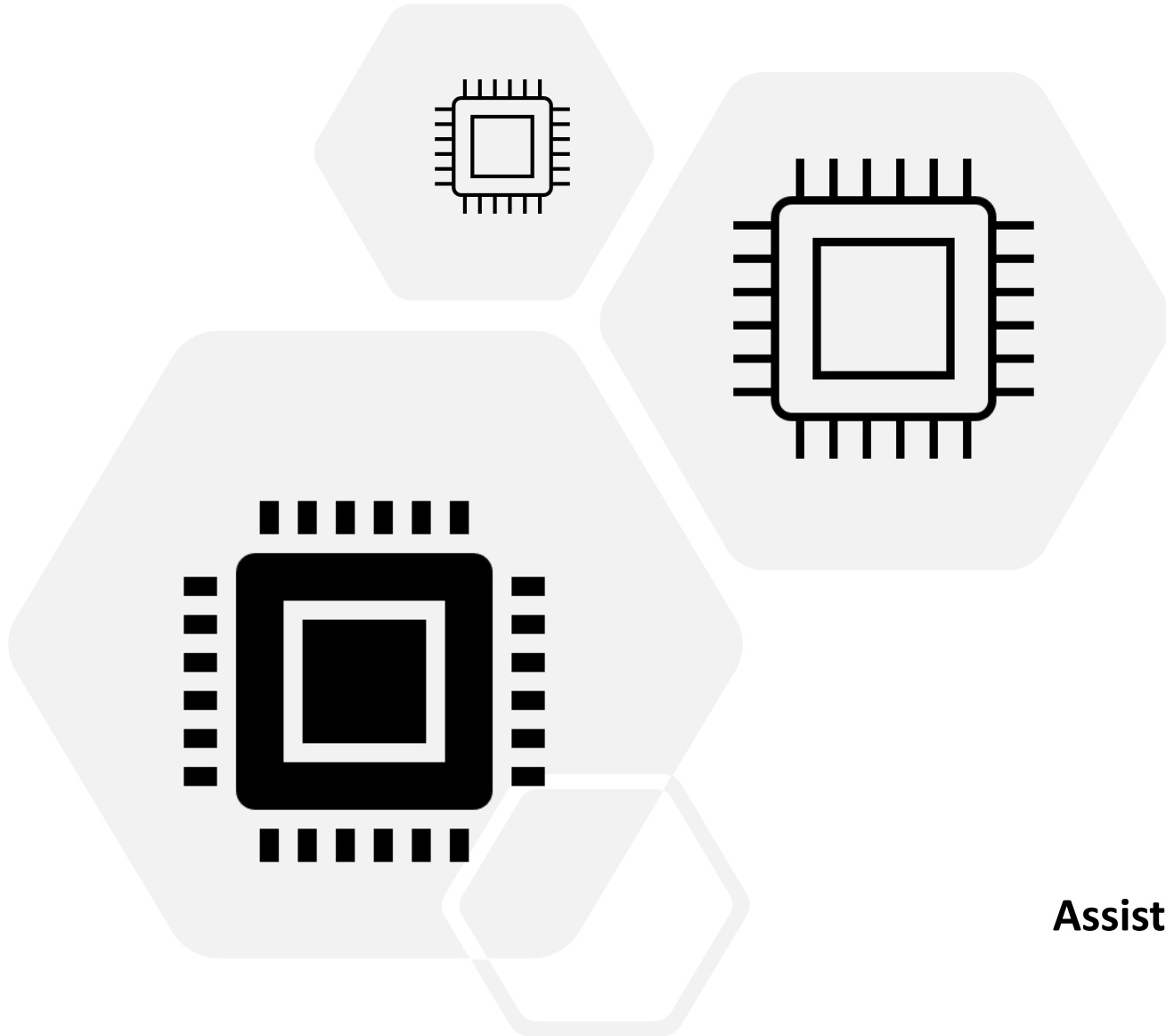


High Performance Computing

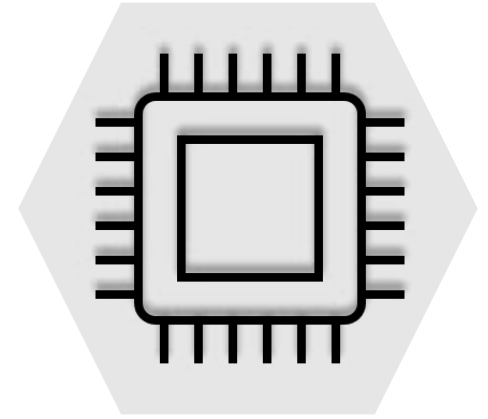


Module 3 - Parallel Algorithm and Concurrency

Shafaque Fatma Syed

Assistant Professor - Dept. of Information Technology
A P Shah Institute of Technology, Mumbai

Topics to be discussed



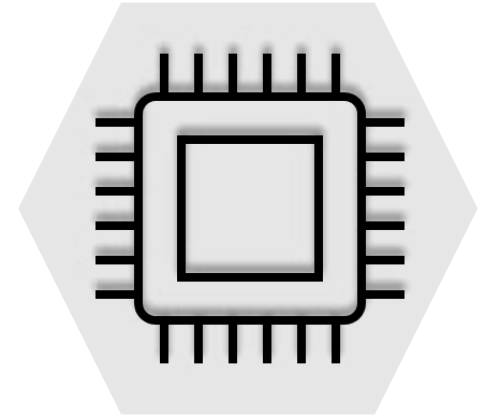
- **Introduction to Parallel Algorithms - Preliminaries**

- **Tasks and Decomposition**
- **Processes and Mapping**
- **Processes Versus Processors**

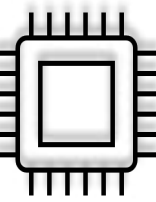
- **Decomposition Techniques**

- **Recursive Decomposition**
- **Data Decomposition**
- **Exploratory Decomposition**
- **Speculative Decomposition**
- **Hybrid Decomposition**

Topics to be discussed



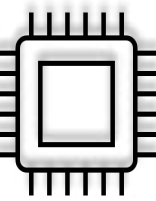
- **Characteristics of Tasks and Interactions**
