

2. Parallel Programming.

* principle of Parallel Algorithms design:

- An algo. provides step by step solnⁿ of the given problem. An algo. accepts from the user and based upon I/P. after performing the defined computations provides the ~~app.~~ O/P.
- The basic steps in the design of HPC algo. are: partitioning of overall computation into smaller computation and assignment of these smaller ~~components~~ computations into diff. processors.
- Decomposition, tasks & Dependency graphs:

* Decomposition:

- The overall computation can be partitioned into ~~small~~ number of small size computation so that these computⁿ execute in parallel.
- The decomposition deals with the approaches of partitioning the overall computⁿ into sub-parts. When a computⁿ is divided into many small tasks, it is referred as fine-grained decomposition.
- On other hand, the coarse-grained decompoⁿ consists a small number of large tasks.

* Tasks:

- In progⁿ, the basic unit of computⁿ is referred as a task and is controlled by OS.
- In context of HPC progⁿ, the task are unit of computation based upon that the overall computⁿ is decomposed.

* Task-Dependency graph:

- Is a directed acyclic graph, Typically a graph is a collection of nodes and edges, the task-dependency graph also contains nodes & edges.
- The node in task-dependency graph is a task whereas edges betⁿ any two nodes represent dependency betⁿ them e.g.

There is edge exists betⁿ 2 nodes T₁ & T₂.

If T₂ must be executed after T₁.

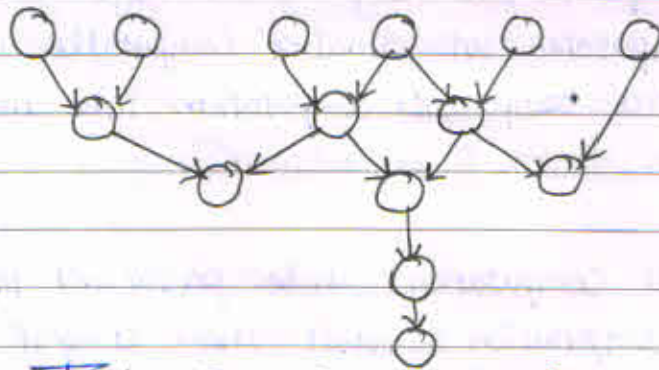
Consider the task-Dependency graph and find-out the following :-

i) Maximum degree of concurrency

ii) Critical path length.

iii) Total amount of work.

iv) Average degree of concurrency.



Task Dependency Graph.

→ 1) Max. degree of Concurrency. →

The max² nos of tasks allowed to execute in 111 is termed as the max² degree of concurrency.

In the given graph, max² degree of concurrency is 6.

2) Critical path length:-

- There are 2 types of node included in the task-dependency graph namely start node and finish node.

- Start nodes are nodes with no incoming edges

whereas nodes with no outgoing edges are called as finish nodes.

- The critical path in task-dependency graph is the longest directed path bet² any pair of start and finish node.

- The quantity referred as critical path length is the sum of weight of nodes on a critical path.

- In the given graph, the critical path length is 5.

3) Total amount of work :-

- Here, the total amount of work is 14 if it is assumed that the each tasks takes one unit of time.

4) Average degree of Concurrency:

- The average degree of Concurrency deals with the avg nos of tasks allowed to execute in parallel.

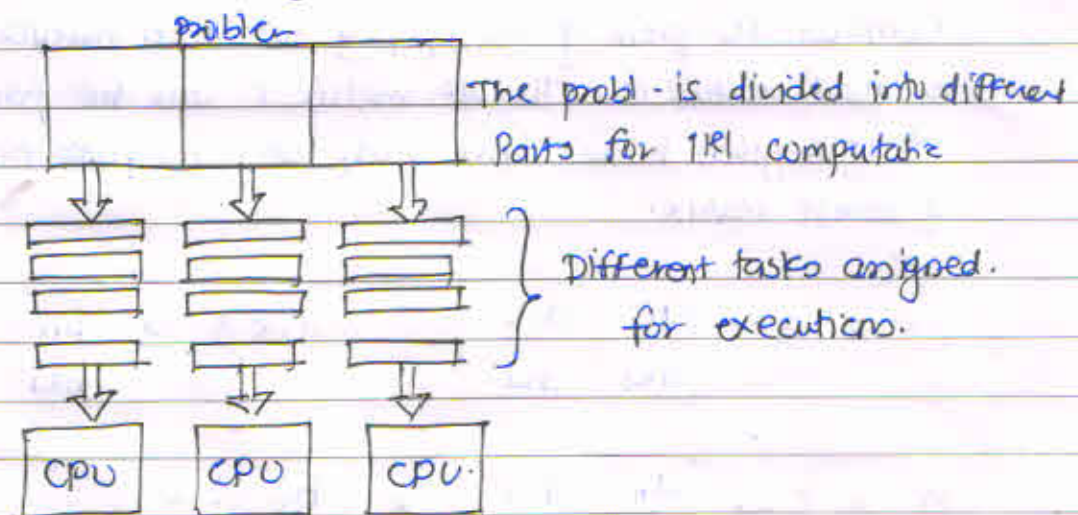
$$\text{Avg degree of concurrency} = \frac{\text{Total amount of Work}}{\text{Critical Path length}}$$

e.g.

$$\text{Avg degree of Concurrency} = \frac{14}{5} = 2.8$$

* Processes and Mapping:

- The environment in which a problem is being solved by means of parallel computation usually has several running processes simultaneously. These processes communicate with each other for need of synchronisation and exchanging of information.
- The program written for such platforms are called parallel programs. A single parallel program can have several processes for handling different tasks simultaneously.



Parallel Execution of tasks.

* Mapping:

In parallel program design, it is required to specify where each task is to execute. As we know process performs a task. In this approach of context is used for assignment of tasks to process is called mapping.

* De-Composition techniques :

— Data Decomposition :-

- The large data sets involved in the problem can be divided into smaller part and process independently by several computer concurrently.
- The data composition is common tech used for concurrent processing of data.
- There are basically 2 steps in this technique.
In 1st step, the overall I/P data is required for performing the defined computation is divided into multiple part.
- In 2nd step, it is based upon the division of overall computation into nos of tasks handled on the partitioned data.
- The actual ops implemented on these tasks are similar but data upon which these ops work are different.

