

$$\text{efficiency} = \frac{T_s}{T_p} \times \frac{1}{P} \quad (\text{how many processors involved in whole process.})$$

(amount of processor utilization.)

- Bounds of efficiency are the upper and lower bound.

$$0 \leq \varepsilon(n, p) \leq 1.$$

$$\begin{aligned} \varepsilon(n, p) &\leq \left[\frac{\sigma(n) + \phi(n)}{\sigma(n) + \phi(n) + k(n, p)} \right] \times \frac{1}{P} \\ &\leq \boxed{\frac{\sigma(n) + \phi(n)}{P\sigma(n) + \phi(n) + pk(n, p)}} \end{aligned}$$

we want an ideal situation where $\sigma(n) = 0$ and $k(n, p) = 0$.

$$\varepsilon(n, p) = 1 \quad (\text{upper bound}), \text{ best case}$$

worst case: everything in parallel.

$$p \rightarrow \infty, \phi(n) = 0.$$

$$\varepsilon(n, p) = 0 \quad (\text{lower bound}), \text{ worst case.}$$

$$\psi(n, p) \leq \frac{\sigma(n) + \phi(n)}{\sigma(n) + \phi(n) + k(n, p)}$$

$$\psi(n, p) \ll \frac{\sigma(n) + \phi(n)}{P\sigma(n) + \phi(n)}$$

Set us define a function of form:

$$f = \frac{\sigma(n)}{\sigma(n) + \phi(n)} \quad (\text{amount of serial part in the code.})$$

$$u(n, p) \ll \frac{\sigma(n) + \phi(n)}{\sigma(n) + \phi(n)}$$

$$\text{if } f\sigma(n) - \sigma(n) = -f\phi(n) - f\phi(n)$$

$$\Rightarrow \phi(n) = \sigma(n) \left[\frac{1-f}{f} \right]$$

$$\Rightarrow \phi(n) = \sigma(n) \left[\frac{1}{f} - 1 \right] \quad \text{for } f = \text{number of processors}$$

$$\psi(n,p) = \frac{\sigma(n)/f}{\sigma(n) + \sigma(n)(\frac{1}{f} - 1)}$$

$$= \frac{1/p}{1 + (\frac{1}{f} - 1)/p} = \frac{T}{f + f - 1} = \frac{T}{p}$$

andahl's law will give an upper bound on speedup for a number of processors and problem size.



For a fixed number of processors, the speedup will increase as no. of processors will increase, this is Andahl's effect (different from Andahl's law.)

General formula of a code on p processors, taking $\sigma(n)$ as time

$$\psi(n,p) \leq \frac{\sigma(n) + \phi(n)}{\sigma(n) + \phi(n) + k(n,p)}$$

$$\frac{\sigma(n) + \phi(n)}{p} + \frac{(n-1)}{p} \geq \frac{n-1}{p}$$

how good is the code on p processors compared to only one processor.

$$(2-17)(n) = (n) \cdot 6$$

given that the time is divided into:
 s → fraction of time spent in the parallel work (performing serial work)

$$s = \frac{\sigma(n)}{\sigma(n) + \frac{\phi(n)}{P}} \quad | \quad 1-s = \frac{\phi(n)/P}{\sigma(n) + \phi(n)/P}$$

$$\sigma(n) = \left[\sigma(n) + \frac{\phi(n)}{P} \right] s$$

$$\phi(n) = \left[\sigma(n) + \frac{\phi(n)}{P} \right] (1-s) P$$

$$\begin{aligned} \psi(n, p) &\leq \frac{\left[\sigma(n) + \frac{\phi(n)}{P} \right] \phi(s) + \left(1 - \frac{\phi(n)}{P} \right) \phi(1-s)}{\sigma(n) + \frac{\phi(n)}{P}} \\ &\leq p + (1-p) s \end{aligned}$$

This will show how the code will perform if it is put on p processors.

Sources of overhead:

- 1) creation of threads (process setup).
- 2) communication time.
- 3) extra computation.
- 4) synchronization (race conditions, etc.).