 - Task Generation, Granularity, and Context
 - Characteristics of Task Interactions
- **Mapping Techniques for Load Balancing**
 - Static and Dynamic Mapping



Preliminaries: Decomposition, Tasks, and Dependency Graphs

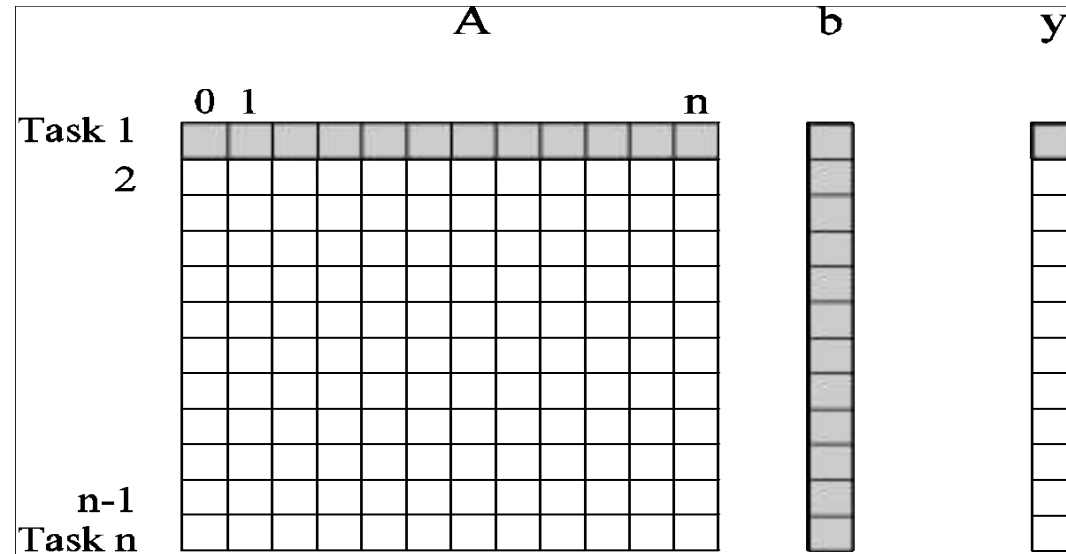
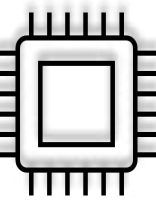
- **Algorithm development** is a critical component of problem solving using computers.
- Specifying a **parallel algorithm** may include some or all of the following:
 - **Identifying** portions of the work that can be performed concurrently.
 - **Mapping** the concurrent pieces of work onto multiple processes running 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.
- **Two key steps in the design** of parallel algorithms:
 - i) Dividing a computation into smaller computations
 - ii) Assigning them to different processors for parallel execution

