

## Parallel Software (MK book, §2.4)

- Issues of writing software for parallel systems

ex) { typical shared-memory program:  
MIMD { start a process → forks threads  
typical distributed-memory program:  
start multiple processes

focus is on SPMD – single program/multiple data ↔

Running a different program on each core (MPMD)

ex) { if (process\_id/thread\_id == 0)  
do this;  
Task parallel { else  
do that;  
  
{ if (process\_id/thread\_id == 0)  
operate on the 1<sup>st</sup> half of the array;  
Data parallel { else  
operate on the 2<sup>nd</sup> half of the array;

SPMD model is suitable for implementing both data and task parallelism.

Task parallel program – dividing tasks among processes/threads.

### Coordinating processes/threads:

1. Divide work among processes/threads  
Load balancing, minimum communication
2. Synchronization among processes/threads
3. Communication among processes/threads

- **Shared-memory system issues**

- vs. {
- Dynamic threads – master thread forks a thread (joins to master thread) when needed;  
→ high cost of thread fork/join (time-consuming operations)
  - Static threads – all are forked at the same time;  
→ less efficient resource usage (idle thread holds resources)

Non-determinism: vague order, ex) which thread terminate earlier/later

ex) {  
`my_val = compute_val (my_rank);`  
`x += my_val; // x is shared_var, my_val is local to each thread`

one exec scenario:

	<u>thread 0</u>	<u>thread 1</u>
time ↓	x = 0	...
	my_val = 7	x = 0
	x = x + 7	my_val = 9
	store x = 7	x = x + 9
	...	store x = 9

finally, shared x = 9, which is different from the desired result (16).

→ problem: simultaneous updates of shared location.

critical section – shared location is updated

Critical section (C.S.) should be mutually exclusive;

How to implement?    - using a mutex semaphore  
                                  - monitor  
                                  - busy-waiting (highly inefficient)

Lock C.S. Unlock
------------------------

C.S. is serialized among processes/threads

Thread safty – some built-in serial func.s malfunction in parallel program

- **Distributed-memory system issues**

multiple processes, each on an independent processor

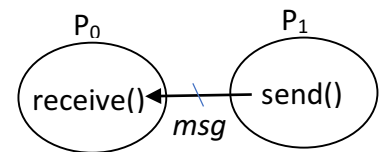
message-passing API (ex, MPI)

send/receive functions

```
ex) char msg[100]; //local to each process
    my_rank = get_rank( );
    if (my_rank == 1)
    { msg = "... from process_1";
      send (msg, MSG_CHAR, 100, 0);
    }
    else if (my_rank == 0)
    { receive (msg, MSG_CHAR, 100, 1);
      cout<<...;
    }
```

*SPMD:  
2 proc's use  
same code,  
but different  
actions*

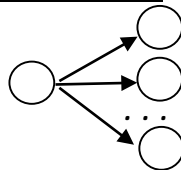
ref)  
It's possible to use  
distributed-mem API on  
shared-mem system;  
Programmer logically  
partitions shared-mem into  
private address spaces;



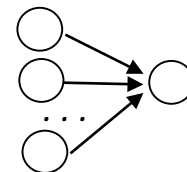
functions: send( ) – 1. Sending proc is blocked until receiver is ready;  
or, 2. sending proc sends data to a buffer and continue op;  
receive( ) – receiving process is blocked until msg is received;

Collective communication

broadcast

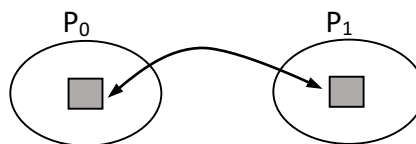


reduction



One-sided communication (ex, remote mem access in DSM)

not easy to use



P<sub>0</sub> updates P<sub>1</sub>'s mem with P<sub>0</sub>'s data  
or, P<sub>0</sub> updates P<sub>0</sub>'s mem with P<sub>1</sub>'s data

Programming hybrid systems – e.g., cluster of multicore processors

It's possible to use both MPI and OpenMP, but difficult to implement;

→ using only MPI for both inter- and intra-node communication.

