

## A Sources of Overhead in Parallel Programming -

different types of overheads -

1) Inter process communication - In any non-trivial parallel task system, intertask communication b/w processing elements are crucial for sharing data.

The time spent communicating data between processing overhead units slows down parallel processing.

2) Idling - In parallel systems, processing elements can sit idle for various reasons like unequal workloads,

If different processing elements have different workloads some processing elements may be idle during part of the time that others are working on a problem. In some parallel programs processing elements may be idle during In some parallel programs processing elements must synchronize at certain points during parallel exec?

3) Excess computation - Parallelizing a problem may necessitate

Excess computation - When the best sequential algorithm for task is hard to parallelize, we're left to work with less efficient parallel option. This means the parallel program does more work than optimal one, resulting in excess computation.

Teacher's Sign.: \_\_\_\_\_

## A) Different Performance metrics

1) Speedup - Speedup measures how much faster a parallel program runs compared to its sequential counterpart. It's calculated as the ratio of execution time of sequential prog. to the execution time of parallel prog.

$$S = \frac{T_s}{T_p}$$

2) Efficiency - measures utilization of resources in parallel system. Ratio of the speedup achieved by parallel program to no. of processor used.

$$E = \frac{S}{P}$$

Processing ele.

3) Overhead - Refers to additional time or resource consumed by parallelization related activities such as communication, sync, & load balancing. minimize overhead to achieve optimal performance.

4) Throughput - Measures the rate at which tasks or computations are completed in parallel system over given period.

5) Resource utilization - Measures extent to which system resources, such as processors, memory, & network bandwidth, are utilized. High resource utilization indicates efficient use of resources.

Teacher's Sign.: \_\_\_\_\_

## A Scalability of Parallel Systems.

Algorithm are said to be scalable if the level of parallelism increases at least linearly with size of problem.

A parallel computer system is said to be scalable if its efficiency can be fixed by simultaneously increasing machine size and the problem size.

## ? The Isoefficiency Metric of Scalability -

- (E)
  - 1) For a given problem size, as we increase the number of processing elements, the overall efficiency of parallel system goes down.
  - 2) The efficiency of a parallel systems increases if the problem size is increased while keeping the number of processing elements constant.

## A Isoefficiency Function -

Parallel execution time can be expressed as a function of problem size overhead function, and the number of processing elements. parallel runtime:

$$T_p = \frac{w + T_o(w, p)}{p}$$

$$\text{Speedup} = S = \frac{w}{T_p} = \frac{w}{w + T_o(w, p)}$$

$$\text{Efficiency } E = \frac{s}{p} = \frac{w}{w + T_o(w)} = \frac{1}{1 + T_o(w)/w}$$

A large isoeficiency function indicates a poorly Scalable System.

Teacher's Sign.: \_\_\_\_\_

## \* Minimum Execu<sup>n</sup> Time - (MET)

MET refers to the shortest possible time required to execute a given task or problem on a parallel system.

- The eq<sup>n</sup> to find the number of processing elements for which  $T_p$  (parallel runtime) is minimum.

$$\frac{d}{dp} T_p = 0.$$

- For example - min execu<sup>n</sup> time for adding n number parallel runtime  $T_p = \frac{n}{P} + 2\log P$

$$-\frac{n}{P^2} + \frac{2}{P} = 0.$$

$$-n + 2P = 0 \quad P = \frac{n}{2}$$

Substitute  $P = n/2$  we get.  $T_p^{mn} = 2\log n$

The processor time product for  $P = P_0$  is  $\Theta(n \log n)$

which is higher than  $\Theta(n)$ . Hence, the parallel system is not cost-optimal for the value of  $P$  that yields minimum parallel runtime.

## \* Minimum cost-optimal parallel time -

- Let  $T_p^{\text{cost-opt}}$  be the min cost-optimal parallel time

- If the isoefficiency func<sup>n</sup> of a parallel system is  $\Theta(P(F(P)))$  then a problem of size  $w$  can be solved cost optimally if and only if  $w = \Omega(F^{-1}(P))$ .

In other words for cost optimality  $P = O(F^{-1}(w))$

- For cost optimal system  $T_p = \Theta(w/P)$  Therefore

$$T_p^{\text{cost-opt}} = \Omega\left(\frac{w}{F^{-1}(w)}\right)$$

Example.

- The iso efficiency func'  $F(p)$  of this parallel sys  $\Theta(p \log p)$  we have  $p \leq n$

$$T_p^{\text{cost-opt}} = \log n + \log \left( \frac{n}{\log n} \right) = 2 \log n - \log \log n$$

$T_p^{\min}$  and  $T_p^{\text{cost-opt}}$  for adding n mo. are same  $\Theta(\log n)$ .

(\*) Amdahl's law - Amdahl's law governs the speedup of using parallel processors on a problem, versus using only one single processor, under the assumption that the problem size remains the same when parallelized.

- based on fixed problem size
- To keep efficiency of system fixed we have to increase both the size of the problem and the no. of processors simultaneously.
- It tells that for a given size, the speedup doesn't increase linearly as the number of processor increases. In fact Speed up tends to saturated.

$$S(n) \leq \frac{1}{f + \frac{(1-f)}{n}}$$

↑  
Serial part.      Parallel part

$f$  is the fraction of opera" in computa" performed sequentially where  $0 \leq f \leq 1$ . The max Speedup achievable by a parallel computer with  $n$  processor is

### Gustafson's Law -

- Instead of having fixed size problem we must assume that we have a fixed execution time. This is because in case of huge problem size, we will need to increase system size i.e. no. of processors.
- As the number of processing elements increases the total workload can also increase proportionally leading constant exec<sup>n</sup> time.
- Exec<sup>n</sup> time is fixed let serial exec<sup>n</sup> time be 's' and the parallel exec<sup>n</sup> time be 'p' for 'n' processor system. Thus total exec<sup>n</sup> time is  $s+nP$ .
- If exec<sup>n</sup> to be done on sequential sys, the exec<sup>n</sup> time will be  $s+nP$  thus speed up will be.

$$S'(n) = \frac{s+nP}{s+P} = \frac{s+nP}{1}$$

assuming time required on parallel system is 1  
i.e.  $s+P=1$  therefore  $P=1-s$

$$\begin{aligned} S'(n) &= \frac{s+n(1-s)}{s+(1-s)} \\ &= \frac{n+s(1-n)}{1} \end{aligned}$$

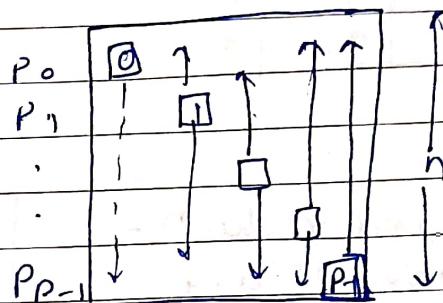
$$S'(n) = n + s(1-n)$$

Com  
ie  
ele  
Com

\* Matrix vector multiplication -

> Rowwise 1-D partitioning -

Required communica<sup>n</sup> - All to All broadcast.