Example:-

Matrix Multiplication :-

Consider the probl. of multiplying two $n \times n$ matrices. A and B to yield matrix C. The O/P matrix C can be partitioned into 4 task given below. Here each task computes one element of result matrix.

$$\text{Matrix A} \rightarrow \begin{bmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \end{bmatrix}$$

$$\text{Matrix B} \rightarrow \begin{bmatrix} b_{11} & b_{12} \\ b_{21} & b_{22} \end{bmatrix}$$

$$\text{Matrix C} \rightarrow \begin{bmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \end{bmatrix} \times \begin{bmatrix} b_{11} & b_{12} \\ b_{21} & b_{22} \end{bmatrix}$$

$$\text{Matrix C} \rightarrow \begin{bmatrix} c_{11} & c_{12} \\ c_{21} & c_{22} \end{bmatrix}$$

$$\begin{bmatrix} a_{11} \cdot b_{11} + a_{12} \cdot b_{21} & a_{11} \cdot b_{12} + a_{12} \cdot b_{22} \\ a_{21} \cdot b_{11} + a_{22} \cdot b_{21} & a_{21} \cdot b_{12} + a_{22} \cdot b_{22} \end{bmatrix}$$

* Recursive Decomposition :

- As we know, divide & Conquer is one of the strategies of solving the Computational problem. It is based on 2 ideas:
1st approach says to solve prob. directly if it is trivial and in 2nd approach decompose the prob. into smaller prob. if it cannot be solved as it is and solve the smaller prob.
 - The recursive decomposition is based on providing concurrency in problems that can be handled in divide and conquer strategy.
 - The prob. to be solved using recursion is divided into multiple sub prob. provided that each of these sub problems is independent each sub prob. can be solved individually by dividing further into other sub problem & solved using recursion.
- In programming a function or procedure calls it-self is called as recursion.

e.g.

```
main ()
```

```
{
```

```
    int i;
```

```
    i = 10;
```

```
    demo(i);
```

```
}
```

```
demo (int count)
```

```
{
```

```
    count--;
```

```
    printf ("The value of the count is %d \n", count);
```

```
    if (count > 0)
```

```
        demo (count);
```

```
    printf ("The count is %d \n", count);
```


* Exploratory Decomposition:

- There are situations where the probl. decomposⁿ goes hand in hand with it's exⁿ. Such problems typically involve the exploration or search of state space & solⁿ.

The search space probl. is divided into smaller parts & each smaller part is searched concurrently till the point at which the expected solⁿ is found.

Example:-

- Consider a state space searching problem, such as finding a solution to puzzle problem.

The steps are:-

- 1) The computⁿ required for the decompⁿ can be divided into multiple tasks, where each task is searching for a different portion of the search space.
- 2) The task of finding the shortest path from initial to final configⁿ now translates to finding a path from one of these newly generated nodes to final configⁿ.
- 3) The 15 Puzzle probl. is typically solved using tree search techⁿ. Starting from initial configⁿ, all possible successors are generated.
- 4) In the 15 puzzle problem, there are 15 tiles numbered 1 through 15 are arranged in 4x4 grid shown in fig.
 - one tile is left blank, so that moves can be made.The 4 possible moves here are represented as moves up, down, left & right. The start & final configⁿ are specified.
- 5) This problem is characterized by objective determining any seqⁿ of moves for a shortest sequence of moves.

The solⁿ of the probl. must be searched from an arbitrary state

1	2	3	4
5	6	↑	8
9	10	7	11
12	13	14	15

1	2	3	4
5	6	7	8
9	10	←	11
12	13	14	15

1	2	3	4
5	6	7	8
9	10	11	↑
13	14	15	12

1	2	3	4
5	6	7	8
9	10	11	12
13	14	15	

Task 4

1	2	3	4
5	6	7	8
9	10	11	
3	14	15	12

1	2	3	4
5	6	7	
9	10	11	8
13	14	15	12

1	2	3	4
5	6	7	8
9	10	11	12
13	14	15	

1	2	3	4
5	6	7	8
9	10		11
13	14	15	12

Task 3

1	2	3	4
5	6	7	8
9	10	11	
13	14	15	12

1	2	3	4
5	6	7	8
9	10	11	
13	14	15	12

1	2	3	4
5		7	8
9	6	10	11
13	14	15	12

1	2	3	4
5	6	7	8
9	4	10	11
13		15	12

Task 2

1	2	3	4
5	6		8
9	10	7	11
12	13	14	15

1	2	3	4
5	6	8	
9	10	7	11
13	14	15	12

1	2	3	4
5		8	8
9	10	7	11
13	14	15	12

1	2		4
5	6	3	8
9	10	7	11
13	14	15	12

1	2	3	4
5	6	7	8
9	10		11
13	14	15	12

Task 1

1	2	3	4
5	6	7	8
9	10	15	11
13	14		12

1	2	3	4
5	6	7	8
9	10	15	11
13		14	12

1	2	3	4
5	6	7	8
9	10	15	11
13		14	12

1	2	3	4
5	6	7	8
9	10		11
13		14	12

1	2	3	4
5	6	7	8
9	10		11
13	14	15	12

* Speculative Decomposition %

- It is used in situation when prog has many options to in terms of branches based upon the output of the other part that precedes it. The situation can be described in terms of specific task. Consider a task T₁ involves in doing the computation C₁ and going to produce output O₁.
 - The next computation to be performed may be decided by based upon O₁ task T₁. In similar fashion, other tasks can start exe. concurrently in the next stage. The situation can better compared with the switch statement in C programming.
 - The switch statement works based upon the value of expr² on which it is based and only the corresponding case statement executes. Some or all of the cases executes in advance in speculative decomposition.
 - The situation in which values of the expr² is known then the result from the comput² to be executed in that case will be kept. The anticipation about the possible computation causes the performance gain in the statement execution.
- e.g. consider foll. Switch statement in seq & parallel version.

Sequential	parallel.
Compute expr;	slave(i)
switch (expr)	{
{	Compute e _i ;
case 1: Compute-e ₁ ;	wait (request);
break;	if (request)
Case 2: Compute-e ₂ ;	send (e _i , o);
break;	}
case 3: compute-e ₃ ;	Master ()
break;	{
---	Compute expr)
}	Switch (expr)
	{
	case 1):
	send (request, 1);
	Receive (a ₁ , i);

	}

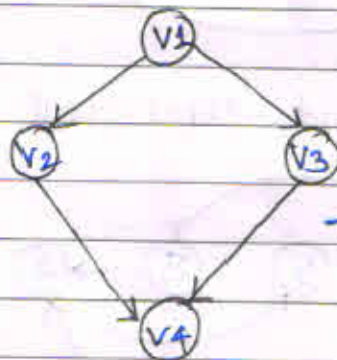
Example: Topological Sorting:

- In directed cyclic graph, topological sorting is used to provide an ordering of the vertices.

e.g. two vertices u & v of a graph G . In this graph, when path exists betⁿ u & v , then v appears after u in the ordering.