### **Input and Output** (MK book, §2.5)

C/C++ with OpenMP, MPI, Pthread

Multiple processes/threads → non-determinism on I/O (stdin, stdout, stderr)

ex)  $\bigcirc P_0 \bigcirc P_1 \bigcirc P_2 \dots \bigcirc P_n$

only 1 process should be responsible for output

### Assumptions/suggestions/rules:

- in distributed-mem programs – only process\_0 will access stdin;
- in shared-mem programs – only the master thread (Thread\_0) will access stdin;
- in both distributed-/shared-mem programs,
  - all processes/threads can access stdout/stderr;
  - non-deterministic
  - Sol: only 1 process/thread should be used for stdout/stderr
- only a single process/thread attempts to access any single file;
  - each process/thread can open its own private file for r/w, but no two processes/threads will open the same file
- debug output should include the rank (id) of process/thread;

## Performance (MK book, §2.6)

*ideal*

$$T_{\text{parallel}} = \frac{T_{\text{serial}}}{P} \quad // \text{linear speedup (ideal), where } P \text{ is number of processors}$$

overheads: increases as # of processes/threads increases

shared-mem – critical sections are executed in serial (with mutex);

*ideal* distr-mem – data transition across IN (slower than local mem access);

$$\text{ideal (linear) speedup } S = \frac{T_{\text{serial}}}{T_{\text{parallel}}} = \frac{T_{\text{serial}}}{\frac{T_{\text{serial}}}{P}} = P$$

in reality,

as P increases, S becomes smaller;

as P increases,  $\frac{S}{P}$  (efficiency of parallel program) becomes smaller;

$$\text{efficiency (speedup per processor) } E = \frac{S}{P} = \frac{\frac{T_{\text{serial}}}{T_{\text{parallel}}}}{P} = \frac{T_{\text{serial}}}{P * T_{\text{parallel}}} = \frac{1}{P} \left( \frac{T_s}{T_p} \right)$$

*effective*

$$\text{effective } T_{\text{parallel}} = \left( \frac{T_{\text{serial}}}{P} \right) + T_{\text{overhead}}$$

$$\text{effective speedup } S = \frac{T_{\text{serial}}}{\text{effective } T_{\text{parallel}}} = p * \underbrace{\left( \frac{T_s}{T_s + P * T_o} \right)}_{< 1}$$

ex) Table 2.4

P	1	2	4	8	16	cores/processors
S	1.0	1.9	3.6	6.5	10.8	
E	1.0	0.95	0.9	0.81	0.68	

UKbook

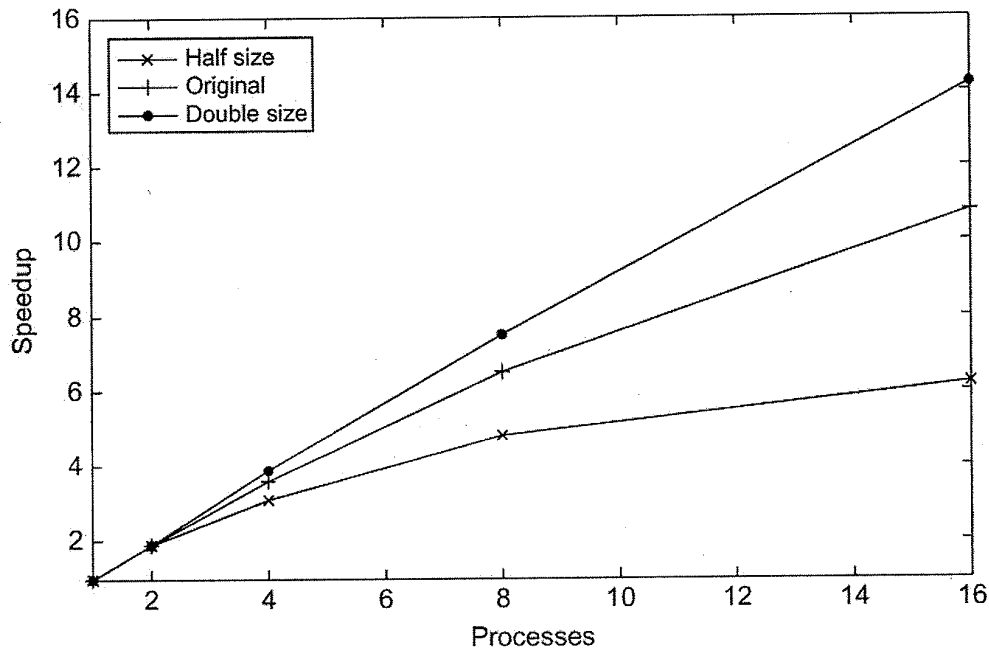


Fig 2.18 Speedups of parallel program on different problem sizes

ideal  $S = P$   
 effective  $S$   
 $= P \cdot \left( \frac{T_s}{T_s + P \cdot T_o} \right)$   
 less than 1  
 $\Rightarrow S < P$