Distribu<sup>n</sup> or full vector all to all broadcast

$P_0$	0	1	...	$P_{p-1}$
$P_1$	0	1	...	$P_{p-1}$
:	0	1	...	$P_{p-1}$
$P_{p-1}$	0	1	...	$P_{p-1}$

vector distribu<sup>n</sup> to  
each process after  
( broadcast )

$P_0$	0
$P_1$	1
:	.
$P_{p-1}$	$p-1$

final distribu<sup>n</sup> of the  
matrix and result vector y.

2) Rowwise 2-D partitioning

Req. communicat<sup>n</sup>.

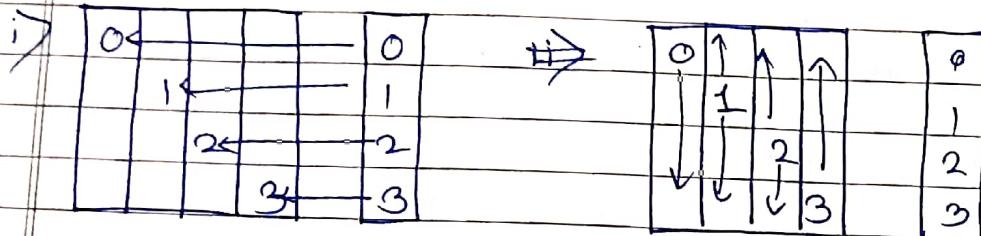
> One to one.

> one to All

> All to one.

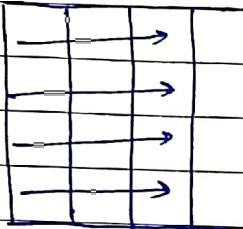
Teacher's Sign: \_\_\_\_\_

1) 1 to 1 communication -  $\Rightarrow$  one to all broadcast.



3) All to one reduction -  
whatever results are  
generated in a row by  
a particular process in a  
given element that will be  
reduced in the last row.

(cost complexity is  $\Theta(n^2 \log n)$ )



**Granularity** - Granularity in parallel computing refers to the size or level at which a problem is divided into smaller tasks or units of work for parallel execution.

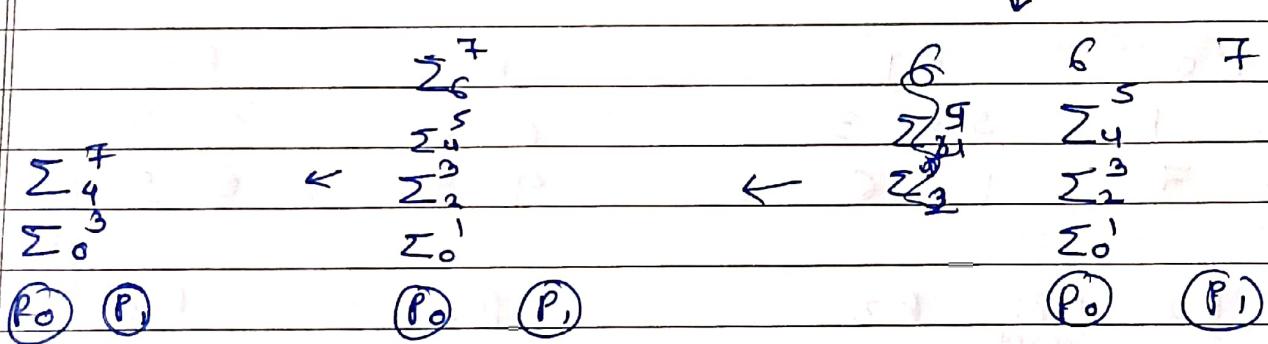
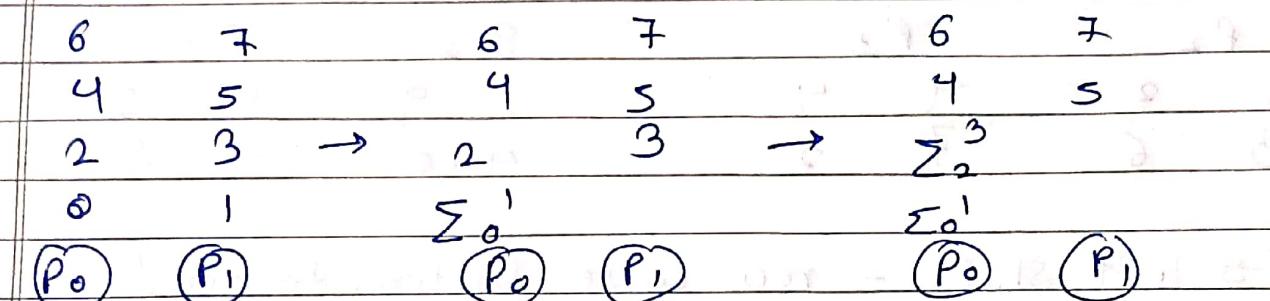
Effects of granularity on parallel systems.

- 1) Task overhead - It involves breaking problem into smaller tasks, which can lead to increased overhead due to task creation, scheduling & synchronization.
- 2) Communication overhead - Granularity can result in higher communication overhead as smaller tasks may require more frequent data exchanges & synchronization between processing elements.

Teacher's Sign.: \_\_\_\_\_

3) Load Balancing - Coarse-grained granularity involves longer tasks, where some processing elements have heavier workloads than others. This can result in underutilization of resources, reducing overall system efficiency.

4) computation decreases by factor of  $n/p$   
 communication increases by factor of  $n/p$ .



$$\Sigma_0^7$$

$(P_0)$   $(P_1)$

$P_0$  &  $P_1$  are physical processing element.

there are 8 virtual processing units. the computation will increase by  $n/p$  ie  $8/2 = 4$   
 ie one physical processor will handle 4 processing elements.

computation cost is  $\Theta(n \log p)$  if it is performed sequentially, cost will be  $\Theta(n)$   
 ie  $\Theta(n) < \Theta(n \log p)$ .