Preliminaries: Decomposition, Tasks, and Dependency Graphs



- The **first step** in developing a parallel algorithm is to **decompose** the problem into tasks that can be executed concurrently.
- A given problem may be **decomposed into tasks** (of same, different, or even interminate sizes) in many different ways.
- The process of dividing a computation into smaller parts, some or all of which may potentially be executed in parallel, is called **decomposition**.
- **Tasks** are programmer-defined units of computation into which the main computation is subdivided by means of decomposition.
- A decomposition can be illustrated in the form of a **directed acyclic graph** with **nodes corresponding to tasks** and **edges indicating that the result of one task** is required for processing the next. Such a graph is called a **task dependency graph**.

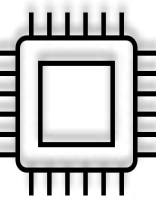
Example: Multiplying a Dense Matrix with a Vector



Computation of each element of output vector y is independent of other elements. Based on this, a dense matrix-vector product can be decomposed into n tasks. The figure highlights the portion of the matrix and vector accessed by Task 1.

Observations:

While tasks share data (namely, the vector b), they do not have **any control dependencies** - i.e., no task needs to wait for the (partial) completion of any other. All tasks are of the same size in terms of number of operations. *Is this the maximum number of tasks we could decompose this problem into?*



Example: Database Query Processing

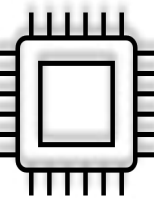
Consider the execution of the query:

**MODEL = ``CIVIC" AND YEAR = 2001 AND
(COLOR = ``GREEN" OR COLOR = ``WHITE)**

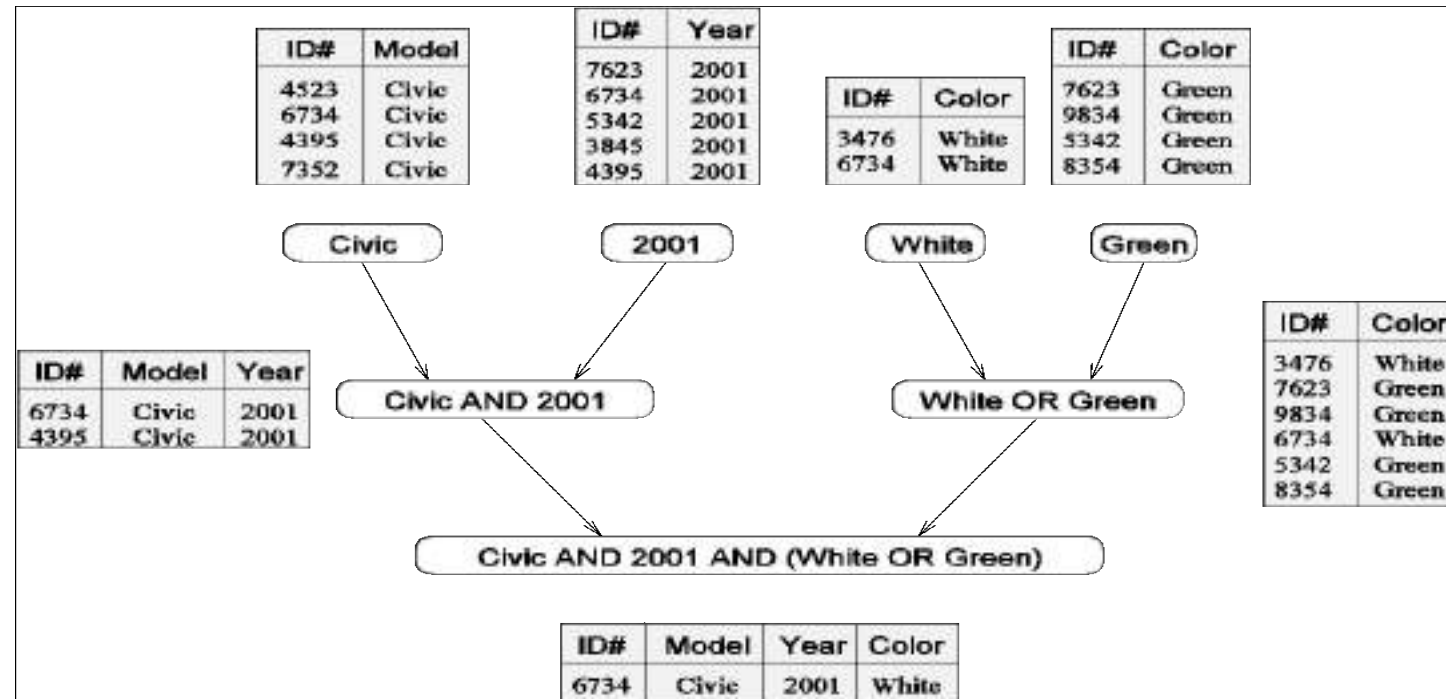
on the following database:

ID#	Model	Year	Color	Dealer	Price
4523	Civic	2002	Blue	MN	\$18,000
3476	Corolla	1999	White	IL	\$15,000
7623	Camry	2001	Green	NY	\$21,000
9834	Prius	2001	Green	CA	\$18,000
6734	Civic	2001	White	OR	\$17,000
5342	Altima	2001	Green	FL	\$19,000
3845	Maxima	2001	Blue	NY	\$22,000
8354	Accord	2000	Green	VT	\$18,000
4395	Civic	2001	Red	CA	\$17,000
7352	Civic	2002	Red	WA	\$18,000

Example: Database Query Processing

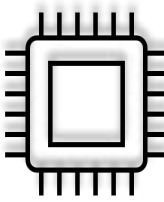


The execution of the query can be divided into subtasks in various ways. Each task can be thought of as generating an intermediate table of entries that satisfy a particular clause.



