

Lecture 11 – Parallel/Processor Domain Decomposition

- **Parallelizing a computer code requires the following steps:**

- Identifying portions of the work that can be performed concurrently
- Mapping the concurrent pieces of work into multiple processes (that run in parallel)
- Distributing the input, output, and intermediate data associated with the program
- Managing accesses to data shared by multiple processors
- Synchronizing the processors at various stages of the parallel program execution

Parallel/Processor Partitioning or Decomposition

MPI, PVM, or OpenMP

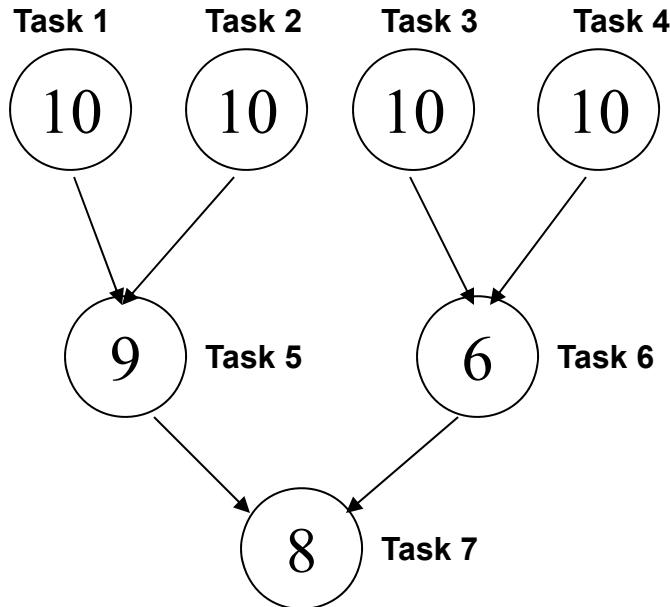
Granularity

- **Granularity:** The number and size of tasks into which a problem can be decomposed
 - A decomposition into a large number of small tasks is considered *fine-grained*
 - A decomposition into a small number of large tasks is considered *coarse-grained*
- **Degree of Granularity:**
 - The maximum number of tasks that can be executed simultaneously in a parallel program at any given time is the *maximum degree of concurrency*
 - The average number of tasks that can concurrently run over the entire duration of program execution is the *average degree of concurrency*

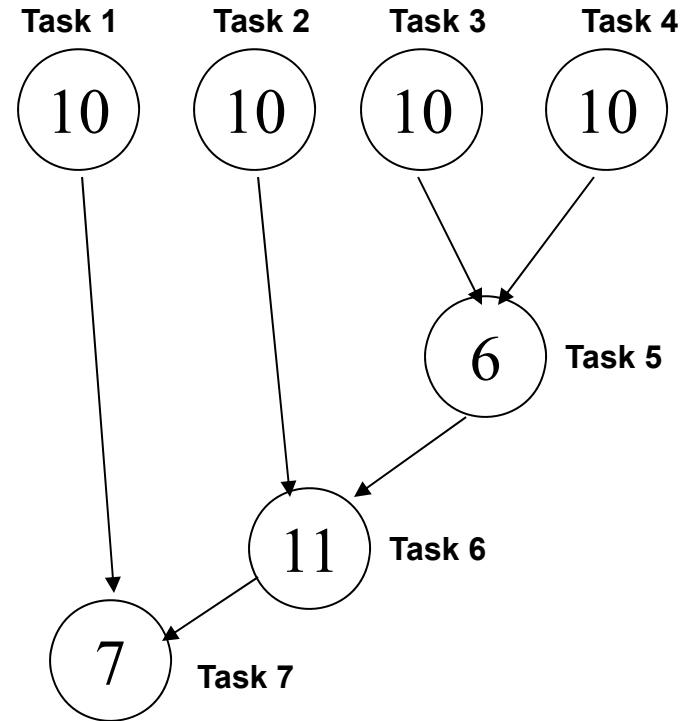
Critical Path

- We can create a *task-dependency graph* of the various tasks performed during program execution
 - A task-dependency graph is an abstraction (flow chart) used to show dependencies among program tasks and their relative order of execution
 - The initial tasks in the graph are called the *start* nodes
 - The final tasks in the graph are called the *finish* nodes
 - The longest path between any pair of start and finish nodes is the *critical path*
 - The sum of the weights (amount of computational work) of the nodes (tasks) along the critical path is known as the *critical path length*
 - The ratio of the total amount of work to the critical path length is the *average degree of concurrency*