Experimentally determined serial fraction:

As number of processors increases towards infinity

$$\tau(n, p) = \sigma(n) + \frac{\phi(n)}{P} + k(n, p)$$

where $k(n, p)$ corresponds to overhead

$$\tau(n, 1) = \sigma(n) + k\phi(n) \rightarrow \text{serial work}$$

$$e = \frac{(p-1)\sigma(n)}{(p-1)\tau(n, 1)} + \frac{pk(n, p)}{(p-1)\tau(n, 1)}$$

$$\Rightarrow e = \frac{p\sigma(n)}{p\tau(n, p)} - \frac{\sigma(n)}{\tau(n, 1)} = p\tau(n, p) - \tau(n, 1)$$

$$\Rightarrow = p\tau(n, p) - \phi(n) - \sigma(n) = p\tau(n, p) - \tau(n, 1)$$

$$pT(n,p) = e(p-1) T(n,1) + T(n,1)$$

$$pT(n,p) = T(n,1) ep - T(n,1)e - eT(n,1)$$

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

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

$$\Rightarrow 1 = \psi_e + \frac{\psi(1-e)}{p}$$

$$\Rightarrow \frac{1}{\psi} = e + \frac{1-e}{p}$$

$$\Rightarrow \frac{1}{\psi} = e + \frac{1}{p}$$

$$\Rightarrow e = \frac{\psi - 1/p}{1 - 1/p}$$

(This is for experimental observations.)

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

p	2	3	4	5	...	8
ψ	1.82	2.5	3.08	3.57	...	4.71
e	0.199	0.1	0.1	0.1	...	0.1

- we are getting a speedup which is not linear.
- Problem: Most of the processors are idle, or there is some overhead.

$$e = \frac{(p-1)\ln p}{(p-1)T(n,1)} + p k(n,p)$$

If $p k(n,p)$ is removed,

$$e = \text{constant}$$

As the value of e is constant, it is no longer possible to further improve the code.

p	2	3	4	5	6	8
ψ	1.87	2.61	3.2	3.73	4.2	4.7
e	0.09	0.75	0.88	0.81	0.64	0.51

In a fixed problem size, efficiency will decrease with increase in number of processors.

* Inefficiency metric :-

$$\psi(n, p) \leq \frac{\sigma(n) + \phi(n)}{\sigma(n) + \phi(n) + (p-1)\alpha(n) + pk(n, p)}$$

It is the ability to increase performance as the number of processors increases.

$$\Rightarrow \psi(n, p) \leq \frac{p[\sigma(n) + \phi(n)]}{\sigma(n) + \phi(n) + (p-1)\alpha(n) + pk(n, p)}$$

$$T_0(n, p) = (p-1)\alpha(n) + pk(n, p) : \text{overhead.}$$

Efficiency :

$$\epsilon(n, p) \leq \frac{\sigma(n) + \phi(n)}{\sigma(n) + \phi(n) + T_0(n, p)}$$

$$\leq \frac{8.05}{1 + T_0(n, p)/T_0(n, 1)} = \frac{0.05}{0.05 + k - 1} = \frac{0.05}{0.05 + 0.01} = 0.98$$

$$\frac{T_0(n, p)}{T_0(n, 1)} \leq \frac{1 - \epsilon(n, p)}{\epsilon(n, p)}$$

$$\boxed{T(n, 1) = \frac{\epsilon(n, p)}{1 - \epsilon(n, p)} T_0(n, p)}$$

Rate at which problem size has to be increased to maintain constant efficiency.

99% }
 97% } Cache hit.
 90% }

Cache hit : getting data from the cache

- If data is not in cache, penalty must be paid to get the data from the RAM.

Cache hit time = 1 cycle.

Miss penalty \approx 100 cycle.

If you not get something at your nearest location, you have to pay a penalty.

$$\text{Average access time} = \frac{99}{100} \times 1 + \frac{1}{100} \times 100$$

$$= 1 + 0.99 = 1.99$$

$$= 1.99 \text{ cycle}$$

If cache hit rate is 97%

$$\text{Average access time} = \frac{97}{100} \times 1 + \frac{3}{100} \times 100$$

$$= 1 + 0.97 = 1.97$$

The access time will reduce by a factor of around 2.

$$\frac{80}{100} + \frac{20}{100} \times 100 = 20.8$$

Miss rate = 1 - cache hit

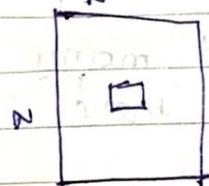
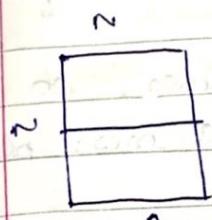
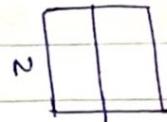
Cache hit time : Time taken for data to be transferred from cache to processor.