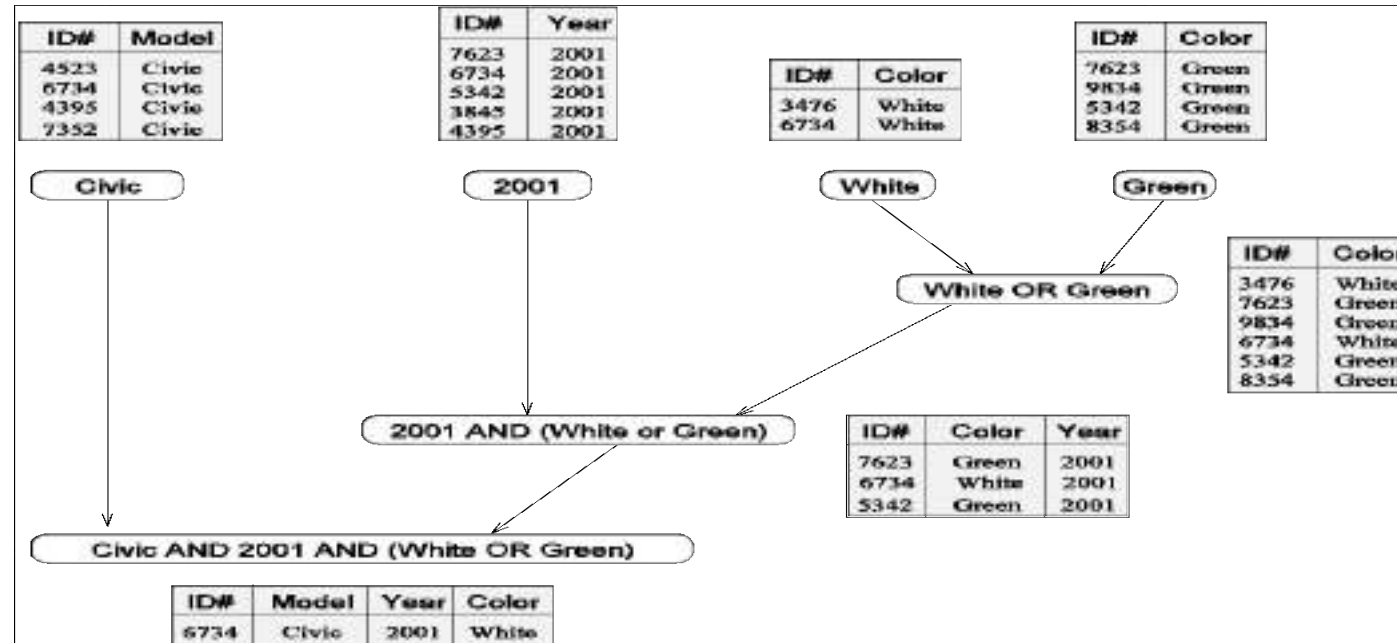
Decomposing the given query into a number of tasks.

Edges in this graph denote that the output of one task is needed to accomplish the next.



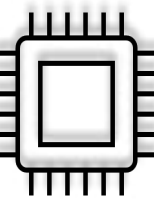
Example: Database Query Processing

Note that the same problem can be decomposed into subtasks in other ways as well.

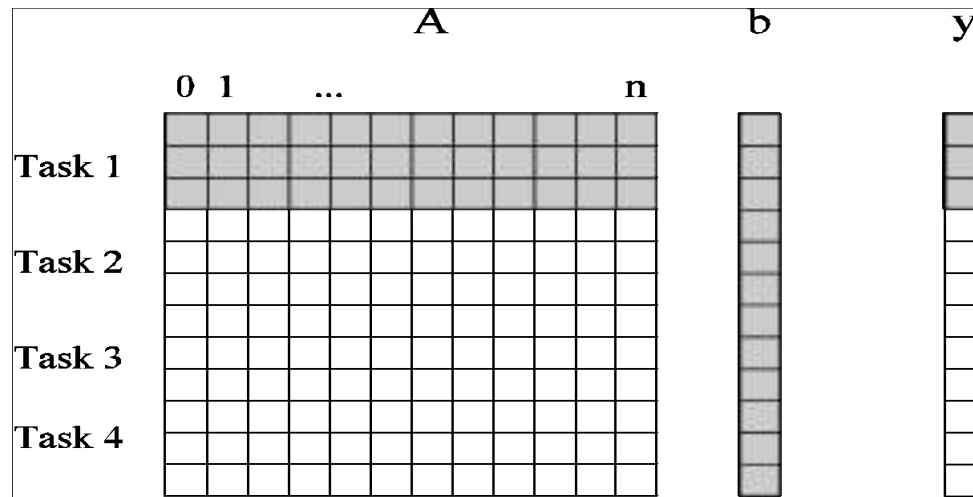


- An alternate decomposition of the given problem into subtasks, along with their data dependencies.
- There are **multiple ways** of expressing certain computations, especially those involving **associative operators such as addition, multiplication, and logical AND or OR**.
- Different task decompositions may lead to significant differences with respect to their eventual parallel performance.