- The graph should acyclic property, otherwise for an edge represented as (u, v) there would be a path from u to v and from v to u , & therefore the ordering can't be obtained.

- Now, consider there are nos of tasks we need to be performed in which some task depend on the others and it is possible to do only one. The task can be organized in the dependency graph.



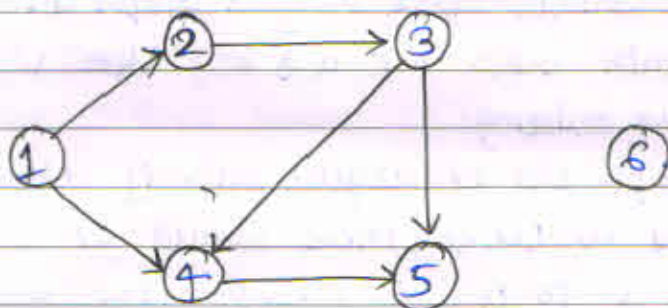
The legal ordering are: V_1, V_2, V_3, V_4
and V_1, V_3, V_2, V_4

* The steps of the logical sorting algo. are:

- 1) Initially compute the in-degrees of all the vertices in the graph.
- 2) Find U as a vertex with degree 0 and store it in the list for ordering.
- 3) At this stage if no such vertex is detected then there is a cycle found and the algorithm stops.
- 4) Remove the vertex U and all edges (U, v) it belongs from the graph.
- 5) At this step, need to update in-degree of the vertices remain.
- 6) Repeat step-2 thro. 4 till the vertices present for processing.

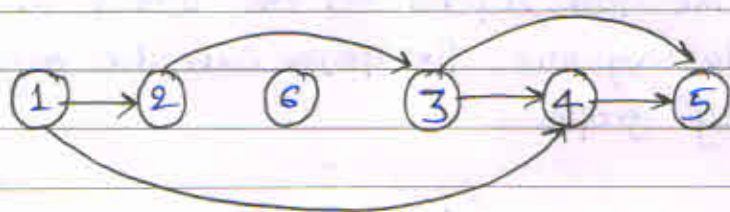
Ex. Topological Sorting :

Consider. given directed graph $G(V, E)$ find linear ordering vertices such that: for all edges (v, w) in E , v precedes w in the ordering.

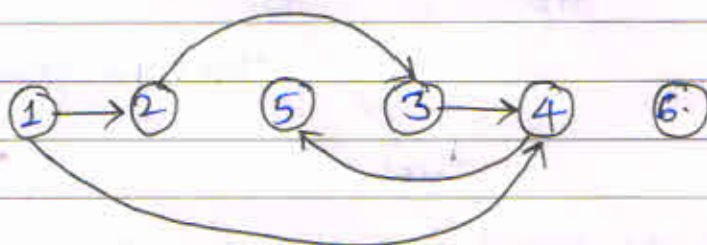


solⁿ :

→ For topological sorting any linear ordering in which all ~~are~~ arrows go to the right is a valid solution.



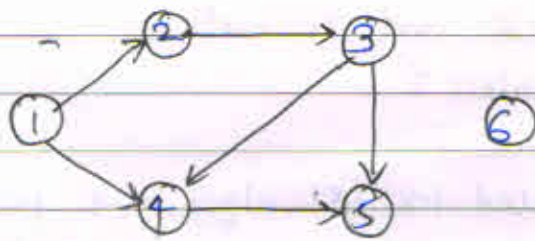
Valid topological sorted graph.



Invalid topological sorted graph

Topological sort Algo :

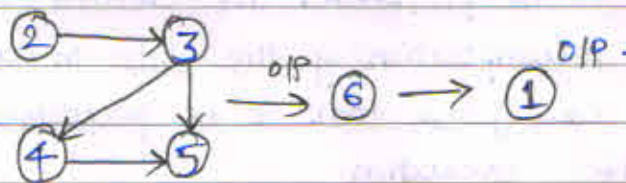
- 1) Identify vertices that have no incoming edges (select 1 vertex)
- 2) Delete the vertex of in-degree 0 and all its outgoing edges from the graph, place it in the o/p.
- 3) Repeat step 1 and step-2, until graph is empty.



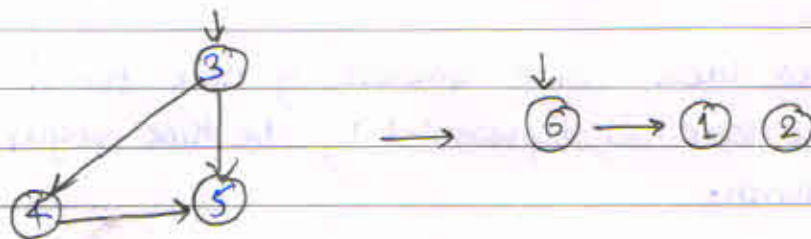
Step 1: Identify vertices that have no incoming edges



Step 2: Delete this vertex from graph.



Repeat step 1 & step-2 until graph is empty.



Repeat step 1 & 2 until graph is empty.



Repeat step 1 & 2 until graph is empty.



finally the topological sorting is shown here

* Characteristics of tasks and Interaction:

1) Task characteristics:

→ • Generation of Task:

The tasks included in the algo. are possible to generate either on prior basis termed as static/dynamic.

* Static task generation:

The tasks are defined before starting the exe. of the algo. and specified order is to be followed by the algo. in exe.
e.g. matrix multiplication.

* Dynamic task generation:

- The tasks to be performed are created dynamically based upon the decomposition of the data in certain situation.

In these cases, the tasks to be performed are not available before algo. execution.

e.g. recursive & Exploratory algo.

* Size of task:

- Each task takes some amount of time for its completion. The size of task is represented by the time required for its completion.

* Knowledge of task sizes:

- The task size knowledge is used for the choice of mapping. The tasks are mapped into processes and the prior knowledge of the task is used in mapping.

e.g. where dynamic decisions are required during processing.