Processor transfer time

L1 - 30, L2 = 20, B = N, n = 10, cache size = 30

L3 - 256.

L4 - 6-8 MB.



Matrix elements = 8 bytes.

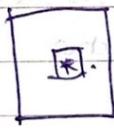
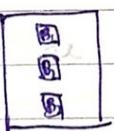
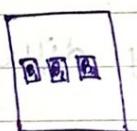
Total no. of operations in square = $O(N^2)$.

Cache block size: A chunk of data is received instead of a single block of data, generally of 64 or 32 bytes.

No. of cache misses = $\frac{N}{8} + n$

Cache size < Problem size (N).

- In the first 3 access, there will be no miss but not for the others.
- In C, the data is stored in row major, so there will be cache misses for any only one element in row but for all elements in the case of column.



(Multiplying 2 blocks, we get

$B_1 \times B_2 \rightarrow B_1 \text{ blocks} \times B_2 \text{ blocks} \rightarrow \text{resultant block.}$

Block size = $R^2 = B \times B$ (size of block)

- Assume that the cache can hold at most 3 blocks but not more than 3 blocks.

- Condition for $3R^2 \leq \text{cache size}$

$$\text{cache miss for each block} = \left(\frac{n}{B} + \frac{n}{B}\right) \frac{B^2}{8} = \frac{2n}{B} \times \frac{B^2}{8}$$

$$= \frac{nB}{4}$$

- one block is reused many times, so only for the first time, there is a miss penalty.

$$\text{Total no. of blocks} = \left(\frac{n}{B}\right)^2, O\left(\frac{n^3}{B}\right)$$

$$\text{Thus, } \left(\frac{n}{B}\right)^2 \times \frac{nB}{4} = \frac{n^3}{4B}$$

$\frac{9n^3}{8}$ $\frac{n^3}{4B}$ (as n is fixed).
value of B will

$$B=4 : \frac{n^3}{16} \text{ change.}$$

By doing block matrix multiplication, the increase in speed = $\frac{9n^3}{8} \times \frac{16^2}{n^3} = 18$.

Thus, algorithm will be 18 times faster.

Now, if $B=8$

$$\frac{9n^3}{8} \times \frac{32}{n^3} = 36 \text{ times faster}$$

This won't always be true if done on a machine practically.

- Thus cache change is very different.

General lesson: Keep working data set very small, this can be achieved by dividing into blocks so that it is possible to fit the data into cache and so the cache hit rate increases.

- Keep the strike sizes small.
- Focus on pattern in innermost loop.
- 1) n processors takes $(\log n)^2$ time in parallel.
serial time = $n \log n$. special case ($p=n$).
 $\text{speedup} = \frac{n}{\log n}$
 $\text{efficiency} = \frac{1}{\log n}$
 $\text{cost} = n(\log n)^2$.
- efficiency will decrease with increase in the problem size.
- the algorithm is not cost optimal by a factor of $\log n$.
- let us make a real assumption where $p < n$.

$$T_p = \frac{(\log n)^2}{p} \cdot n(\log n)^2$$

$$\text{new speedup} = \frac{p}{\log n}$$

the speedup decreases as the problem size increases for a fixed value of the problem size because the algorithm is not cost optimal.

* sorting a list \rightarrow serial $\rightarrow n \log n$.

Algorithm	A_1	A_2	A_3	$= A_4$
No. of processes	n^2	$\log n$	n	\sqrt{n}
T_p	1	n	\sqrt{n}	$n \log n$

$$203.61 = (0.01 \times 50.2) + (0.01 \times 39.9) + (5 \times 8)$$

$$203.61$$

	$n \log n$	$\log n$	$n \log n$	\sqrt{n}	
efficiency	$\frac{\log n}{n}$	1	$\frac{\log n}{n}$	1	
cost	n^2	$n \log n$	$n \log n$	$n^{1/2}$	$n \log n$

(a)

A₄ should be chosen.

A₁ won't be chosen because the cost is very high.

A₃ won't be chosen as speedup is high. We may choose between A₂ and A₄.

(b)

+ Parallel program on 2 process system.

- Problem size = w.

- serial \rightarrow 1 processor, cache = 64 kB, cache hit rate = 80%

- latency of cache = 2 ns, latency of DRAM = 100 ns.