Granularity of Task Decompositions

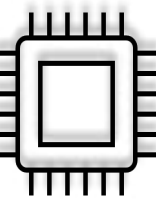


- The number of tasks into which a problem is decomposed determines its *granularity*.
- Decomposition into a large number of tasks results in *fine-grained decomposition* and that into a small number of tasks results in a *coarse grained decomposition*.



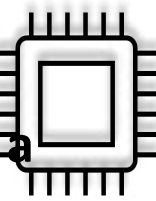
Each task in this example corresponds to the computation of three elements of the result vector.

Degree of Concurrency

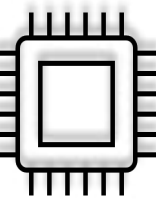


- The number of tasks that can be executed in parallel is the *degree of concurrency* of a decomposition.
- Since the number of tasks that can be executed in parallel may change over program execution, *the maximum degree of concurrency is the maximum number of such tasks at any point during execution. What is the maximum degree of concurrency of the database query examples?*
- The *average degree of concurrency* is the average number of tasks that can be processed in parallel over the execution of the program. *Assuming that each tasks in the database example takes identical processing time, what is the average degree of concurrency in each decomposition?*
- The degree of concurrency increases as the decomposition becomes finer in granularity and vice versa.

Critical Path Length

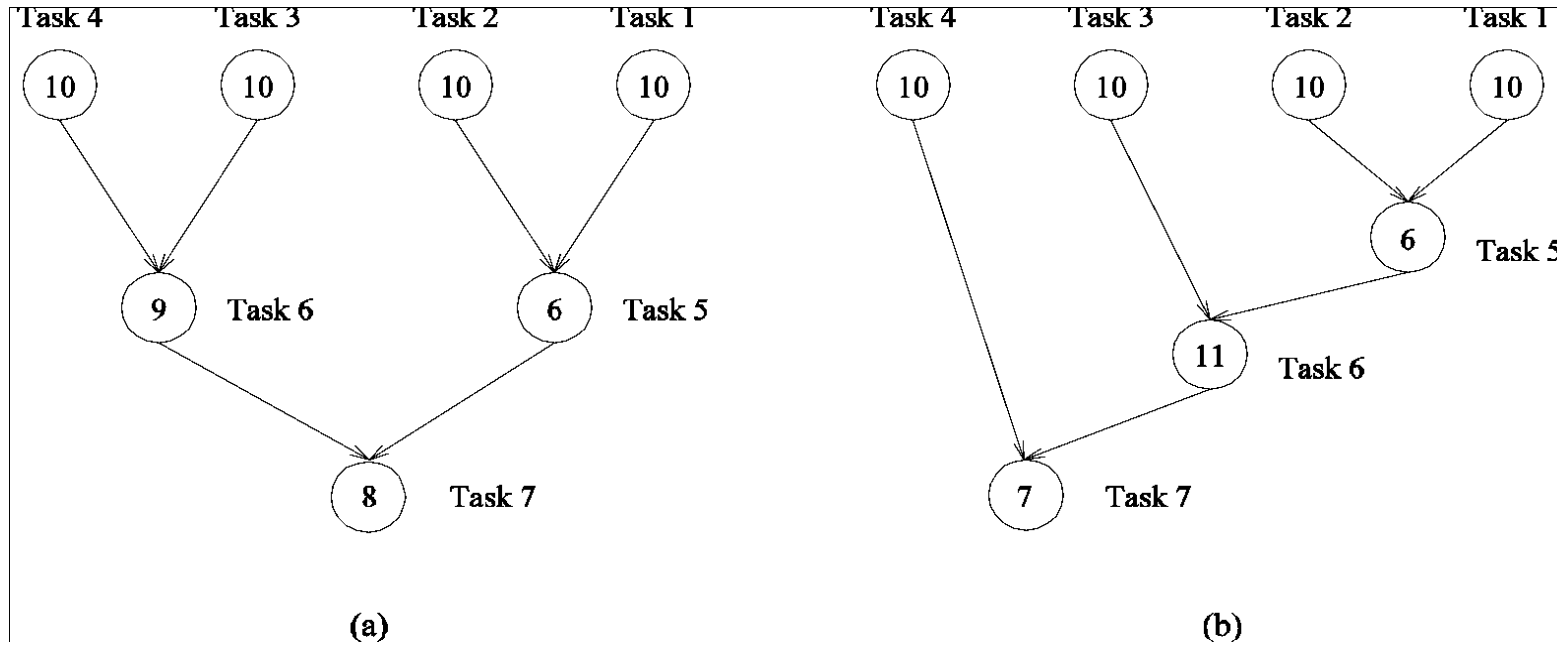