Examples



- **Critical Path = 2 (4-6-7) or (1-9-8)**
- **Critical Path Length = 27 (1-9-8)**
- **Total Amount of Work = 63**
- **Average Degree on Concurrency = $63/27 = 2.33$**



- **Critical Path = 3 (4-5-6-7)**
- **Critical Path Length = 34**
- **Total Amount of Work = 64**
- **Average Degree on Concurrency = $64/34 = 1.88$**

Granularity

- At first, it appears that the time required to solve a problem can be reduced by increasing the granularity of decomposition and perform more tasks in parallel.
- This is not the case in general, however.
- There is an inherent bound on how fine-grained a decomposition of a problem permits (due to the nature of the problem, for example: the number of intervals, number of data points, number of cells, number of blocks, etc.)
- Another major factor that limits speedup from further decomposition is the *interaction* between tasks (processors)

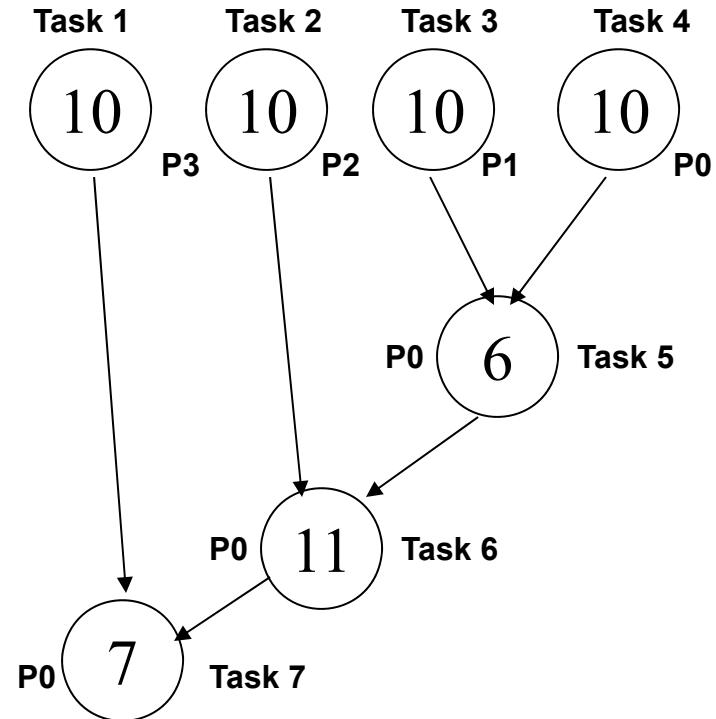
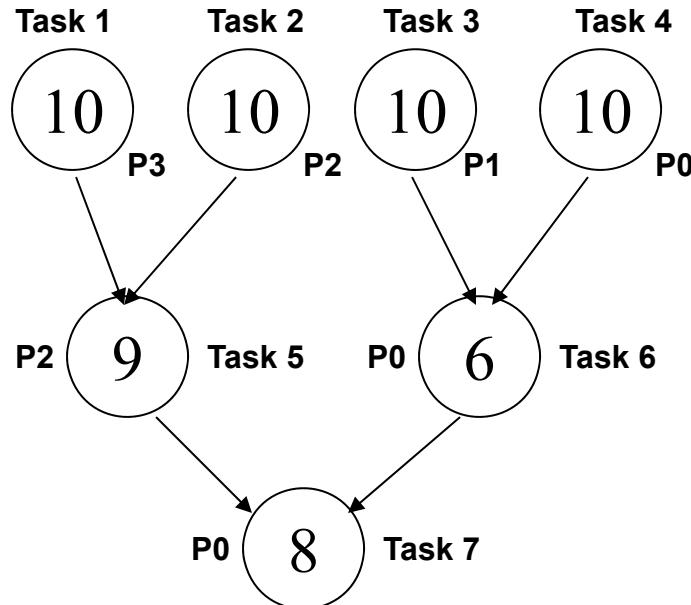
Interaction

- **The pattern of interaction among the tasks can be described in a *task-interaction graph***
 - The nodes in a task-interaction graph are the tasks
 - The edges connecting the nodes represent the interactions. The direction of the edge-lines represent the flow of data (from one task to another).
 - The nodes and edges can be assigned weights based upon the amount of work performed.
- **The task-interaction graph looks much the same as the task-dependency graph but with weights assigned to the interactions (edges).**

Mapping

- A **process** is a computing agent that performs tasks (much like a processor)
- The procedure by which tasks are assigned to processes for execution is called ***mapping***
- A good mapping should
 - Maximize the use of concurrency by mapping independent tasks onto different processes
 - Minimize the total completion time by ensuring that processes are available to execute tasks on the critical path as soon as possible
 - Minimize the interaction among processes by mapping tasks with a high degree of mutual interaction onto the same process