Memory bound computation \rightarrow 1 flop/memory access

Parallel \rightarrow each processor does $w/2$.

- cache hit rate = 90%, 8% from RAM, 2% from the remote DRAM (communication overhead)
 \rightarrow latency of 400 ns.

Comment on expected speedup.

A: expected speedup = 2.42

The speedup is a superlinear speedup.

This is due to a combined cache effect

Serial:

$$(0.8 \times 2) + (0.2 \times 100) = 21.6 \text{ ns.}$$

$$\text{Performance} = \frac{1}{21.6 \text{ ns}} = 46.3 \text{ Mflops (serial)}.$$

(effective memory access time)

Parallel performance:

$$(0.9 \times 2) + (0.08 \times 100) + (0.02 \times 400) = 17.8 \text{ ns.}$$

$$\text{Performance} = \frac{1}{17.8 \text{ ns}} \text{ (parallel)}$$

$$\text{Processing time} = 2 \times 17.8 = 35.6 \text{ ns}$$

$$\text{Performance per processor} = \frac{1}{17.8} = 56.18 \text{ MHz}$$

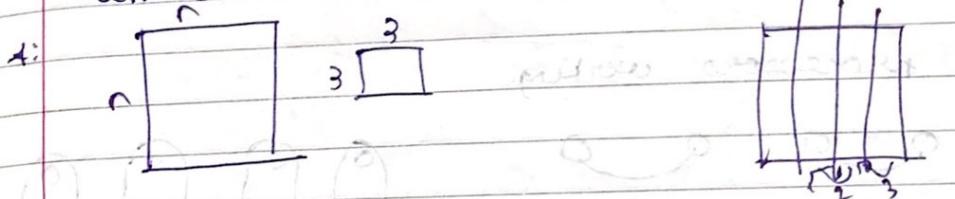
$$\text{Processor Performance} = 2 \times 56.18$$

$$\text{Speedup} = \frac{56.18 \times 2}{46.3} = 2.42$$

* Image processing (convolution).

- $n \times n$ pixel image
- convolute with 3×3 mask.
- each multiply-add \Rightarrow takes (t_c) time
- distributed system \Rightarrow each node has its subimage and applies the mask on subimage.
- each pixel \Rightarrow one word communication
- communication latency $\rightarrow t_s + t_w$.
- $L \downarrow$

comment on efficiency?



Efficiency in terms of n, p, t_s, t_w, t_c

$$\text{Serial time} = n^2 q t_c$$

$$\text{Parallel time} = \frac{q n^2 t_c}{p} + \lceil 2(t_s + t_w n) \rceil$$

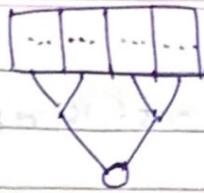
we always assume that $p=n$.

$p < n$: change granularity of problem.

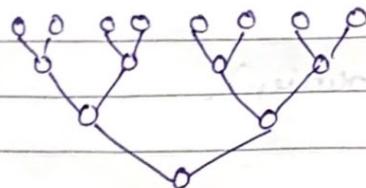
effect of granularity on performance.

n numbers on $n-p$ processors.

$n=16, p=4$ (power of 2). valid only for power of 2.



Design 1: $\frac{n}{p} + \log p$



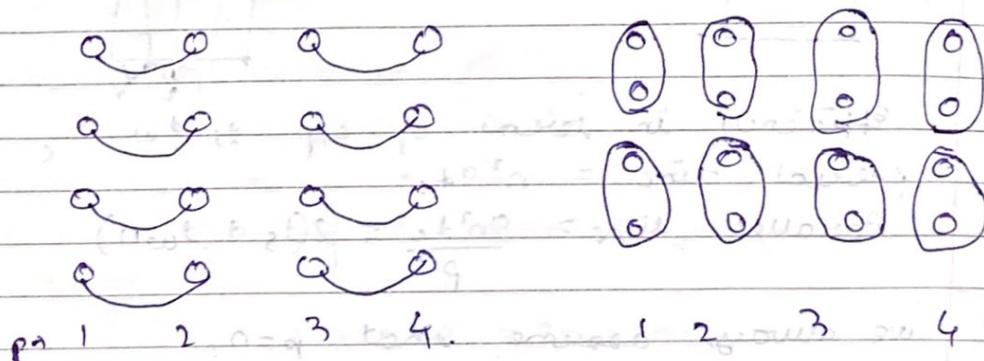
n elements = p processors. $\lceil \log n \rceil$

n elements $\rightarrow p$ processors with the same design.

cost $\rightarrow O(n + p \log p) \sim O(n)$. this is cost time as long as $n \sim p \log p$.