Teacher's Sign.: \_\_\_\_\_

## Matrix Matrix multiplication

$$A = \begin{bmatrix} 2 & 1 & 5 & 3 \\ 0 & 7 & 1 & 6 \\ 9 & 2 & 4 & 4 \\ 3 & 6 & 7 & 2 \end{bmatrix}, B = \begin{bmatrix} 6 & 1 & 2 & 5 \\ 4 & 5 & 6 & 5 \\ 1 & 9 & 8 & -8 \\ 4 & 0 & -8 & 5 \end{bmatrix}$$

$$P_0 \quad P_1 \\ \begin{matrix} 2 & 1 \\ 0 & 7 \end{matrix} \quad \begin{matrix} 5 & 3 \\ 1 & 6 \end{matrix}$$

$$P_0 \quad P_1 \\ \begin{matrix} 6 & 1 \\ 4 & 5 \end{matrix} \quad \begin{matrix} 2 & 5 \\ 6 & 5 \end{matrix}$$

$$P_2 \quad P_3 \\ \begin{matrix} 9 & 2 \\ 3 & 6 \end{matrix} \quad \begin{matrix} 4 & 4 \\ 7 & 2 \end{matrix}$$

$$P_2 \quad P_3 \\ \begin{matrix} 1 & 9 \\ 4 & 0 \end{matrix} \quad \begin{matrix} 8 & -8 \\ -8 & 5 \end{matrix}$$

left shift - row wise bottom to up  
 Up shift - column wise right to left.

$$P_0 \quad P_1 \\ \begin{matrix} 2 & 1 \\ 0 & 7 \end{matrix} \quad \begin{matrix} 5 & 3 \\ 1 & 6 \end{matrix}$$

$$P_0 \quad P_1 \\ \begin{matrix} 6 & 1 \\ 4 & 5 \end{matrix} \quad \begin{matrix} 2 & 5 \\ 6 & 5 \end{matrix}$$

↑ Up shift!

$$P_2 \quad \text{left shift} \quad P_3 \\ \begin{matrix} 9 & 2 \\ 3 & 6 \end{matrix} \quad \begin{matrix} 4 & 4 \\ 7 & 2 \end{matrix}$$

$$P_2 \quad P_3 \\ \begin{matrix} 1 & 9 \\ 4 & 0 \end{matrix} \quad \begin{matrix} 8 & -8 \\ -8 & 5 \end{matrix}$$

↓ .

$$P_2 \quad P_3 \\ \begin{matrix} 4 & 4 \\ 7 & 2 \end{matrix} \quad \begin{matrix} 9 & 2 \\ 3 & 6 \end{matrix}$$

$$P_0 \quad P_0 \quad P_3 \\ \begin{matrix} 6 & 1 \\ 4 & 5 \end{matrix} \quad \begin{matrix} 8 & -8 \\ -8 & 5 \end{matrix}$$

~~$$P_0 \quad P_1 \\ \begin{matrix} 2 & 1 \\ 0 & 7 \end{matrix} \quad \begin{matrix} 5 & 3 \\ 1 & 6 \end{matrix}$$~~

$$P_2 \quad P_3 \\ \begin{matrix} 1 & 9 \\ 4 & 0 \end{matrix} \quad \begin{matrix} 2 & 5 \\ 8 & 5 \end{matrix}$$

Teacher's Sign.: \_\_\_\_\_

## Final module. -



$$\begin{matrix} P_0 & P_1 \\ 2 & 1 & 5 & 3 \\ 0 & 7 & 1 & 6 \end{matrix}$$

$$\begin{matrix} P_0 & P_1 \\ 6 & 1 & 8 & -8 \\ 4 & 5 & -8 & 5 \end{matrix}$$

$$\begin{matrix} P_2 & P_3 \\ 4 & 4 & 9 & 2 \\ 7 & 2 & 3 & 6 \end{matrix}$$

matrix A

$$\begin{matrix} P_2 & P_3 \\ 1 & 9 & 2 & 5 \\ 4 & 0 & 6 & 5 \end{matrix}$$

matrix B

Multiply  $P_0$  with  $P_0$ ,  $P_1$  with  $P_1$ , and so on.

$$C_0 = \begin{matrix} 1 & 6 & 7 \\ 2 & 8 & 3 & 5 \end{matrix}$$

$$C_1 = \begin{matrix} 1 & 6 & -2 & 5 \\ -4 & 0 & 2 & 2 \end{matrix}$$

$$C_2 = \begin{matrix} 2 & 0 & 3 & 6 \\ 1 & 5 & 6 & 3 \end{matrix}$$

$$C_3 = \begin{matrix} 3 & 0 & 3 & 7 \\ 4 & 2 & 3 & 9 \end{matrix}$$

Perform left shift up shift  $\sqrt{P}$  times in our case  
 $P=4$  then  $\sqrt{P}=\sqrt{4}=2$  perform 2 times.

Prev model :-

$P_0$ left $P_1$	$P_0$	$P_1$	mat
2 1 shift 5 3	5 3	2 1	
0 7 ← 1 6      ⇒	1 6	0 7	A
P_2 left $P_3$	P_2	P_3	
4 9 6 4 ← 9 2	9 2	4 4	
7 2 3 6	3 6	7 2	

$P_0$	$P_1$	$P_0$	$P_1$	mat
6 1	8 -8	1 9	8 -8	
4 5	-8 5	4 0	-8 5	B
P_2 ↑ up shift	P_3 ↑ up shift	P_2	P_3	
2 5	2 5	6 1	2 5	
1 0	6 5	4 5	6 5	

now multiply

Teacher's Sign.: \_\_\_\_\_

matrix B.

$$\begin{array}{cc}
 P_0 & P_1 \\
 6 & 8 -8 \\
 4 & -8 5
 \end{array}
 \Rightarrow
 \begin{array}{cc}
 P_0 & P_1 \\
 1 & 9 \\
 4 & 0
 \end{array}
 \quad
 \begin{array}{cc}
 P_0 & P_1 \\
 2 & 5 \\
 6 & 5
 \end{array}$$

$$\begin{array}{cc}
 P_2 & P_3 \\
 1 & 9 \\
 4 & 0
 \end{array}
 \quad
 \begin{array}{cc}
 P_2 & P_3 \\
 1 & 8 -8 \\
 4 & -8 5
 \end{array}$$

now perform corresponding multiplications & adding with previous values -

$$C_0 = B_3 \cdot 52 + C_1 = 26 \cdot -14 \\
 53 \quad 44 \qquad \qquad \qquad 2 \quad 57.$$

