allotted to each thread
 - 2) if thread i completes its work first,
then next 8 are given to i
(NOT ROUND ROBIN)

* different iterations take different times

iii) Schedule (Guided, CHUNK-SIZE)

- Split the first $\frac{n}{2}$ iterations across all
- Then the next $\frac{n}{4}, \frac{n}{8}, \frac{n}{16}$

Date: /

- * for scenarios where the workload decreases over the iterations

(iv) Schedule (auto, chunk-size)

- * when you don't know anything about the iterations.

* Data Parallelism - parallel for

* Functional

#pragma omp sections
{

* non-iterative

* creates threads

#pragma omp section
{ }

* each section is allocated to a thread & run in parallel

#pragma omp section
{ }

3

Eg: Sequential code

int A[max], i, sum = 0;

double average;

for (i=0 ; i < max ; i++)

 sum = sum + A[i]

average = sum / max

Date: /

Parallel code

```
int A[max], i, sum = 0;  
double average;
```

* #pragma omp parallel for reduction (+:sum)

{

```
for (i=0 ; i < MAX ; i++)
```

```
    sum = sum + A[i];
```

}

```
average = sum / max
```

= aggregation

reduction (+:sum)

at the end, the
partial sums
computed by each
thread is accumulated
at the master
thread.

- better alternative to locking & unlocking a shared variable "sum"
- here, we're just leaving it to the compiler to sum across all "sum"s and give it to the master thread.

Date: 15 / 9 / 23

- * Synchronization:
 - 1) `#pragma omp critical`
`{ //block }`
 - 2) `#pragma omp atomic`
`//one statement`
 - 3) `#pragma omp barrier`
 - 4) `#pragma omp parallel for nowait`
 - nowait \Rightarrow no barrier is added at the end
 \Rightarrow if we need a barrier, add it explicitly
 - by default parallel for has a synchronization barrier at the end
 - 5) `#pragma omp flush [list]`

`omp_set_lock()`

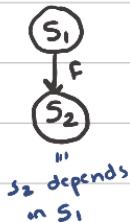
`omp_unset_lock()`

* Dependence Analysis

(a) Data dependence

$$S_1: a = b \times 2$$

$$S_2: d = a + c$$



Flow dependence

RAW - Read After Write

Date:

$$S_3 : e = g \times 2$$

$$S_4 : g = f + h$$



Anti-dependence

WAR - Write after Read

$$S_5 : e = g \times 2$$

$$S_6 : c = f + g$$



Output dependence

WAW - Write after write

(b) Control dependence

$$S_7 : \text{if } (a == 5) \text{ then}$$

$$S_8 : b = 9$$



Eg: for i = 2 to 7 do

$$S_1 : a[i] = a[i] + c$$

$$S_2 : b[i] = a[i-1] + b[i]$$

endfor

* dependency across iterations

* compiler does thorough analysis
+ optimization to execute this loop in parallel.