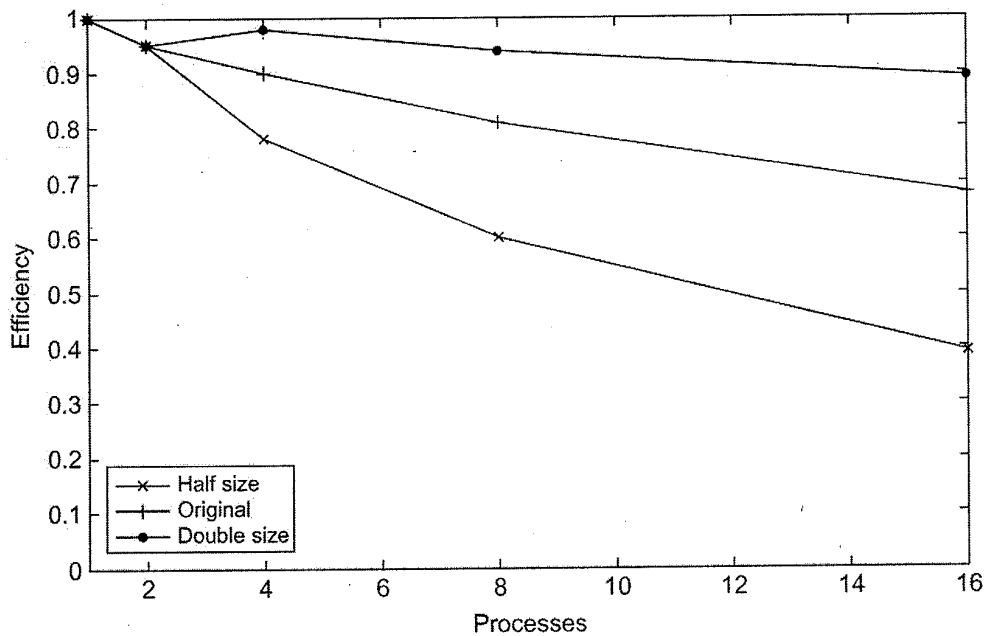
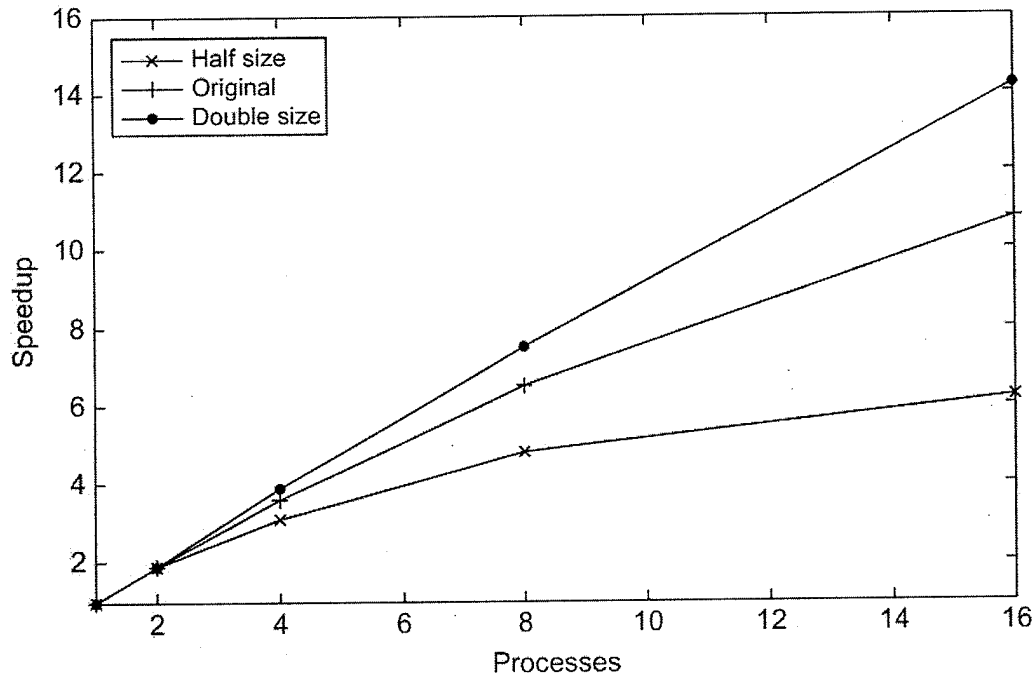
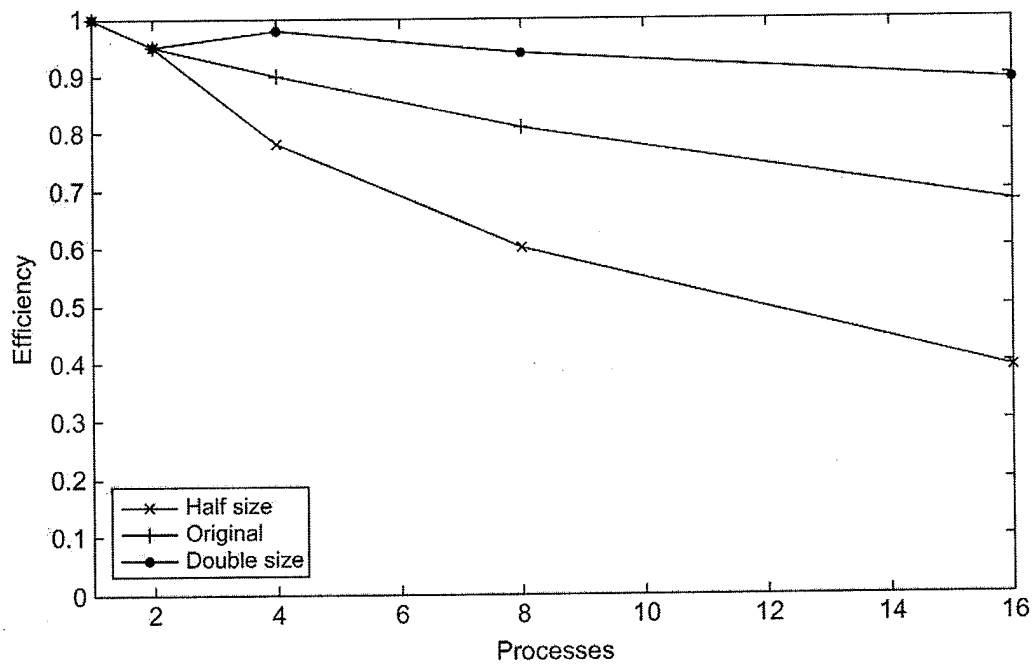