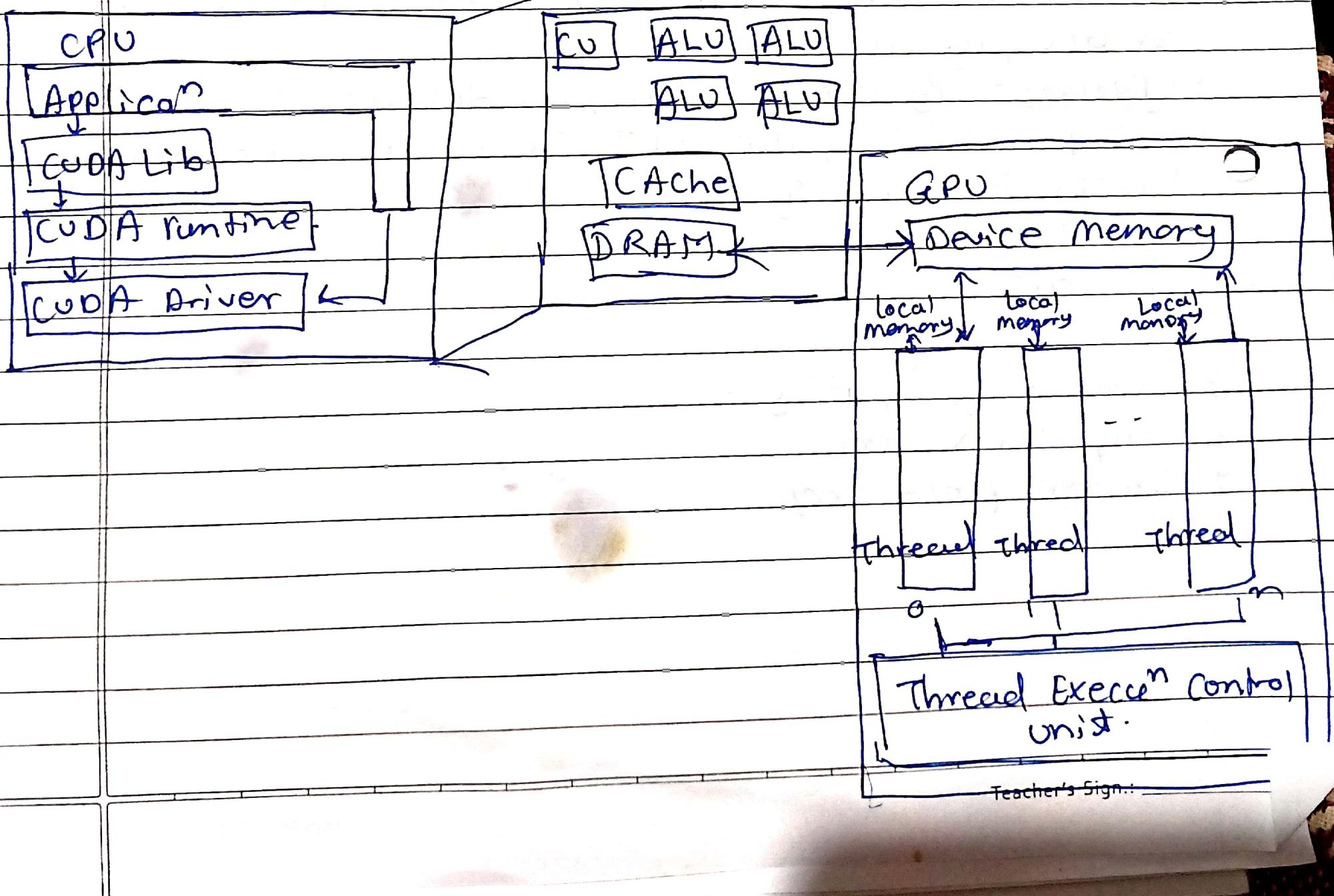
$$C_2 = 82 \cdot 55 + C_3 = 30 \cdot 25 \\
 57 \quad 96 \qquad \qquad \qquad 82 \quad -7.$$

## CUDA Architecture

### What is CUDA?

- Compute Unified Device Architecture is a parallel computing platform and application programming interface (API) developed by NVIDIA for general computing on GPU
- CUDA provides programmers with instructions that allow GPU acceleration for data-parallel computation
- CUDA can speed up computing applications across a wide range of domains including deep learning, analytics.
- Different programming languages supported by CUDA
  - CUDA C/C++ : primary language for CUDA programming. Extends C/C++ with GPU-specific features, such as kernel functions that execute on GPU & syntax for managing memory & launching parallel threads.
  - 2) CUDA Fortran - CUDA Fortran compiler for scientific computing, it extends Fortran with features for GPU programming.
  - 3) Python - PyCUDA library provides access to NVIDIA's CUDA parallel computer API. To write CUDA kernels directly in python syntax.
  - 4) MATLAB
- Applications
  - 1) Deep learning
  - 2) Scientific computing -
  - 3) Computer vision.
  - 4) Image processing.