- A feature of a **task-dependency graph** that determines the **average degree of concurrency** for a given **granularity** is its **critical path**.
- In a task-dependency graph, let us refer to the nodes with no incoming edges by start nodes and the nodes with no outgoing edges by finish nodes.
- The longest directed path between any pair of start and finish nodes is known as the critical path.
- A directed path in the task dependency graph represents a sequence of tasks that must be processed one after the other.
- *The length of the longest path in a task dependency graph is called the critical path length.*
- The sum of the weights of nodes along this path is known as the critical path length, where the weight of a node is the size or the amount of work associated with the corresponding task.
- The ratio of the total amount of work to the critical-path length is the average degree of concurrency. Therefore, a shorter critical path favors a higher degree of concurrency

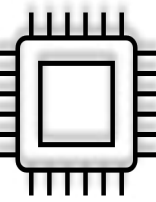


Critical Path Length

Consider the task dependency graphs of the two database query decompositions:

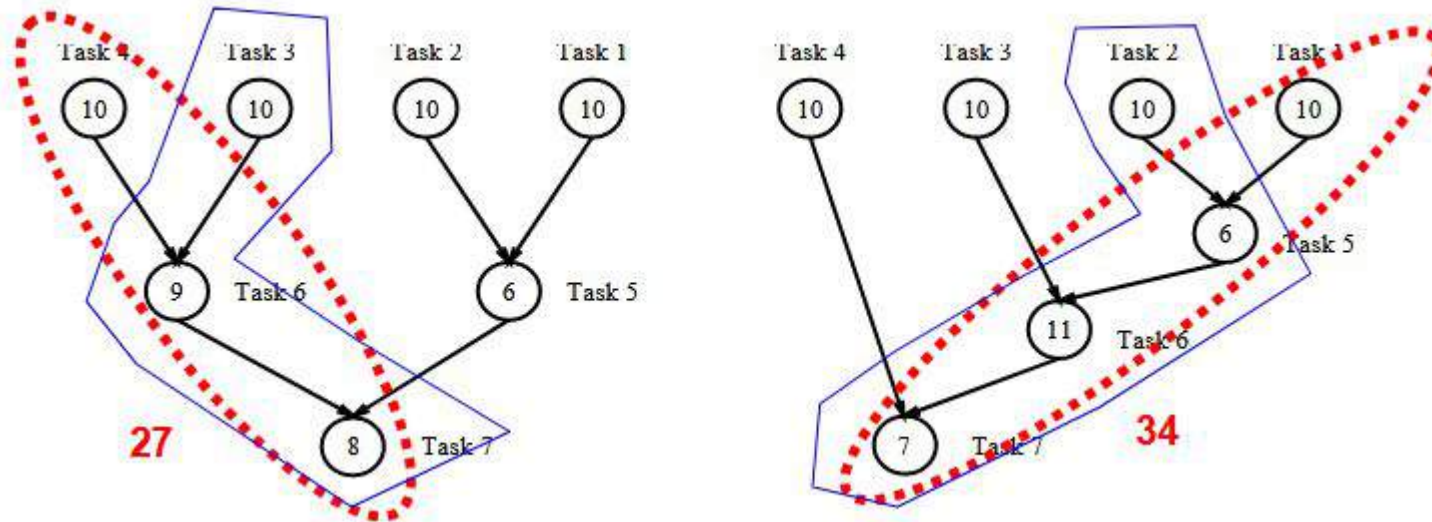


What are the critical path lengths for the two task dependency graphs? What is the maximum degree of concurrency? What is the average degree of concurrency?