Examples: Mapping 7 Tasks onto 4 Processes



- A maximum of 4 processes can be used efficiently. The maximum degree of concurrency is 4
- The last 3 tasks can be mapped arbitrarily to satisfy the constraints of the task-dependency graph
- Makes more sense to map the tasks connected by an interaction edge onto the same process since this prevents an inter-task interaction from becoming an inter-processes interaction

Decomposition Techniques

- **There are different techniques for decomposition depending on the type of problem that is being solved**
 - Recursive decomposition
 - Data-decomposition
 - Exploratory decomposition
 - Speculative decomposition
- **Recursive and data-decomposition are the most general**
- **Exploratory and Speculative are for special-purpose problems**

Recursive Decomposition

- **Recursive decomposition used for problems that can be solved with a divide-and-conquer strategy**
 - Example: Sorting
 - First, divide the problem into a set of independent sub-problems
 - Each of these sub-problems is solved by sub-dividing again by another number of sub-problems, and so on...

Data Decomposition

- **Data decomposition is the most general way to break a problem into sub-problems**
 - The data on which the computations are performed is partitioned
 - The data partitioning is used to induce a partitioning of the computations in the tasks
 - For instance, break the data (intervals, cells, nodes, blocks, etc.) into partitions and perform the tasks of a computation on the partitions of the data. This is the technique typically used in simulations and engineering computations and will be what you use in Projects-4 and 5.
 - Data decomposition can be performed based upon:
 - Output data
 - Input data
 - Output and input data

Exploratory Decomposition

- **Exploratory decomposition is used for problems whose underlying computations correspond to a search of a space for solutions.**
 - Partition the search space into smaller parts
 - Search each one of these parts concurrently until the desired solutions are found
 - Optimization or search problems could fall into this category

Speculative Decomposition

- **Speculative decomposition is used when a program may take on one of many possible computational branches depending on the output of other computations that precede it.**
 - While one task is performing the computation whose output is used in deciding the next computation, other tasks can concurrently start the computations of the next stage
 - Example: switches (branches) of logic
 - While one task is performing the computation that will eventually resolve the switch, other tasks could pick up the multiple branches of the switch in parallel
 - When the input for the switch has finally been computed, the computation corresponding to the correct branch would be used while that corresponding to the other branches would be discarded

Task Characteristics and Interaction, Mapping of Tasks to Processes

- **So far we have discussed the different type of decomposition techniques**
 - Recursive decomposition
 - Data-decomposition
 - Exploratory decomposition
 - Speculative decomposition
- **Now let's discuss ways of characterizing the tasks, mapping those tasks to processes, and load-balancing the processes**

Task Characteristics

- **The following characteristics of the tasks have a large influence on the suitability of the mapping scheme:**
 - Task generation
 - **Static task generation**: all the tasks are known before the algorithm is executed
 - **Dynamic task generation**: the actual tasks are not known a priori (the tasks change with the computation)
 - Task sizes
 - The relative amount of time required to complete it. Depends on the uniformity of the various tasks (i.e. are they all the same)
 - Knowledge of task sizes
 - Size of data associated with tasks
 - The size of the data for a task is very important since it determines the load or weight of that task. Load-balancing of processes will depend on “averaging” the weights of tasks across the processes

Interaction Characteristics

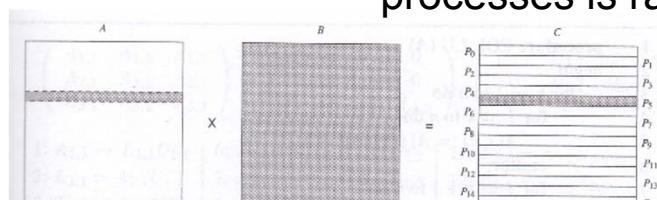
- **Static vs Dynamic**
 - **Static**: The interactions happen at predetermined times and the stage of the computation at which each interaction occurs is known (your projects fall into this category)
 - **Dynamic**: The timing of interactions or the set of tasks to interact with cannot be determined prior to execution
- **Regular vs Irregular**
 - **Regular**: The interaction pattern has some structure that can be exploited for efficient implementation (your projects fall into this category)
 - **Irregular**: There are no regular interaction patterns (eg. sparse matrix-vector multiplication)
- **Read-only vs Read-Write**
 - **Read-only**: Tasks only require a read-access to the data shared among many concurrent tasks
 - **Read-write**: Multiple tasks need to read and write on some shared data (your projects fall into this category)
- **One-way vs Two-way**
 - **One-way**: Only one of a pair of communicating tasks initiates an interaction and completes it without interrupting the other one
 - **Two-way**: The data or work needed by a task or a subset of tasks is explicitly supplied by another task or subset of tasks (your projects fall into this category)