- If $n = p$ processors, $\log p$ steps will be done to converge to 1 processor.

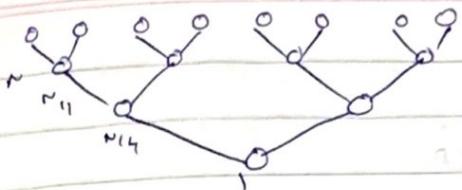
4 processors working



cost $\in O(n \log p)$ (with p processors)

this algorithm is non-cost optimal.

$$n=16, p=4. \frac{16}{4} \log_2 4 = 4+2=8$$



simulated id on p
virtual processors.

$$p \ll n$$

After each stage, there is a deactivation.

$$n = 2^{k+1} \rightarrow n/p$$

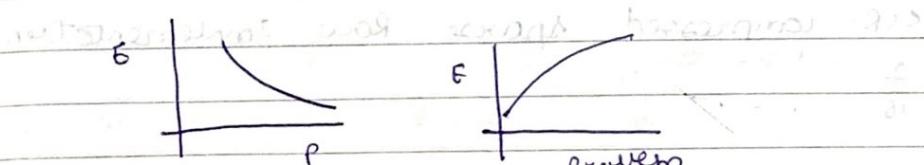
- 1) replaces $\log p$ steps after the $\log n$ steps at the original. (step $\mapsto n/p$)
- 2) step 2 $\mapsto n_{12}/p_{12} = n/p$.
- 3) step 3 $\mapsto n/p$.

$$\boxed{\frac{n}{p} \log p + \frac{n}{p}}$$

cost = $p \left[\frac{n}{p} \log p + \frac{n}{p} \right] = n \log p + n$: this is non-cost optimal.

$$\frac{\partial E}{\partial p} = \frac{S}{p} = \frac{T_S}{p T_p}$$

$$\Rightarrow E = \frac{1}{1 + \frac{T_S}{T_p}} \quad | T_0 = p T_p - T_S$$



$$N = \frac{E}{1-E} T_0(n, p) \rightarrow \text{Problem size, inefficiency function.}$$

Problem size: Minimum number of computers in serial case.

loading n numbers on p processors.

$$w = k \log p \quad , \quad P = \frac{p \log p}{p \log p}$$

$$T_0 = p^{3/2} + p^{3/4} w^{3/4}$$

w = problem size.

p = no. of processors.

a) what rate should problem size be increased to maintain same efficiency?

b) step 1:- treat the problem individually.

$$w = kp^{3/2}$$

$$w = kp^{3/4} w^{3/4}$$

$$w = kp^{3/4}$$

$$w = k^4 p^3$$

algorithm

Inefficiency

relation

Architecture

c) sparse matrix

vector

$$\begin{array}{cccc|c|c|c|c} 3 & 0 & 1 & 0 & | & * & | & ? \\ 0 & 0 & 0 & 0 & | & * & | & A \times B \\ 0 & 2 & 4 & 1 & | & * & | & \\ 1 & 0 & 0 & 1 & | & * & | & \\ \end{array} \quad \boxed{\text{CSR}}$$

CSR - compressed sparse row implementation.

$\frac{7}{16}$

Data: 3 1 2 4 1 1 1 : 7 elements.

Col pointers: 0 2 1 2 3 0 3 :

Rowids: 0 2 2 5 7 : look into the serial matrix

No. of rows = 5 - 1 = 4.

$$n + n + n = O(n)$$

n

$$n \times n = O(n^2)$$

$$10000 \times 10000 = 100000000$$

- a huge drop in speedup when there is a parallel generation of random numbers. Why?
 - look at the seed, figure out problem. (note - cards).
- current model which we are using is single program multiple data.
- Multiple Program multiple data works on different data.

SPMD:-

- ① Instruction level parallelism
 - ② Thread level parallelism
 - ③ Processor level parallelism
 - * There is a fixed IP for each processor.
- In MPMD, different processors execute different programs on different data.

The major difference in this case is that processors can access all memories, but in the other case it can only access local memories.