Fig 2.19 Efficiencies of parallel program on different problem sizes

$E = \frac{S}{P}$  (ideal case,  $E = 1$ )  
 $\because S = P$



**Fig 2.18** Speedups of parallel program on different problem sizes



**Fig 2.19** Efficiencies of parallel program on different problem sizes

## Time checking (measuring parallel program run time)

Elapsed time = CPU time + I/O time      *//>time for elapsed time*

Run time of parallel program – varies at different run

wall-clock time from start to finish of any interested part

ex) double start, finish;

```
....  
start = Get_current_time( );  
....  
....  
finish = Get_current_time( );  
cout<<finish – start;
```

get wall-clock time  
ex) MPI:  
    MPI\_Wtime( );  
OpenMP:  
    Omp\_get\_wtime( );

Global time checking for multiple processes/threads:

ex) shared double global\_elapsed;

private double my\_start, my\_finish, my\_elapsed;

barrier( ); *//synchronize all processes/threads*

```
my_start = Get_current_time( );
```

```
....
```

```
....
```

```
my_finish = Get_current_time( );
```

my\_elapsed = my\_finish – my\_start; *//for each local process/thread*

*//find the maximum across all processes/threads*

global\_elapsed = Global\_max(my\_elapsed); *//updates shared location*

if (my\_rank == 0) *//if my\_elapsed > current\_max*

cout<<global\_elapsed;

Reported run time should exclude I/O time.