* Data sizes:

- The required data for performing a task should be made available when mapping it into processes.

- The overhead associated with the data movement can be reduced when the size of data as well as its memory location are available.

* Mapping Techniques†

- The mapping techⁿ is used to map tasks into processes so that 11e1 exe. can be performed.
- The overall computation associated in problem to be solved. divided into nos of task.
- These task are mapped onto processes with the goal of completing exe. in the minimum amount of time.

There are 2 techniques†

- 1) static mapping
- 2) Dynamic mapping.

1) static Mapping -

In this techⁿ, the mapping of task onto process is performed before exe. of algo. This indicates that tasks are distributed among available processes prior to exe. of algo.

- scheme for static mapping†

- There are 3 diff scheme
- 1) Mapping Based on Data partitioning,
 - 2) Task graph partitioning.
 - 3) Hybrid strategies.

2) Dynamic mapping†

- In this techⁿ, the mapping of task onto the processes is performed during the exe. of the algorithm.

This indicate, the tasks are distributed among available processes when algo. actually executes.

- There are basic situation where dynamic mapping is applied, The 1st case is if tasks are generated dynamically then this mapping techniques is used.
- Dynamic techⁿ is applied in situation where large amount of data is associated with tasks. In this situation, dynamic mapping moves data among available processes.

There are 2 diff. classes of dynamic mapping.

- 1) Centralized dynamic mapping
- 2) Distributed dynamic mapping.

* parallel Algorithm model :-

" It is used to describe the strategy for partitioning data and how these data are processed -

A model is used to provide appropriate structuring of 1st algo. on the basis of 2 tech.

- 1) sele² of partitioning & mapping tech
- 2) proper use of strategy for interaction minimization.

- There are basically six diff. types of model.

- 1) Data 1st model
- 2) The task graph model
- 3) The task pool model
- 4) Master/slave model.
- 5) pipeline / producer-consumer problem
- 6) Hybrid model.

1) Data 1st model :-

- In this, once the data have been distributed 1st operations can be performed. In this way, the op^{er} performed by each task are similar but the data on which these op^{er} works are different.
- The overall task to be performed may be divided into diff phases. The data upon which the tasks work in diff phases may be diff.
- The 2 popular paradigm share-address space and msg-passing. can be used for implem² of data-1st algo.
- The 1st matrix multiplication is example of data-1st computation.

2) The task graph model :-

- As we have already studied task dependency graph is used to describe computation in any 1st algo.
- The task-dependency graph explicitly used in mapping of some algo.
- The data interaction cost among the tasks can be reduced by the use of interrelationship among the task in task dependency graph.
- The data movement cost is optimized by the use of mapping of task to process statically.
- The 1stism described with task dependency graph where each task is an independent task is called as task 1stism.
- e.g. parallel quick-sort algo.

3] Task pool model:

- The task pool model is also called as work pool model.
- it uses dynamic mapping approach for task assignments.
- In this, centralized or decentralized mapping can be adopted.
- The task-pool model can be used in msg-passing paradigm when data associated with the tasks is smaller than computation associated with the tasks.
- e.g. loop parallelization using chunk scheduling, tree search etc.

4] Master-slave model:

- In this, one processor is designated as master & others as slave and master is responsible to coordinate the activities among all slave processes.
- In master, it generates work to be performed & assign it to the nos of workers/slaves processes.
- The allocⁿ of task is depend upon size of task.
- In this, if work allocated to slave, then it completed by slave & slave is allocated to other jobs. it performs diff. task.
- master is responsible to synchronize the activities of the slaves after each phase.
- The master-slave model is generally suitable for shared address space & msg passing paradigm.

5] The pipeline or producer-consumer model:

- This model is based on the passing on stream of data thro. processes arranged in a succession.
- The data flow thro. successive processes & at that each process does some opⁿ on it. This is called stream fltism.
- Whenever, a new data arrives a new task execution initiates by the process in the pipeline.
- This model is called as producer-consumer model.
- bcz pipeline act as chain of producer & consumers.

* parallel programming with CUDA.

- It is based on the coverage of all diff. activities required to perform while handling HPC computing using CUDA.

Processor Architecture:

- The cuda architecture is used for construct^g the CUDA capable GPUs. The appl^y written can be executed on any card which support this architecture.

Each GPU has slightly diff. features bcz of diff. specification.

- When kernel is invoked, each thread block executes on a "multi-processor". This multiprocessor contains the resources to support certain num of threads.
- Each multiprocessor consist of:
 - 1) scalar processor cores
 - 2) 2 special fp^u unit for transcendental
 - 3) 1 multithreaded instr^y set
 - 4) On chip shared memory.
- One/more thread blocks are assigned to multiprocessor. during exe- g a kernel. The cuda runtime handles the dynamic scheduling of thread block on a group of multiprocessors.
- The scheduler will only assign thread block to multiprocessor when enough resources are available to support the thread block.
- Each block is split into SIMD (single instr^y multiple data). group of threads called "warps".
- The SIMD unit creates, manages, schedules and executes 32-threads simultaneously to create a warp.
- Every warp is synchronous, and therefore care must be taken to ensure that certain threads within warp do not take abnormally longer. compared to other threads in that same warp. , bcz warp will only execute as fast as the slowest thread.

* Interconnect :

The multiple GPU systems are designed bcz of foll. 2 reasons :

1) Problem Domain size :

In the case, where data sets are too large and ~~not~~ unable to fit into the memory of a single GPU.

2) Throughput and efficiency :

The throughput of an app. can be increased by processing it in the multiple GPU concurrently.

- The inter-GPU comm. is required to design properly while converting an app. into multiple GPU system.

- The efficiency of data transfer among the GPU depends on how GPUs are connected together within node or across a cluster.

There are 2 types of connectivity in multi-GPU systems:

* Multiple GPU connected over the PCIe bus in a single node.

* Multiple GPU connected over the ~~PCI~~ n/w switch in a cluster.

Fig. shows topology for cluster with two computing nodes.



- GPU1 and GPU2 are connected via PCIe Bus on node 1.

- GPU3 and GPU4 are connected via PCIe Bus on node 2.

- These two nodes (node 1 & node 2) are connected to each other through InfiniBand Switch.

* Communication:-

- The commⁿ in platform is in term of peer-to-peer commⁿ.
The CUDA peer-to-peer API is used to enable direct inter-device commⁿ.
- The CUDA 4.0 or higher version is required for peer-to-peer commⁿ. The CUDA P2P API supports two modes that allow direct commⁿ betⁿ GPUs.