Performance depends on only data in cache.

Cache coherency:- If a piece of data in one processor is updated, it must be updated in the others.

Rule :- The problem must be broken into independent tasks with independent memory.

Challenges are of idling time and giving a wrong value.

Buffered Master mess

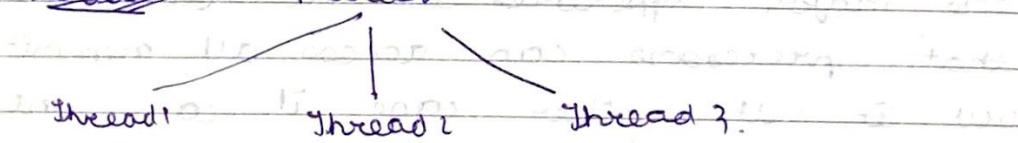
	Buffered	non-buffered	
working	(2) ✓	(1) ✓	
non-working			

Non-working : There is no idling time, but in working there is idling time.

- communications overheads will overlap.

- In non-working operation, it is unsafe to send data.
- Anything can be written.

Process



When 1 process has more than 1 thread, there is a hybrid code.

4 stages in MPI program :-

- 1) MPI include file.
- 2) Initialize MPI environment.
- 3) Do work and make message passing calls.
- 4) Terminate MPI environment.

`MPI_Init(&argc, &argv);`

`MPI_Finalize();`

Initialization :- It will initialize the MPI environment.

Termination :- Once we terminate, no MPI routine is called.

MPI processes are identified by their ranks.

- Rank ranges from: 0, 1, 2, ..., nproc - 1.
- nproc do not change during computation.

communicator :- A group of processes that will communicate with one another.

- They are used as arguments to all message transfer MPI routines.

ways to check if messages are sent correctly or not.

- 1) Point to point
- 2) Collective.

- + send and receive.
- 1) Destination / source
- 2) size of message.
- 3) data type of message (integer, float, etc.).
- 4) Protocol.
- 5) status
- 6) 1 tag so that receiver knows which message is to be received.

We are sending data: $A[0:20, 0:30]$

- Create argument, convert into ip and send

MPI: Integration using the ~~expression~~ rule.

Q) steps:-

1) division of work and time at once.
problem size = n

n in Processors = P

$$\frac{n}{P}$$

$A[0:n_0]$: several subprocessors belonging to one huge unit.

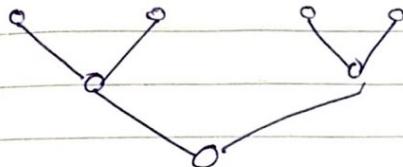
$$A\left[0:\frac{n}{P}\right], A\left[\frac{n}{P}+1: \frac{2n}{P}\right]$$

2) Idea on what data to work on.

3) collective communication is used.

Point to point communication

- copied from user memory to system buffer.
- sent from buffer over network to buffer of receiving process.
- Receiver copies data from system buffer to local user memory space.



Figures
infinite -
band switch

The switches that are used are infinite band switch.

Communication protocols are blocking and non-blocking. State variables can give information about the errors.

- * Architecture : Memory hierarchy, Pipeline, Importance of cache.
- * Data dependency.
- * Design of parallel algorithms. (complexity analysis.) ($O(N^k S)$) .
- * Parallel patterns \rightarrow Reduction, scan, convolutions, matrix multiplication, block multiplication, complexity analysis.
- * Shared memory programming \rightarrow OpenMP.
- * Performance Modelling.
- * Algorithmic analysis (parallel performance)
- * MPI.

* Inefficiency analysis:-

$$T_0 = p^{3/2} + p^{3/4} w^{3/4} = T_0(w, p)$$

$$w = k T_0(w, p)$$

w = problem size

p = no. of processes

$$w = k \tau k p^{3/2}$$

$$w = k p^{3/4} w^{3/4}$$

$$\cancel{w} w^{1-3/4} = k p^{3/4}$$

$$\Rightarrow w^{1/4} = k p^{3/4}$$

$$\Rightarrow w = k p^3$$

Thus, the problem size must be increased in $O(p^3)$, where p is no. of processes.

* deadlock: cannot move out unless all working statement not done.

* In this, circular locks can be broken down.

use cluster only for MPI rank numbers



rank	process 1	process 2
0	4	8

* ~~mpicollect -12~~

* ~~(example working) application condition~~