Teacher's Sign.: \_\_\_\_\_



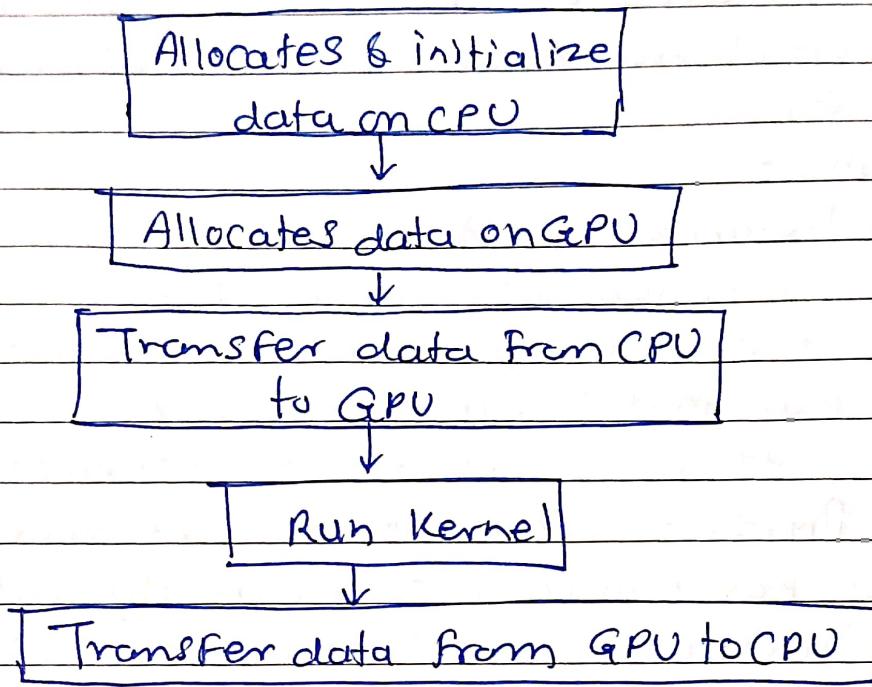
## Architecture

Date \_\_\_\_\_  
Page \_\_\_\_\_

- 1) CPU - Host comprises of CPU and its memory  
host code is executed by CPU. It allocates memory, launches kernel and handles data transfer between CPU and GPU.  
It contains -
  - 1) CUDA Libraries - built on top of CUDA runtime  
These are collections of pre-written functions & routines that provide optimized implementation.
  - 2) CUDA Runtime - built on CUDA Driver.  
Provides API for developers to interact with GPU
  - 3) CUDA Driver - Low-level software component responsible for communicating directly with GPU.  
Provides interface for OS to interact with GPU.
- 2) GPU -
  - 1) CPU allocates memory & prepares data for processing on the GPU. Data transferred from CPU's main memory to GPU's memory.
  - 2) CPU launches CUDA Kernel Funcn. Kernel launches are initiated by CPU cmd executed on the GPU.
  - 3) GPU executes Kernel Function in parallel on multiple threads organized into blocks & grids.
  - 4) After Kernel execu", the results are stored in GPU's memory. CPU requests the result by GPU's memory.

Teacher's Sign: \_\_\_\_\_

## \* Describe CUDA-C program flow



- 1) Host (CPU) - Allocates memory on GPU transfer data between CPU & GPU launches kernel functions and get result
- 2) Allocate data on GPU - CPU allocates memory on GPU for input & output data. Input transferred from CPU's main memory to GPU's memory.
- 3) Kernel - Kernel Function is a CUDA-C Function that executes in parallel on the GPU. GPU launches multiple threads to execute Kernel Function.
- 4) Retrieve result - CPU retrieves the result by transferring data back from the GPU's memory to CPU's main memory.

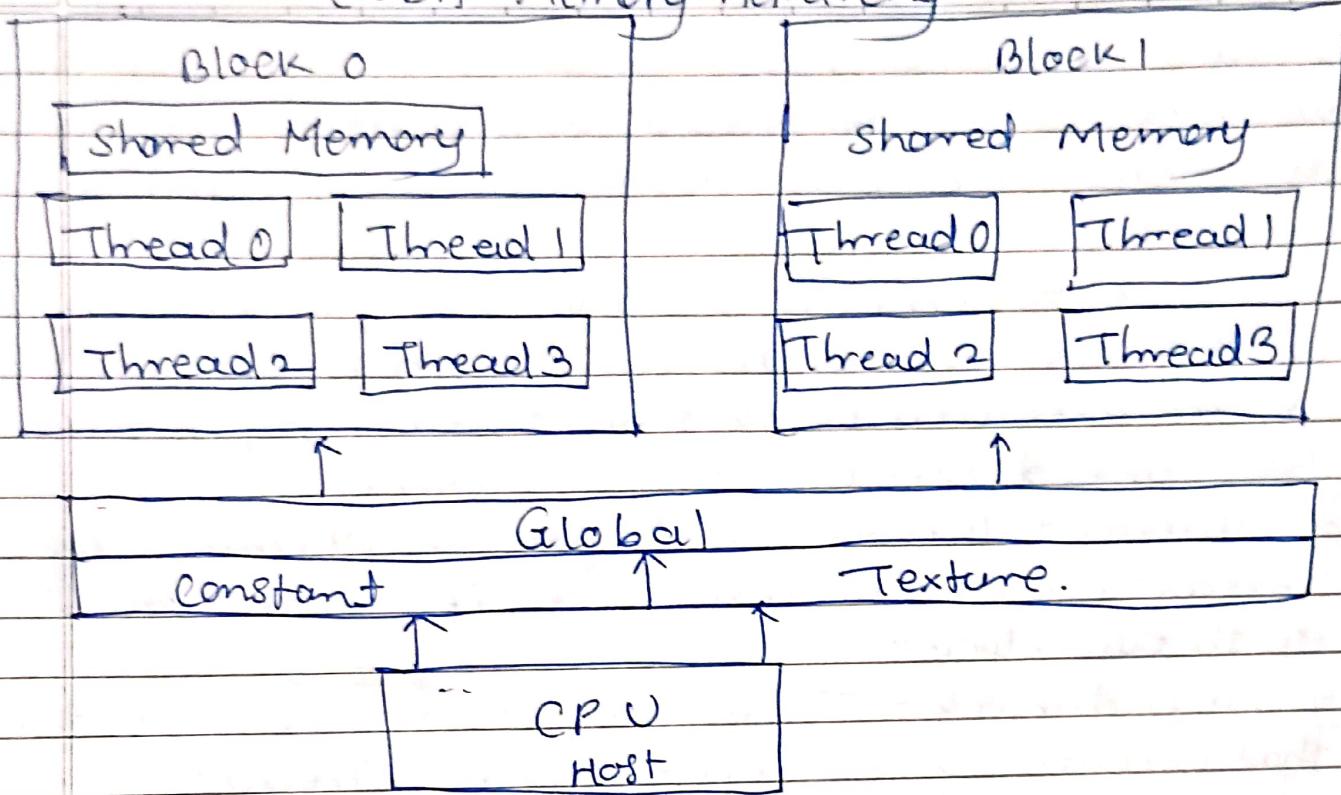
Teacher's Sign.:

## # Terms in CUDA

- 1) GPU (Graphical Processing Unit) - GPU is specialized processing device designed to accelerate the rendering of images for display devices. Used for parallel computing.
- 2) Kernel - Kernel is a func<sup>n</sup> that runs in parallel on the GPU. Kernel are written in CUDA C / C++ & are executed by multiple threads in parallel. Each thread executes the same kernel.
- 3) Thread - Thread is a single exec<sup>n</sup> unit in GPU. Threads execute in parallel on GPU. Some organized in thread blocks.
- 4) Thread block - Thread block is group of threads that execute some kernel code. Threads within some block can communicate and synchronize using shared memory.
- 5) Grid - grid is colle<sup>n</sup> of thread blocks that execute same kernel func<sup>n</sup>. represents the highest level of parallelism in CUDA, organizing multiple thread blocks per exec<sup>n</sup>.
- 6) Device - It refers to GPU used for parallel computation. CUDA programs can run on multiple devices, each with its own memory & processing unit.
- 7) Host - Host is CPU, host initiates & manages exec<sup>n</sup> of CUDA programs including data transfer.
- 8) Global memory - it's the main memory space accessible by all threads in GPU, used for storing data that needs to be shared among threads.
- 9) Shared memory -
- 10) CUDA Runtime API - Provides set of functions for managing device & memory & kernel exec<sup>n</sup>, CUDA runtime API interact with CUDA enabled GPU's.