## Maximum speedup from parallelization

### Amdahl's law

$$\text{Speedup} = \frac{1}{(1 - Fe) + (\frac{Fe}{Se})}$$

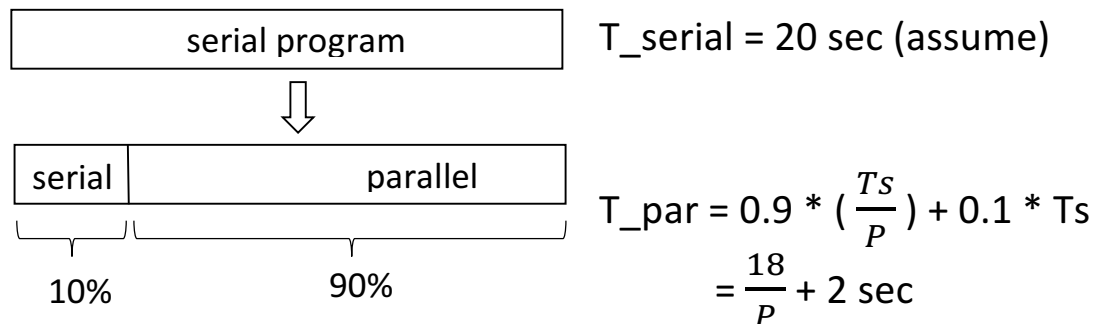
ex) 10% of a task is serial, 90% of the task is parallelizable;

using P processors, ideal (linear) speedup = P

$$\rightarrow \text{maximum speedup } S = \frac{1}{(1 - 0.9) + (\frac{0.9}{P})} = \frac{1}{0.1 + (\frac{0.9}{P})} = \frac{1}{0.1} = 10$$

when P gets larger  $\Rightarrow 0$

analysis:



$$\rightarrow S = \frac{T_s}{T_p} = \frac{20}{\frac{18}{P} + 2}$$

as P gets larger,  $\frac{18}{P} \Rightarrow 0$

$$\text{so, } S \leq \frac{20}{2} = 10$$

## Scalability

Assume that  $n$  (problem size) and  $P$  (# of processes/threads) yields efficiency  $E$ .

If we increase  $\left[ \begin{array}{l} n \rightarrow k * n \\ P \rightarrow k * P \end{array} \right.$

and this yields same efficiency  $E$ , then the parallel program is scalable;

more specifically, the parallel program is weakly scalable.

strongly scalable:

only  $p \rightarrow k * P$  (regardless of the problem size  $n$ ) yields same efficiency  $E$

weakly scalable:

both  $p \rightarrow k * P$  and  $n \rightarrow k * n$  yields same efficiency  $E$   
(same increasing factor  $k$  for  $P$  and  $n$ )

ex) problem size =  $n$

$$T_s = n$$

$$T_p = \left(\frac{n}{P}\right) + 1$$

$$\rightarrow E = \frac{S}{P} = \frac{\frac{T_s}{T_p}}{P} = \frac{T_s}{P * T_p} = \frac{n}{P \left(\frac{n}{P} + 1\right)} = \frac{n}{n + P}$$

ref)

$$\text{Ideal } T_p = \frac{T_s}{P}$$

$$\text{Ideal } S = \frac{T_s}{T_p} = P$$

$$\text{Ideal } E = \frac{S}{P} = 1$$

if we increase  $P \rightarrow k * P$ ,

$$E' = \frac{n}{n + K * P} \neq E \quad \text{--- program is not strongly scalable}$$

if we increase  $P \rightarrow K * P$  and  $n \rightarrow k * n$ ,

$$E' = \frac{k * n}{k * n + k * P} = \frac{n}{n + P} = E \quad \text{--- program is weakly scalable}$$

## Parallel program design (MK book, §2.7)

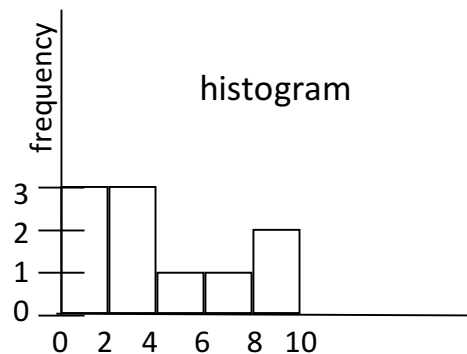
Design steps:

1. partitioning – divide computation into small tasks;
2. communication – among the small tasks;
3. aggregation – combine small tasks into a bigger task;  
→ reduces communication cost
4. mapping – assign tasks to processes/threads with minimum communication and load-balancing way.

ex) computation for making a histogram

ex) 1, 3, 4, 2, 3, 5, 7, 10, 2, 9

bin#	range	count
[0]	1~2	3
[1]	3~4	3
[2]	5~6	1
[3]	7~8	1
[4]	9~10	2

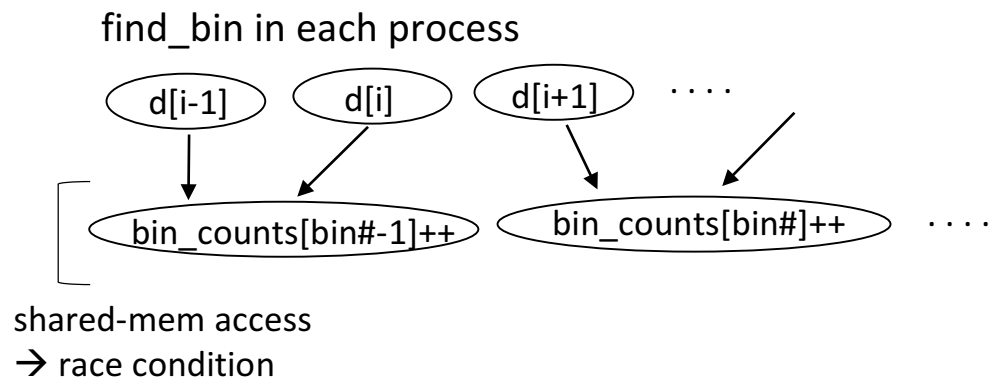


serial algorithm:

```
for each input data element  $i$ ,  
{ determine the bin (it is belong to); // find_bin(  $i$  ) → bin# (index)  
  increment the count of that bin;    // bin_counts[bin#]++;  
}
```

## parallel schemes

### 1. Shared bin-count way:



Processes/threads update shared-memory location (`bin_counts[ ]`)

→ race condition: the shared location is C.S. and should be protected by mutex.

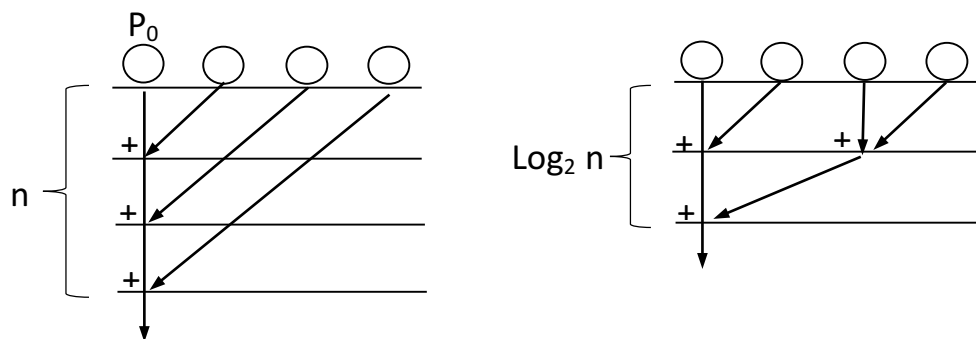
### 2. Local-sum way: -- data parallel on distributed-mem system

Input data is divided by the number of processes/threads;

each process implements the serial algorithm on the part of the data;

then, local `bin_counts` are combined into the global one.

- (1) by a process/thread  
 (2) by tree-reduction way
- both needs communication (send/receive data)*



## Writing and running parallel programs

- Small shared-memory system:
  - 1 op sys schedules threads on the cores/processors;
- large shared-memory system:
  - batch scheduler – user requests # of cores, etc.
- distributed-memory system (e.g., clustered system):
  - host computer allocates nodes among users;
  - or, batch (job queue), check-out way, etc.

## Assumptions used in this course:

Homogeneous MIMD system

SPMD programming

At most one process/thread on a processor/core