* peer to peer access:-

Directly load & store addresses within a CUDA kernel and across GPUs.

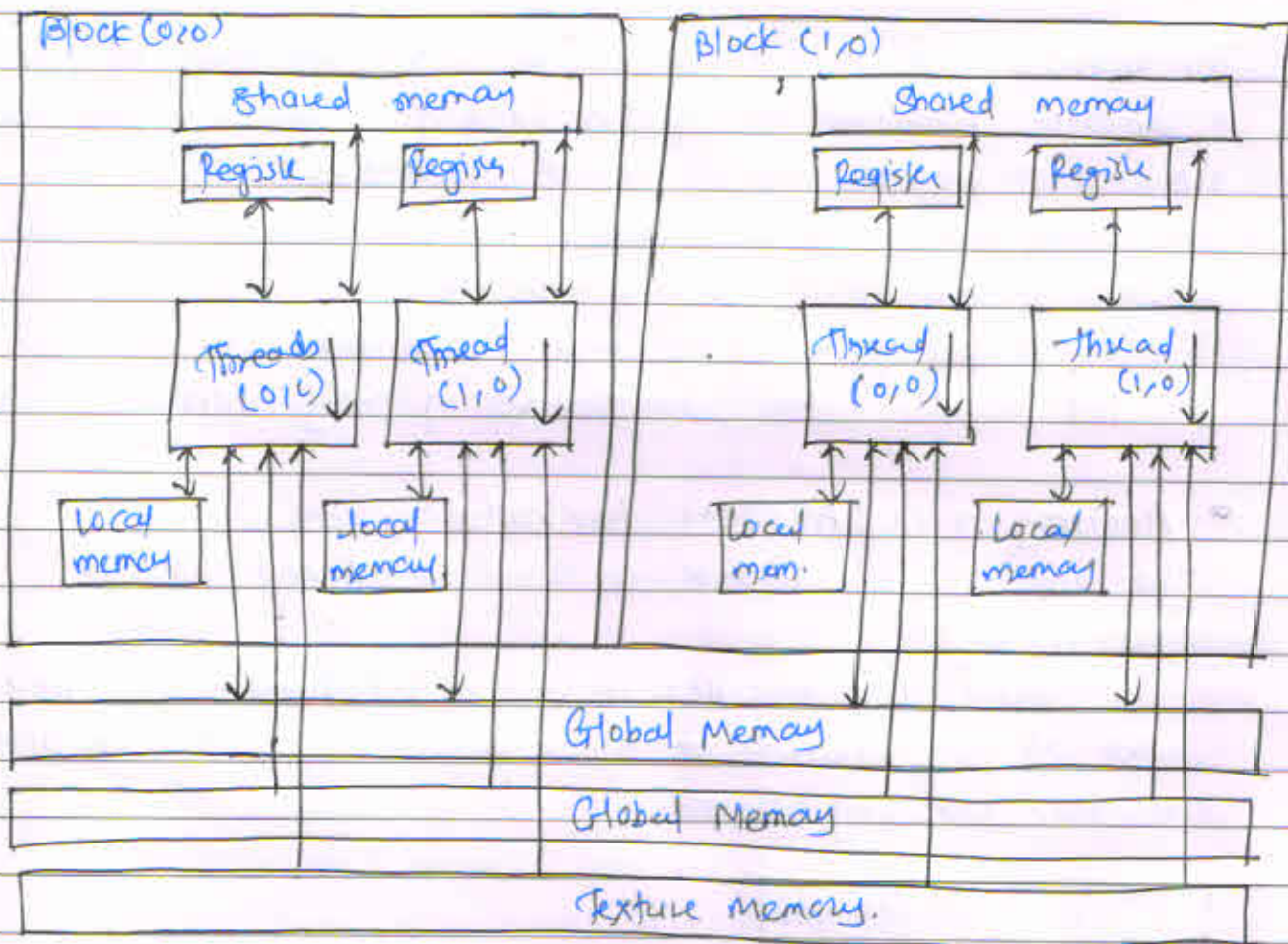
* peer to peer transfer:-

Directly copy data betⁿ GPUs.

- The direct peer to peer access is not supported when two GPUs are connected to diff root nodes within a system.

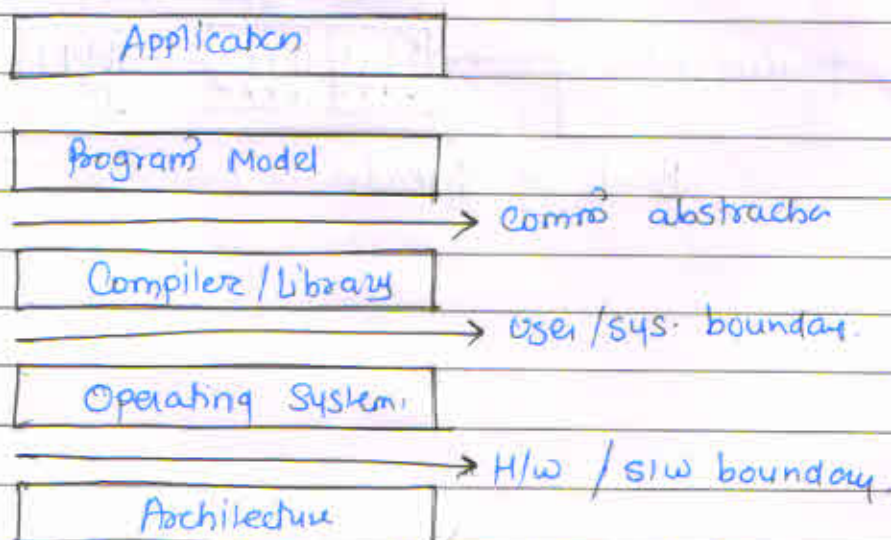
* Memory Organizations:-

- CUDA uses segmented memory architecture that allow appl^s to access data in global, local, shared, constant and texture memory.
- There are several level of memory & memory orgⁿ are arranged in the hierarchical fashion and each level has distinct Read and write characteristics.
- Every primitive thread has access to private 'local' memory as well as registers. Every thread in thread block has access to unified 'shared memory', shared among all threads for the life of that thread block. Finally all thread have read/write access to global memory.
- There also exist a read-only 'constant' & texture memory in the same locⁿ as global memory.
- Global memory is not cached, thro. mem. transactions may be 'coalesced' to hide the high memory access latency.
- The read only constant memory resides in the same location as global memory. this mem. may be cached.
- The read only texture mem. resides in same locⁿ of global mem. and is also cached.