Teacher's Sign.: \_\_\_\_\_

# CUDA Memory hierarchy



- > Device pointers point to GPU memory, host pointer points to host memory
- > Each thread has private local memory and shared memory of each thread is visible to all threads of some blocks.
- > Global Memory - There's large amount of global memory. All running threads can read & write global memory and so can the CPU. The functions `CudaMalloc`, `CudaFree`, `CudaMemcpy` all deal with global memory.
- > Constant Memory - It has its own cache. Not related to L1 & L2 cache of global memory. All threads have access to some constant memory but they can only read can't write. CPU sets values in constant memory before launching kernel

Teacher's Sign:

- > Shared Memory - Shared memory is very fast. Shared memory usage is used to enable fast communication between threads in a block. Bank conflicts can slow down access down. It's fastest when all threads read from different banks.
- > Caches - There's also an L1 cache per multiprocessor and an L2 cache which is shared between all multiprocessors.
- > Local Memory - This is also part of main memory of the GPU so it's slow. Local memory is used automatically by threads when we run out of registers or when registers can't be used.
- > Texture memory - It is specialized read-only memory optimized for texture fetching, suitable for texture mapping in graphics application.

### ~~1~~ CUDA-C Kernel Function execn -

```
#include <stdio.h>
__global__ void multiply(int a, int b, int* res) {
    *res = a * b;
}

int main() {
    int a = 5, b = 10;
    int* dev_res;
    cudaMalloc((void**) &dev_res, sizeof(int));
    multiply<<<1, 1>>>(a, b, dev_res);
    int res;
    cudaMemcpy(&res, dev_res, sizeof(int),
              cudaMemcpyDeviceToHost);
    printf("%d * %d = %d\n", a, b, res);
    cudaFree(dev_res);
    return 0;
}
```

Teacher's Sign.: \_\_\_\_\_

- 1) Kernel declaration - The multiply func<sup>n</sup> is declared as a CUDA kernel with the --global-- qualifier. This means it will run on the GPU and can be called from the CPU.
- 2) Kernel invoke - The kernel is launched from the CPU with one block and one thread ( $\langle\langle 1, 1 \rangle\rangle$ ) (as only one thread is needed since we are multiplying two numbers).
- 3) Memory Access ( $*res = a * b$ ) - The Thread multiplies two input integers ( $a & b$ ) and stores it in the "res" in memory.
- 4) Data Transfer - `CudaMemcpy` result copied from GPU to CPU after kernel execu<sup>n</sup>
- 5) Output - Print output & Free the memory (`CudaFree(dev_res)`).

# Hello world.

```
#include <stdio.h>
__global__ void hellocuda() {
    printf("hello world");
}
int main() {
    hellocuda <<1,10>>();
    CudaDeviceSynchronize();
    return 0;
}
```

3  
nvcc hello-world.cu  
./a.out.

## \* Block Dimensions -

Block dimension refers to the number of threads per block. It defines size of a block which

is a group of threads that execute the same kernel function. It is specified using a `<numThreadPerBlock>` argument when launching.

\* Grid Dimension - Refers to no. of blocks in grid. It defines the overall size of the grid, which consists of multiple blocks that execute some function or kernel function, specified using `<numBlocks>` argument when launching a kernel.

## \* CUDA - C++ for vector addition -

```
#include <cuda.h>
__global__ void Add_Vec( float* a, float* b, float* c, int N )
{ int i = blockIdx.x * blockDim.x + threadIdx.x ;
  if (i < N) {
    c[i] = a[i] + b[i];
  }
}
```

```
int main () { int N = 1024;
```

```
float* ha = (float*) malloc (N * sizeof (float));
float* hb = (float*) malloc (N * sizeof (float));
float* hc = (float*) malloc (N * sizeof (float));
```

```
float* da, *db, *dc;
```

```
(cudaMalloc (&da, ha, N * sizeof (float)), cudaMemcpy
```

```
(cudaMalloc (&db, hb, N * sizeof (float))),
```

```
(cudaMalloc (&dc, hc, N * sizeof (float)), cudaMemcpy HostToDevice),
```

```
(cudaMalloc (&dc, hc, N * sizeof (float)), cudaMemcpy HostToDevice));
```

```
int Threads_Per_Block = 256;
```

```
int Blocks = (N + Threads_Per_Block - 1) / Threads_Per_Block;
```

```
vector Add_Vec << blocks, Threads_Per_Block >>
```

```
(da, db, dc, N); printf(h-c[i]);
```

```
(cudaFree (da));
```

```
(cudaFree (db));
```

```
(cudaFree (dc)); return 0; }
```

Teacher's Sign: \_\_\_\_\_

```

#include <stdio.h>
#include <cuda.h>

__global__ addInt(int *a, int *b, int *c) {
    int i = blockIdx.x * blockDim.x + threadIdx.x;
    if (i < 1) { c[i] = a[i] + b[i]; }
}

int main() {
    int a = 5, b = 10, c;
    int *d_a, *d_b, *d_c;
    cudaMalloc(&d_a, a, cudaMemcpyHostToDevice);
    cudaMalloc(&d_a, &a, sizeof(int));
    cudaMalloc(&d_b, &b, sizeof(int));
    cudaMalloc(&d_c, &c, sizeof(int));
    cudaMemcpy(d_a, &a, sizeof(int), cudaMemcpyHostToDevice);
    cudaMemcpy(d_b, &b, sizeof(int), cudaMemcpyHostToDevice);
    cudaMemcpy(d_c, &c, sizeof(int), cudaMemcpyHostToDevice);

    addInt <<< 1, 1 >>> (d_a, d_b, d_c);
    cudaMemcpy(&c, d_c, sizeof(int), cudaMemcpyDeviceToHost);
    printf("Sum: %d", c);
    cudaFree(d_a);
    cudaFree(d_b);
    cudaFree(d_c);
    return 0;
}

```

\* Kernel launch - Launching kernel involves executing the kernel function on GPU with CUDA configuration.

done by using `<<<..>>>` syntax.

arguments → Grid size - Specifies no. of blocks in grid.