Mapping

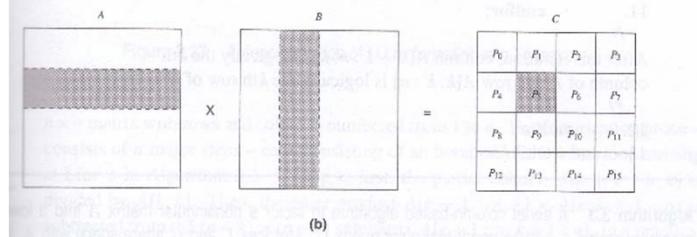
- **An efficient mapping of tasks onto processes must strive to**
 - Reduce the amount of time processes spend in interacting with each other
 - Reducing the total amount of time some processes are idle while the others are engaged in performing tasks
- **This can be difficult**
 - These two objectives can often be in conflict with each other
 - Tasks resulting from a decomposition may not be all ready for execution at the same time. There may be a task dependency.
 - Poor synchronization among interacting tasks can lead to process idling

Mapping

- **Static Mapping: Distribute the tasks among processes prior to execution (what your projects do)**
 - Data Partitioning: Tasks are closely associated with portions of the data by the “owner computes rule”.
 - Array Distribution
 - Block distributions: Distribute an array and assign uniform contiguous portions of the array to different processes.
 - Cyclic and block-cyclic distributions: Partition an array into many more blocks than the number of available processes. Then assign the partitions and the associated tasks to processes in a round-robin manner so that each process gets several non-adjacent blocks. Randomized block distributions are similar except that the assignment of the partitions to the processes is random.

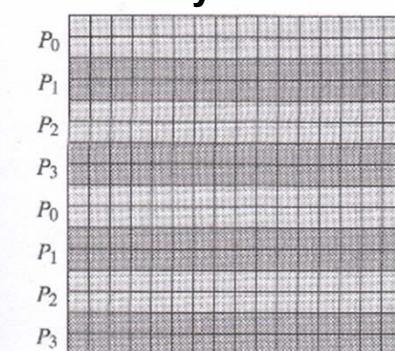


Matrix multiplication partitioning



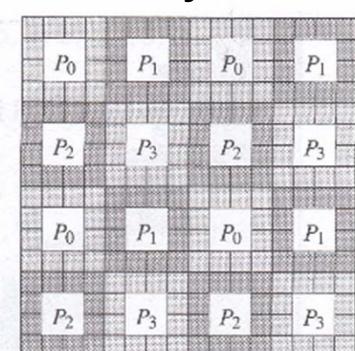
(b)

Cyclic



(a)

Block-cyclic



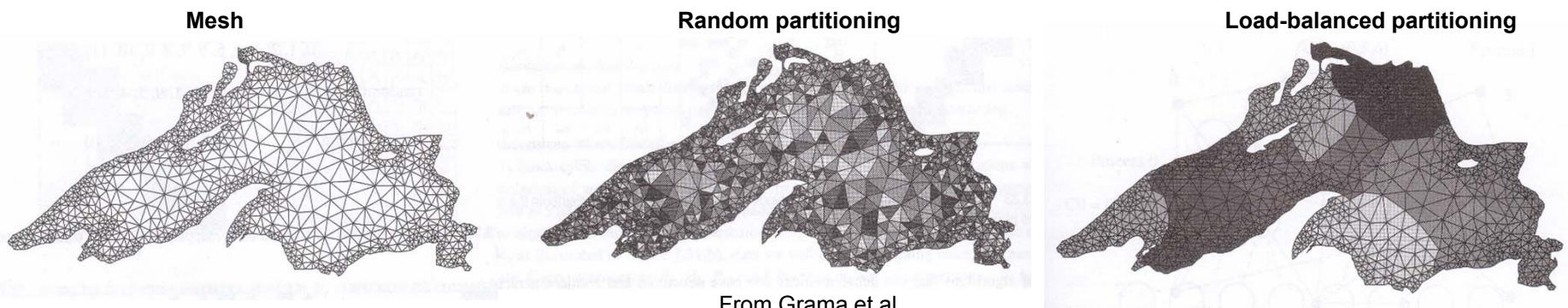
(b)

From Grama et al.