* CUDA programming models for High performance Computing architecture:

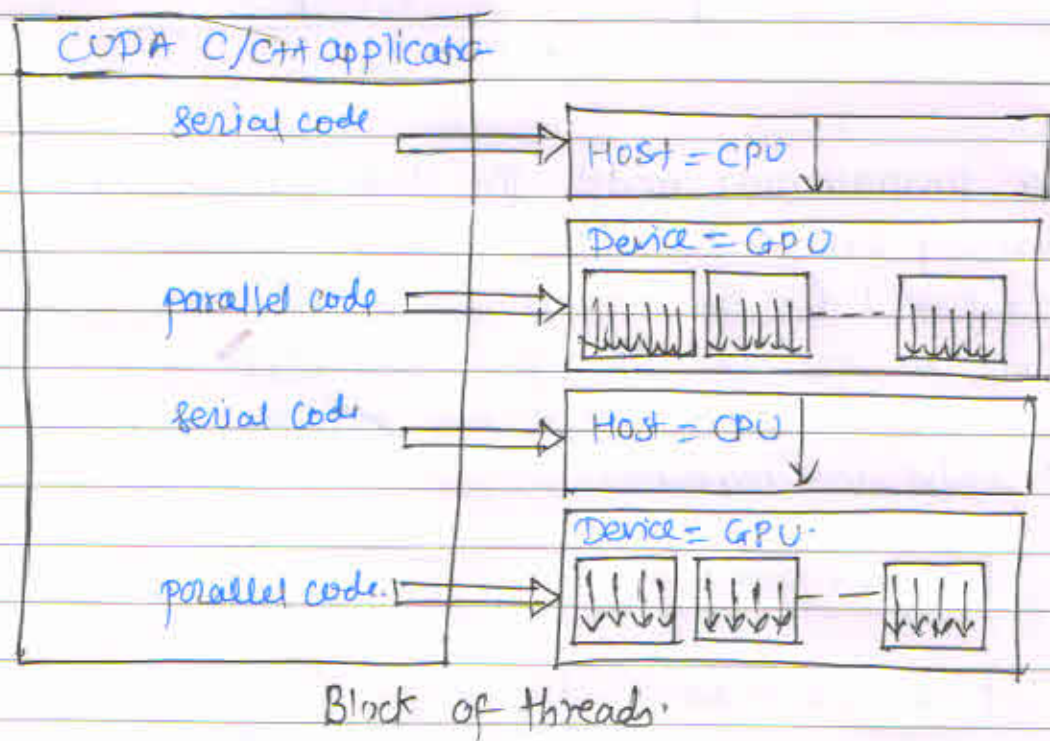
- The model used for progⁿ in particular archite. acts as a bridge betⁿ an applⁿ and its implemⁿ on available h/w.
- Fig. shows layer of abstraction lies betⁿ the prog. and progⁿ model implementation.



- The boundary betⁿ prog & progⁿ model impleⁿ is commⁿ abstraction. A compiler / libraries using sys. calls to access h/w services. are used to provide such a commⁿ abstraction.
- The components of the prog. share inforⁿ and coordinates their activities by the instrⁿ provided thro prog.
- The CUDA progⁿ model share many abstraction with other. H/W progⁿ models. it has some features of fully utilization of GPU.

• Anatomy of CUDA C/C++ Application:

- The applⁿ for the archite. contains serial and H/W code shown in fig.
- The serial code executes in a host (CPU) threads whereas the H/W code executes in many devices (GPU) threads across multiple processing elements.



* IBM Cell Broadband Engine (CBE) :-

→ The new processor created by IBM with the collaboration of Sony & Toshiba is called as cell.

It is combination of 64-bit dual threaded IBM PowerPC processor and eight nos of 'Synergistic Processing Elements (SPEs)'

- The IBM cell broadband engine (CBE) is basically heterogeneous multicore chip. It is combⁿ of 2 types of processors.

- The processor PPE (PowerPC Processing Element) and SPE are two diff processor combined in CBE.

- There are 8 nos of SPEs and one PPE included in CBE system.

- The other component of CBE are :-

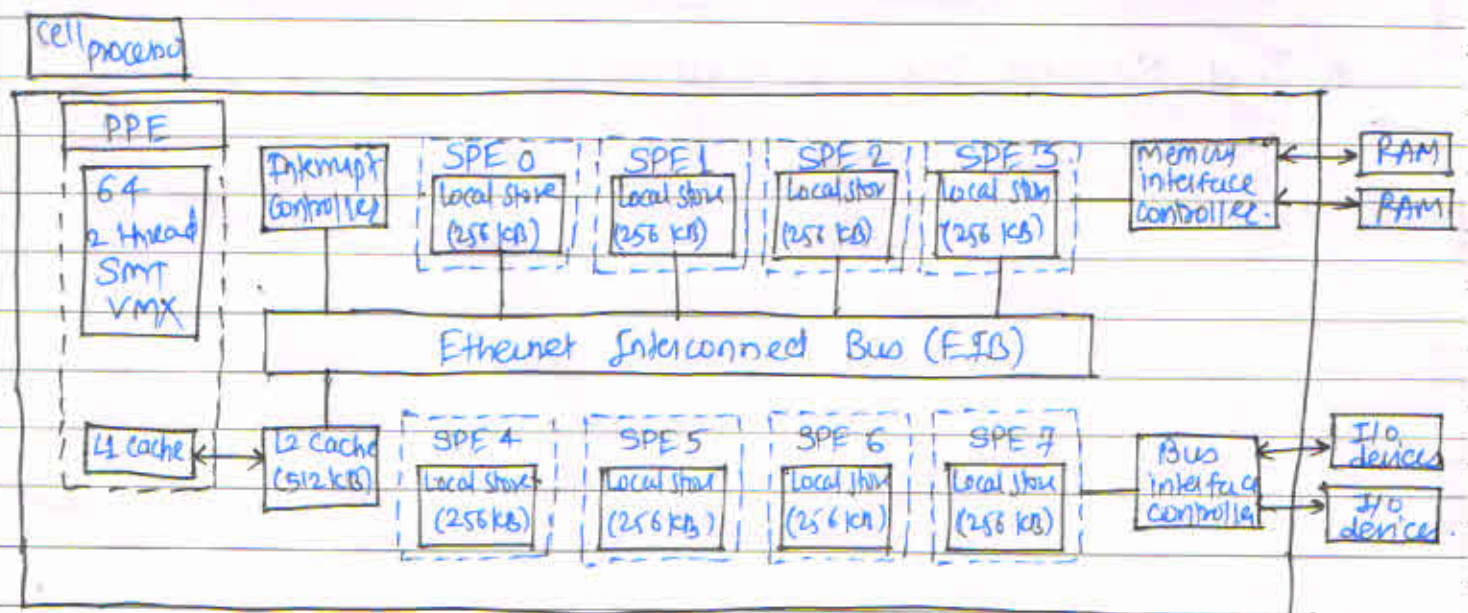
One high speed memory controller, one high-bandwidth element interconnect bus, high speed memory & I/O interfaces.