2) Block dimension - Specifies no. of threads per block

3) Dynamic shared memory - Allows allocation of shared memory per block dynamically.

Teacher's Sign.: \_\_\_\_\_



Scanned with OKEN Scanner

## Thread Hierarchy in Cuda memory

Date \_\_\_\_\_  
Page \_\_\_\_\_

- 1) Threads - Individual unit of work executing the kernel function, each thread has its own Thread ID (threadIdx). can access, register, Shared memory & global memory.
- 2) Blocks - Groups of threads that cooperate & share data through shared memory.  
Each block has unique block ID (blockIdx)
- F Threads within can synchronize their execution  
— Syncthreads
- 3) Grid - collection of thread blocks launched on the GPU. Configured during kernel launch by specifying the no. of blocks in each dimension.

### \* Parallel DFS -

- DFS is a widely used algorithm for graph traversal essential in applications such as Search & pathfinding.
- Parallel DFS involves exploring distinct branches of the search tree simultaneously, potentially yielding speedup.
- Basic parallel DFS - Multiple threads or processes concurrently explore separate branches of the search tree, dividing the graph into subgraphs. This method is effective when branches are of similar size.
- Workpile-based parallel DFS - Using workpile data structure, multiple threads access nodes for exploration. Threads process nodes, expand neighbours & new nodes to the workpile. It avoids graph partitioning.
- Recursive parallel DFS - Threads are assigned subsets of the graph to explore, initiating DFS from specific nodes & recursively exploring adjacent nodes.
- Task-based parallel DFS - Dynamically creating tasks or subproblems, threads independently explore assigned subgraphs. Threads generate tasks for unexplored nodes, dynamically scheduling them.

### \* Parallel BFS -

- BFS is a graph traversal algorithm that explores all vertices of a graph in breadth first order. Starts from a source node and visits neighbours before moving to their neighbours.
- Level-Synchronous parallel BFS - The graph is divided into levels, with each level containing vertices at a specific distance from source. Multiple threads

Teacher's Sign.: \_\_\_\_\_

or processes work in parallel, each responsible for visiting vertices at a particular level.

Initially Source vertex is assigned to 1<sup>st</sup> level & subsequent levels are processed in parallel.

3) Workpile-based parallel BFS - Similar to parallel DFS, a workpile data structure manages nodes to be explored. Multiple threads access workpile concurrently, processing nodes until it's empty.

4) Recursive parallel BFS - Threads are assigned subsets of graphs, initiating BFS from specific nodes & exploring adjacent nodes recursively.

5) Hybrid approach -

### \* Parallel Bubble Sort -

→ Parallelized version of classic bubble sort algorithm. It involves dividing the sorting process among multiple threads or processes to speed up operations.

2) Divide & Conquer - In parallel bubble sort, the list to be sorted is divided into smaller sublists, and each sublist is sorted independently by threads.

3) Parallel Comparison & Swapping - Each thread or process compares adjacent elements within its sublist and swaps them if necessary, similar to sequential bubble sort.

4) Synchronization - Synchronization mechanisms are used to ensure that threads or processes coordinate their sorting operations correctly.

5) Combining Results - Once all sublists are sorted independently, the sorted sublists are combined or merged to produce final list (sorted).

Teacher's Sign.: \_\_\_\_\_

## \* odd-even Transportation in Seq. bubble sort.

Date: \_\_\_\_\_  
Page: \_\_\_\_\_

odd phase.

1	2	3	4	5	6	7	8
3	2	3	8	5	6	4	1

Start from odd index and pair with immediate next element in sequence and then swap according to condition.

1	2	3	4	5	6	7	8
even	2	3	8	5	6	1	4

even phase.

odd	2	3	3	5	1	8	1	6	4
-----	---	---	---	---	---	---	---	---	---

swap

even	2	3	3	5	1	8	4	6
------	---	---	---	---	---	---	---	---

swap.

odd	2	3	1	3	5	4	8	6
-----	---	---	---	---	---	---	---	---

even	2	3	1	3	4	5	6	8
------	---	---	---	---	---	---	---	---

odd	2	1	3	3	4	5	6	8
-----	---	---	---	---	---	---	---	---

so 1 2 3 3 4 5 6 8. → sorted array

In sequential bubble sort it is working 1 element per process.

Parallel Formulation of odd-even Transposition  
Assume there are Two threads Thread 1 & Thread 2 and an array

$$\text{arr} = [3, 1, 4, 2, 5].$$

- > In odd phase Thread 1 is assigned to arr[1] and Thread 2 is assigned to next odd index ie arr[3] In odd phase Thread assigned to odd index.
- Thread 1 has arr[1] = 3 it compares with immediate next element in array ie arr[2] Compares their values & swaps them if necessary
- So Thread 1 compares 3 with 1  
(arr[1] > arr[2], swap) we get [1, 3, 4, 2, 5]
- Same for Thread 2

Thread 2 compares 3 with 4

(arr[3] > arr[4], swap) we get.

$$[1, 3, 2, 4, 5]$$

- > In Even phase Thread 1 is assigned to even index.  
So Thread 1 is assigned to arr[2] compares with arr[3] & Thread 2 is assigned to arr[4] compares with arr[5]

Thread 1 compares 3 with 2

(arr[2] > arr[3], swap) we get.

$$[1, 2, 3, 4, 5].$$

Thread 2 compares 4 with 5.

(arr[4] < arr[5], no swap) we get.

[1, 2, 3, 4, 5] → Sorted array.

$$\text{Time cost} = \Theta\left(\frac{n}{p} \log \frac{n}{p}\right) + \Theta(n) + \Theta(n).$$

Teacher's Sign.: \_\_\_\_\_

## ❖ Parallel merge sort -



→ Parallelized version of sequential merge sort to improve performance on multicore processors.

Sorting process is divided among multiple threads.

→ Consider unsorted array [38, 27, 43, 9, 82, 10]

→ Initial splitting (divide phase) -

Divide array into smaller subarrays.

It can be performed in parallel by dividing input array among multiple threads.

→ Sorting - Sort each subarray independently.