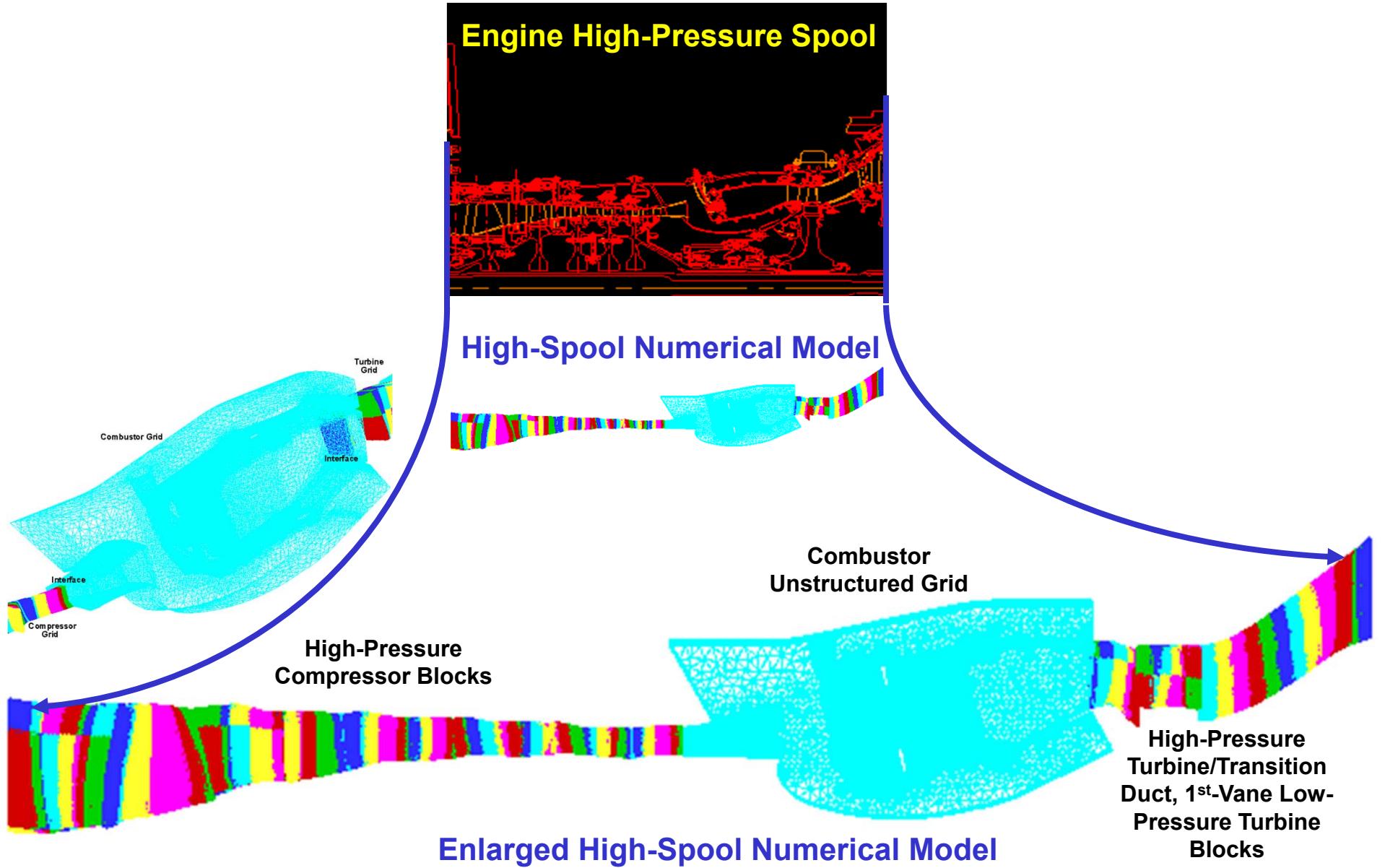
15

Static Data Partitioning

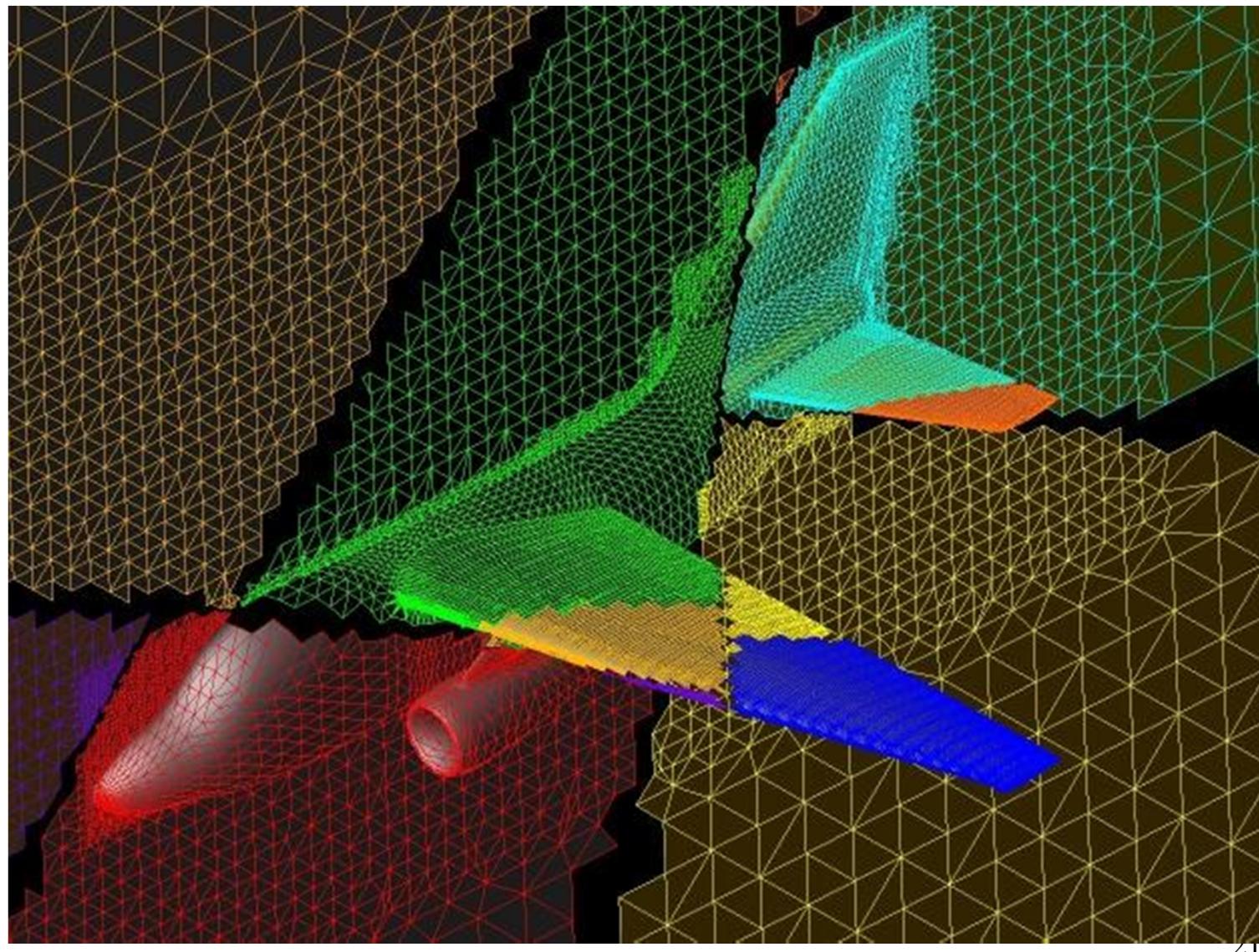
- Graph Partitioning: Many algorithms operate on sparse data with irregular interaction among the data elements.
 - Many engineering simulations of physical phenomena fall into this category. Here, the amount of computation at each mesh point (cell) is about the same.
 - Theoretically, this problem can be easily load balanced by simply assigning the same number of mesh points (cells) to each process.
 - However, a high interaction overhead may occur if the distribution of mesh points (cells) does not strive to keep nearby mesh points together.
 - Generally need to partition the mesh into p parts such that
 - each part contains roughly the same number of mesh points (nodes)
 - The number of edges that cross partition boundaries (i.e. those edges that connect points belonging to two different partitions) is minimized