- All components along with 2 processor are integrated on a chip. In this way, it is considered as hybrid nine-core processor.

- The processor included in CBE have 2 diff. instrⁿ sets as compare to same instrⁿ set supported by most of the CPUs.

- The processor element in the CBE are optimized for certain types of operⁿ but these element can be used for general purpose computing also.

- The IBM CBE is commonly used as a processor for Sony's play station 3 system.



* Nvidia Tesla GPU :-

- - GPU (Graphics Processing Unit) is basically used for 3-D applications. A GPU is single chip processor used to provide lighting effects and transforms objects every time a 3-D scene is redrawn..
- The GPU is used to assist CPU on mathematically intensive tasks, require in performing graphics related activities. The CPU is freed from the overheads occurred from graphics related mathematical intensive tasks and the free cycles of CPU can be used for other tasks.
 - The GPU is used along with the CPU to accelerate scientific, analytics, engineering etc. applis. The overall computing is accelerated by the GPU and therefore this field of computing is called GPU accelerated Computing.
 - The basic difference betⁿ a CPU & GPU is the use of nos of cores included. The CPU contains less nos of cores for sequential serial processing whereas a GPU contains more than thousands of smaller and efficient cores for handling multiple tasks simultaneously.
 - Nvidia's brand name for the product, they develop for stream processing and general purpose GPUs is called as Nvidia Tesla.
 - The GPUs from G80 series onwards are utilized in the Nvidia Tesla Products.

* Intel Nehalem Micro-architecture :-

→

* Intel Nehalem Micro-architecture :

→ It is one of the successors of the intel's core architecture.

Intel core architecture uses multiple cores in single die for improvement of performance over single core architecture.

- The intel's Nehalem architec. provides improvements in core-to-core commⁿ by the use of point to point topology.

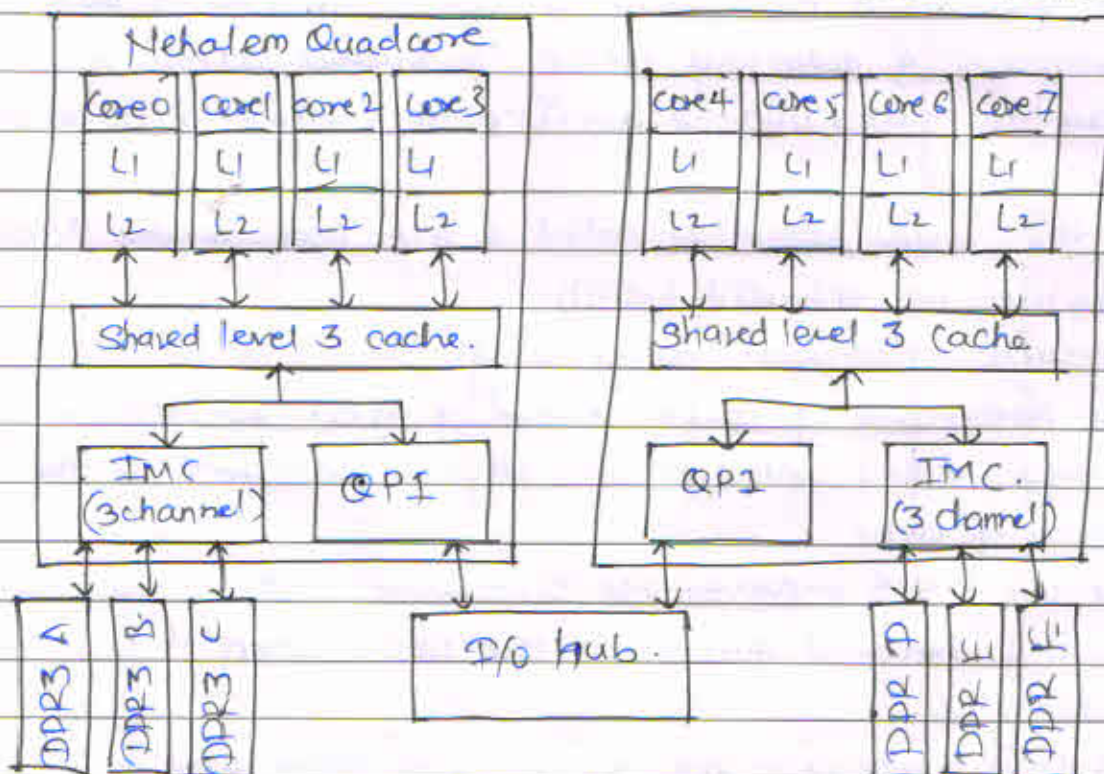
The use of point to point topology provides direct commⁿ ~~bet~~ among the microprocessor cores and the sys. mem. can be accessed in more diff. ways.

1) Architecture Approach in Nehalem -

- The nehalem archi. approach is much more modular as compare to the core archi. The modular archit. in nehalem makes it much more flexible & customizable to the appliⁿ.

- The main blocks in Nehalem includes microprocessor core.

(contain own L2 cache), a shared L3 cache, a bus controller, a graphics core and IMC (integrated Mem. controller)
QPI - quick path interconnect



- fig. shows e.g. of eight-core Nehalem processor in which two QPI bus controller included.

This configⁿ is used in the processor attached with intel Nehalem multiprocessor system.

2) Branch prediction:-

- The loop in code exe are detected by loop stream detector designed specifically for the core architecture.
- The loop included in code is detected during its exe & instr^s are stored in the special buffer so that it is not required to fetch it continuously from the cache. In this way the branch prediction success rate is increased for loops and overall performance improvement is achieved.