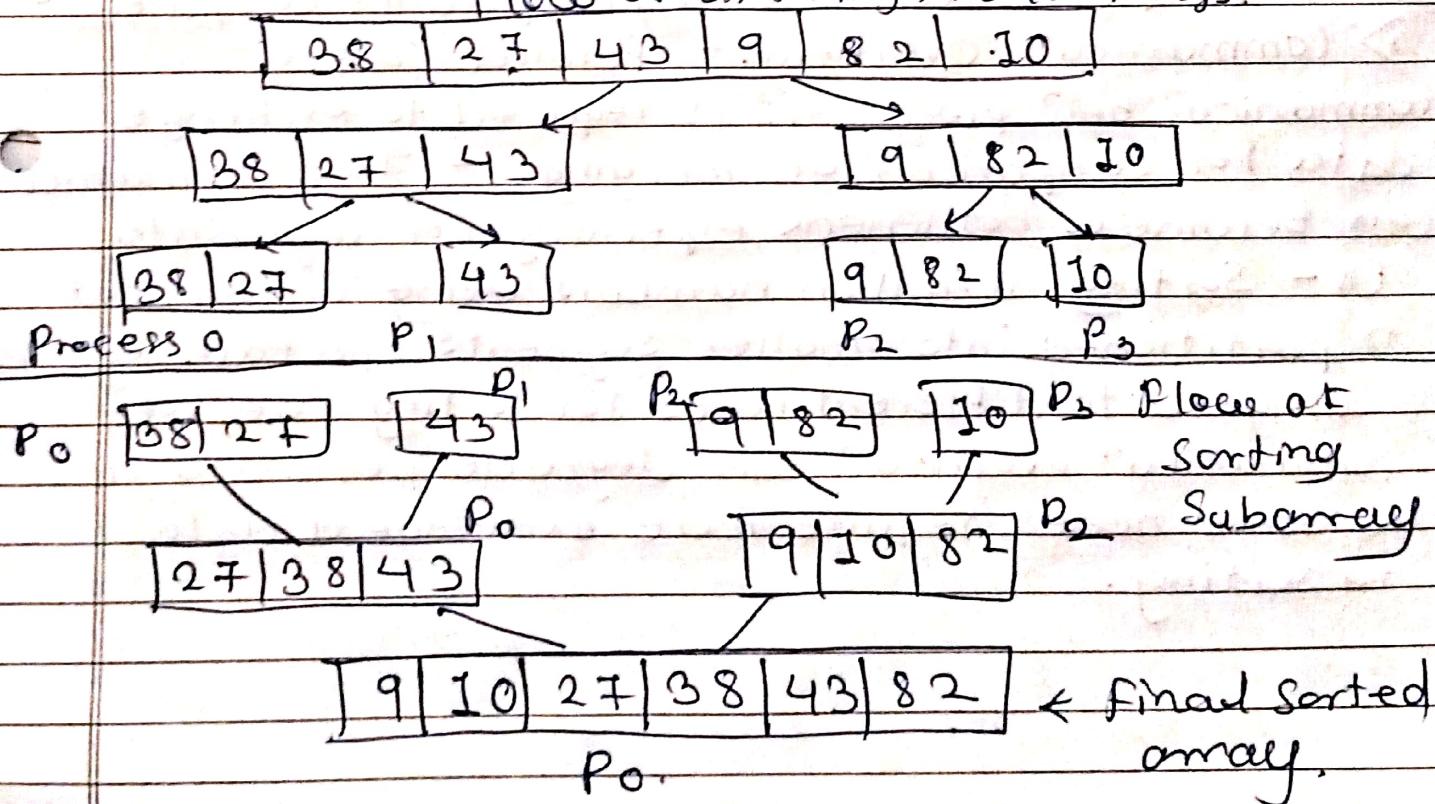
→ Using parallel sorting programming model like OpenMP or MPI, to sort each half recursively.

→ Merge - Once both halves are sorted they need to be merged into single sorted list.

This can be done sequentially by any processor.

For example, two processors can merge [27, 38] & [9, 43] simultaneously resulting [3, 27, 38, 43].

Flow of dividing in subarrays.



Teacher's Sign.: \_\_\_\_\_

C++

```
void parallel_merge(int* data, int low, int high){\n    if (low < high){\n        int mid = (low + high) / 2;\n\n        #pragma omp parallel sections num_threads(2)\n        {\n            #pragma omp section\n            {\n                parallel_merge(data, low, mid);\n            }\n            #pragma omp section\n            {\n                parallel_merge(data, mid+1, high);\n            }\n        }\n        merge(data, low, mid, high);\n    }\n}\n\nvoid merge(int* data, int* left, int* right).\n{\n    // Function of merging.\n}
```

- What are issues in Sorting on parallel computers -
- Communication overhead - Parallel Sorting, communication between processors is required to exchange data for comparison or merging. - This communication can become a bottleneck, especially for large data.  
Ex - Sorting a million numbers, using the dataset is partitioned into smaller segments and each segment is assigned to different nodes for sorting. Excessive communication overhead can degrade performance as some nodes require more exchange of data for sorting.

- 3) Load balancing - The size of subproblems
  - If assigned to processors affects performance if subproblems are too small, the overhead of managing threads can outweigh the benefits of parallelization
  - Ex - while performing sorting on large dataset dataset partitioned into smaller segments due to variation in the sizes of segments or different processing speeds between nodes. Some nodes may finish earlier and may remain idle.
- 4) Synchronization challenges - Parallel sorting algo often require synchronization points where processors need to wait for each other before processing.
  - Ex - Odd-even Transposition sort, threads need to compare elements with neighbors, if not properly sync threads might access the same data.
- 5) Limited parallelism in sorting Algo - Some sorting algo have inherent limitations when parallelized. Ex - insertion sort.

- 6) Kubernetes - It is an open source Container Orchestration platform, It automates deployment, scaling & management of containerized application.
  - Features.
  - 1) Container - Kubernetes automates the deployment, Scaling & management of containerized apps.
  - 2) Scalability - Easily scale application up or down based on demand. Kubernetes can automatically allocate resources & manage containers.

Teacher's Sign.



Scanned with OKEN Scanner

- 3) Self-healing - If container crashes, kubernetes automatically restarts it.
- 4) Load balancing - Kubernetes can distribute incoming traffic across multiple instances of application.
- 5) Storage - Integrates with various storage solutions states allocating persistent storage for containerized app.

### Application:

- 1) Modern web Apps
- 2) Cloud native Dev
- 3) Machine learning
- 4) HPC

- ❖ Application of GPU
- 1) Scientific computing
- 2) Machine learning
- 3) Data Science Analytics
- 4) Media & Entertainment.
- 5) AI
- 6) Forecasting