Example of Static Data (Structured) Partitioning

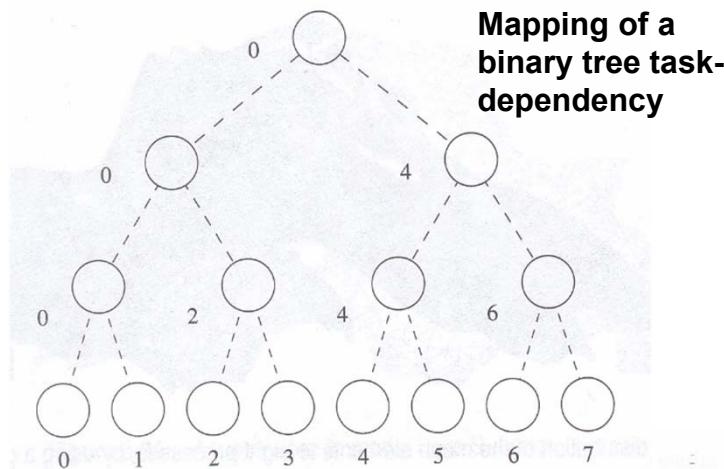


Example of Static Data (Unstructured) Graph Partitioning



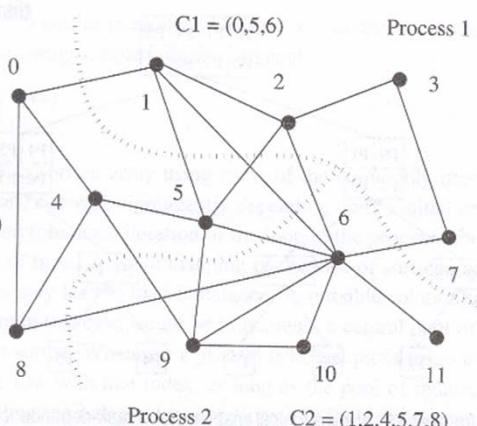
Static Data Partitioning

- **Task Partitioning:** Mapping based on partitioning a task-dependency graph where the nodes (tasks) are mapped onto processes

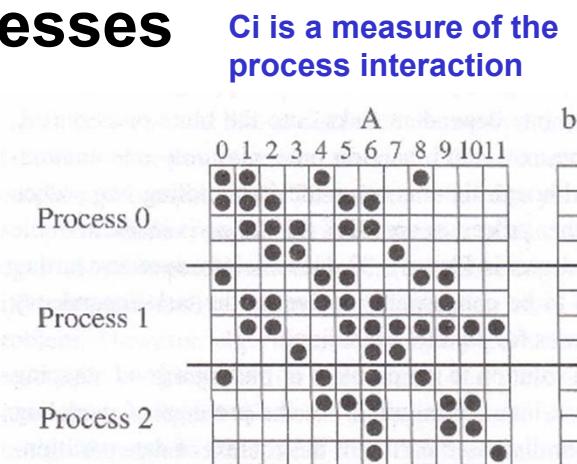


Look at any task of a row or column and look at what other nodes are in that row or column to create the task map. Then partition the map to load balance (ie make the number of Ci's the same)

Mapping of a sparse matrix-vector multiplication by partitioning the task-interaction graph



Ci contains the indices of b that process i needs to access from other processes → fewer interactions