3) Out of order exe:-

- The purpose of out of order exe is to minimize pipeline efficiency by filling up the pipeline stalls with the useful instr^s.
- It was supported by core archite. also but in the Nehalem imple^d approach used to allow more instr^s to be ready for immediate exe.

4) Instr^s set:-

- It has seven new instr^s added by the intel to take advantages of data level parallelism. It support SIMD. facility.
- Due to SIMD, intensive appl^s (i.e. Multimedia) can be performed efficiently.
- Newly added instr^s is called as ATA (Appl^s targeted Accelerator) by intel bcz enhanced facility.

5) Transactor Lookaside buffer (TLB):-

- The performance of cache depend of TLB. The TLB is used to map virtual addresses into physical addresses in the memory or cache.
- The use of TLB improves the performance. In the sense that a page is accessed quickly in the cache when it is mapped in the TLB.

6) Cache & Coherency:- (Note: you can write theory about it)

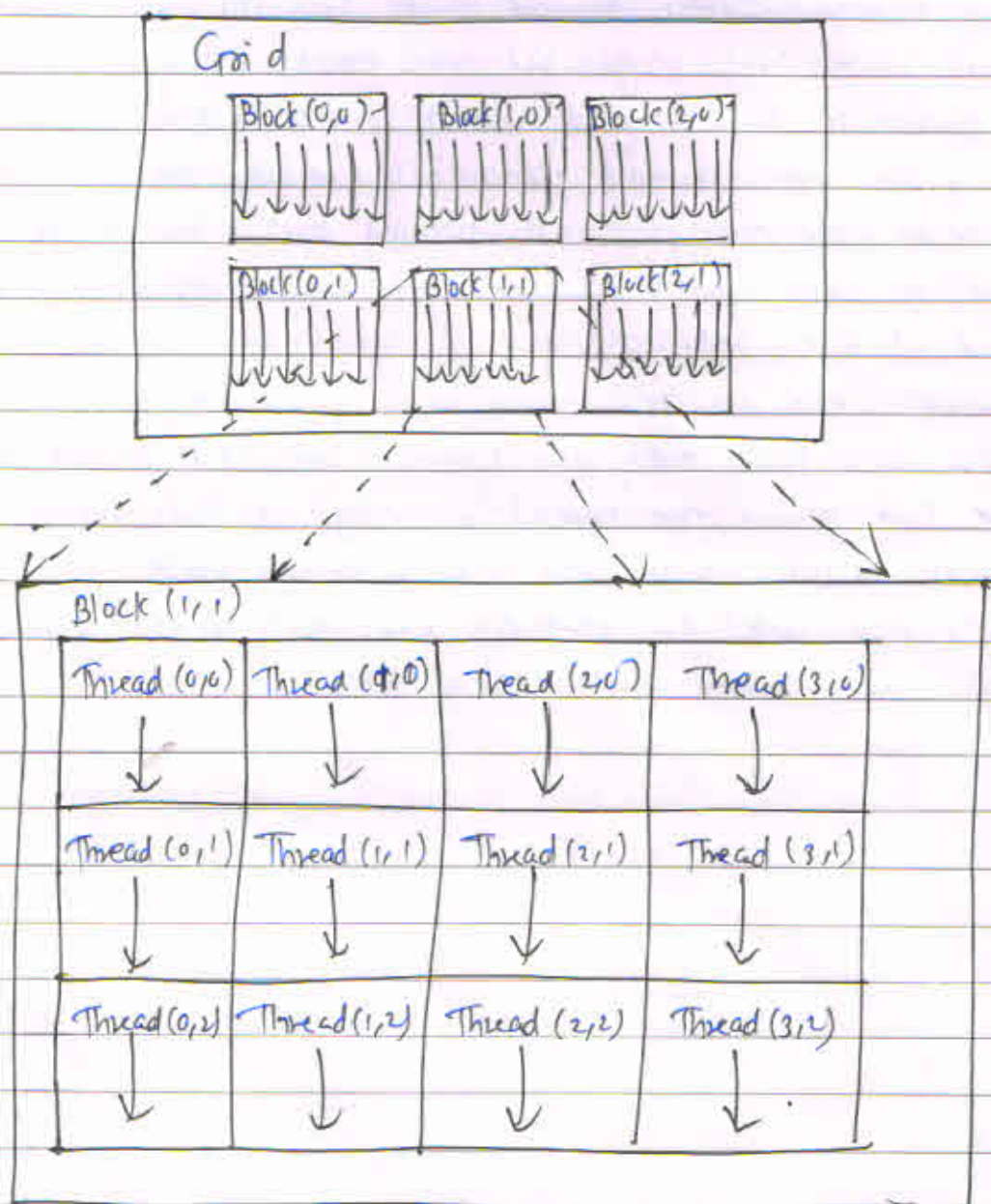
7) Memory Controller:-

- It is located at diff loc as compare to the core processor. Nehalem archi. integrates mem. controller to the processor or in the die instead of the off-chip on motherboard.
- due to on-chip mem. controller, it becomes totally free & independent of h/w.
- It provides improvement by running as fast on h/w platform.

* Thread Organization:

- The CUDA processing paradigm uses an approach is called kernel. A kernel is actually a sub-routine or mini-prog.
- The device like NVIDIA graphics card used inside the host sys. runs kernel.

The kernel prog. is executed simultaneously by no. of primitives thread.



- There are no. of batches of primitives threads organized into thread blocks. A thread block contains a specific no. of primitives threads, chosen based on the amount of available shared memory as well as mem. latency hiding tech characteristics desired.

- The max g thread in thread block is also limited by archi. to total g 512 thread per block. Each thread within thread block can communicate efficiently using shared memory scoped to each thread block.
Using this shared mem, all thread can also sync. within thread block.
- Every thread within thread block has its own thread ID. Thread blocks are organized into 1D, 2D or 3D array.
- As shown in fig. a grid which is collection g thread block g the same thread which all executes the same kernel. The thread block are physically limited to 512 threads per block. & bcz g that grids are used for computing a large no- g thread block in 111.
- Thread block in grid may not synchronize with one another bcz they may not communicate via shared memory.
- The diag shows the thread hierarchy. In this represent a given kernel contain 3×2 grid g thread block.
- There are total 72 threads executing in the said kernel where each thread block is a 4×3 .