Mapping of a sparse matrix-vector multiplication

$$C0 = (4,5,6,7,8)$$

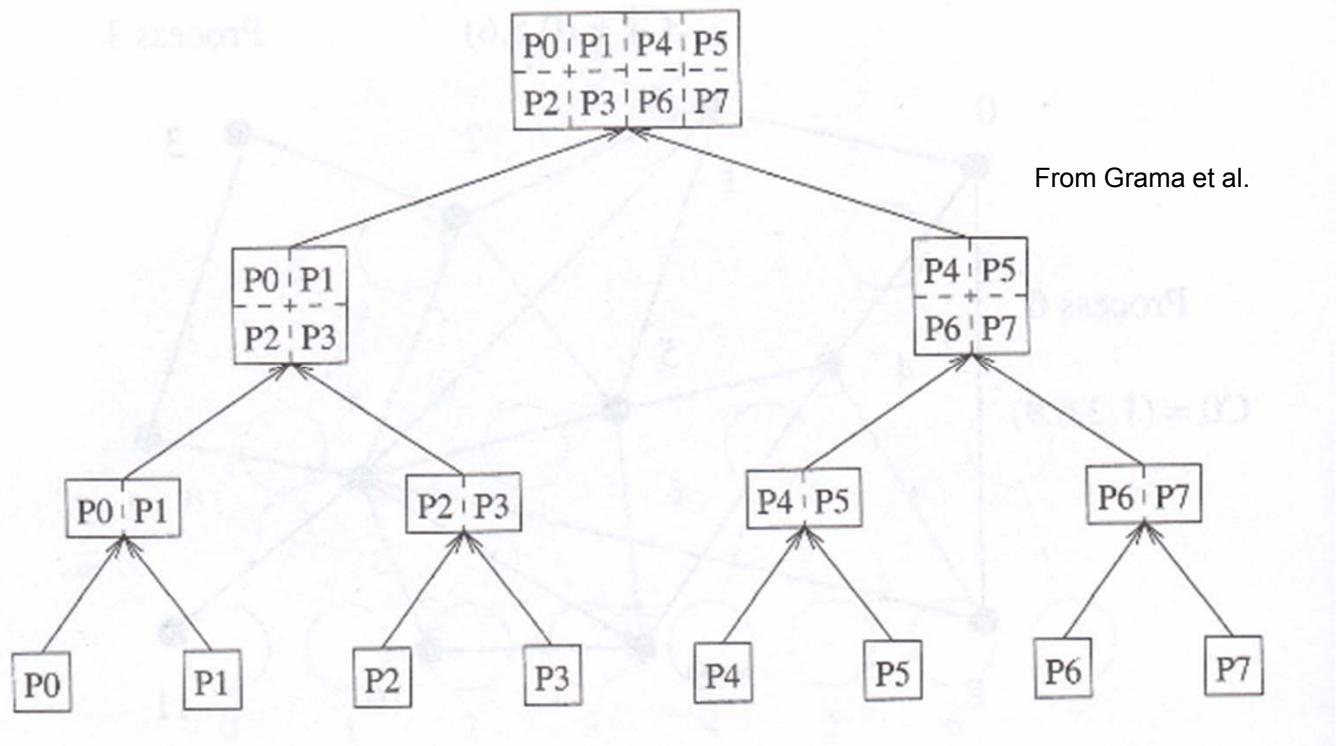
$$C1 = (0,1,2,3,8,9,10,11)$$

$$C2 = (0,4,5,6)$$

Ci contains the indices of b that process i needs to access from other processes (naïve)

Static Data Partitioning

- Hierarchical Mappings can remove some of the load balancing problems with task partitioning by sub-partitioning the task at each level.



Dynamic Mapping

- **Dynamic mapping of data to processors is usually required when**
 - Static mapping results in large imbalance
 - Task-dependency graph is highly dynamic
- **Example might include**
 - Grid-embedding adaptation (refinement)
 - Some optimization or sorting problems
- **Dynamic mapping classified as**
 - Centralized: all executable tasks are maintained in a common central data structure
 - Master-worker relation: When a process has no work, it takes a portion of available work from the central data structure of master
 - Generally easier to implement
 - Distributed: set of executable tasks are distributed among processes which exchange tasks at run time to load balance
 - Each process can send or receive work from any other process
 - Can be quite complex to keep track of process loads and redistribution

Methods for Containing Process Interaction Overhead

- **In general spatial and parallel/process domain decomposition, there are techniques used to minimize communication (overhead) costs between processes**
 - Maximize data locality
 - Minimize volume of data exchange
 - Minimize frequency of interactions
 - Minimize contention and “hot-spots”
 - Overlap computations with interactions as much as possible

Methods for Containing Process Interaction Overhead

- **Maximize data locality:**
 - Use techniques that promote the use of local data (cached) or data that has been recently fetched
 - Minimize the volume of non-local data that are accessed
 - Maximize the reuse of recently accessed data
 - Minimize the frequency of accesses
- **Minimize the volume of data exchange:**
 - Similar to maximizing the temporal data locality
 - Reduces the need to bring more data into local memory or cache
 - Proper spatial decomposition is important!
 - Use local data to store intermediate results and perform the shared data access to only place the final results (example: use of halos or accumulation operators)
- **Minimize the frequency of interactions**
 - Reduce the startup costs associated with interactions
 - Restructure the algorithm so that shared data are accessed and used in large pieces
 - Similar to increasing spatial locality of data access

Methods for Containing Process Interaction Overhead

- **Minimize contention and “hot spots”**
 - Contention occurs when multiple tasks try to access the same resources concurrently
 - Multiple simultaneous transmission of data over the same interconnection link
 - Multiple simultaneous accesses to the same memory block
 - Multiple processes sending messages to the same process at the same time
 - Contention can be reduced by redesigning the parallel algorithm to access data in contention-free patterns
- **Overlapping Computations with Interactions**
 - Processes often spend time waiting for shared data to arrive or to receive additional work
 - Can reduce by:
 - Initiate an interaction early enough so that it can complete before it is needed

Homework 5

- Go to the webpage:

<http://glaros.dtc.umn.edu/gkhome/views/metis/>

and read about the metis and parmetis static graph partitioning libraries

I have downloaded metis and parmetis. The zipped tar-files are located under the “Codes” directory on smartsite. These tar-files contain the code (metis5.1.0 and parmetis4.0.3), manual, makefiles, etc. to create the libraries on just about any platform.

Scan through the manuals of metis and parmetis to learn about their capability and how you might use them in the future