Critical Path Length

Consider the task dependency graphs of the two database query decompositions:



What are the critical path lengths for the two task dependency graphs? What is the maximum degree of concurrency? What is the average degree of concurrency?

Solution to average degree of concurrency

Graph A:

Phase 1: Task 1, Task 2, Task 3 and Task 4 can be executed parallelly (if you have more than 3 processors). So the total amount of work in this phase is the sum of the weights of those four nodes: $10+10+10+10 = 40$.

Phase 2: Task 6 and Task 5 can be executed parallelly (if you have more than 1 processor). So the total amount of work in this phase is: $9+6 = 15$.

Phase 3: You can execute only Task 7, so the total amount of work here is 8.

The maximum number of concurrency is $\max(40, 15, 8) = 40$. The average degree of concurrency is $(40+15+8)/(10+9+8) = 63/27 = 2.33$

Graph B:

Phase 1: Task 1, Task 2, Task 3 and Task 4 can be executed parallelly (if you have more than 3 processors). So the total amount of work in this phase is: $10+10+10+10 = 40$.

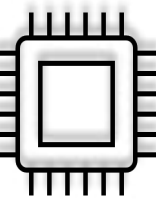
Phase 2: You can execute only Task 5, so the total amount of work here is 6.

Phase 3: You can execute only Task 6, so the total amount of work here is 11.

Phase 4: You can execute only Task 7, so the total amount of work here is 7.

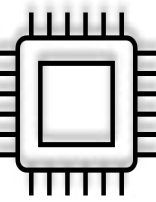
The maximum number of concurrency is $\max(40, 6, 11, 7) = 40$. The average degree of concurrency is $(40+6+11+7)/(10+6+11+7) = 64/34 = 1.88$

Limits on Parallel Performance



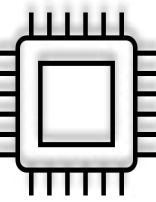
- It would appear that the parallel time can be made arbitrarily small by making the decomposition finer in granularity.
- There is an **inherent bound on how fine the granularity of a computation can be**. *For example, in the case of multiplying a dense matrix with a vector, there can be no more than $O(n^2)$ concurrent tasks.*
- There is another important practical factor that limits our ability to obtain unbounded speedup (ratio of serial to parallel execution time) from parallelization-**the interaction among tasks running on different physical processors.**
- **Concurrent tasks may also have to exchange data with other tasks.** This results in communication overhead. The tradeoff between the granularity of a decomposition and associated overheads often determines performance bounds.

Task Interaction Graphs



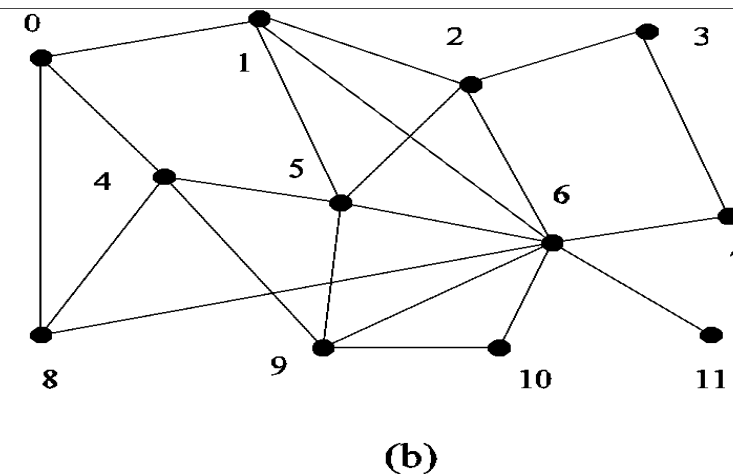
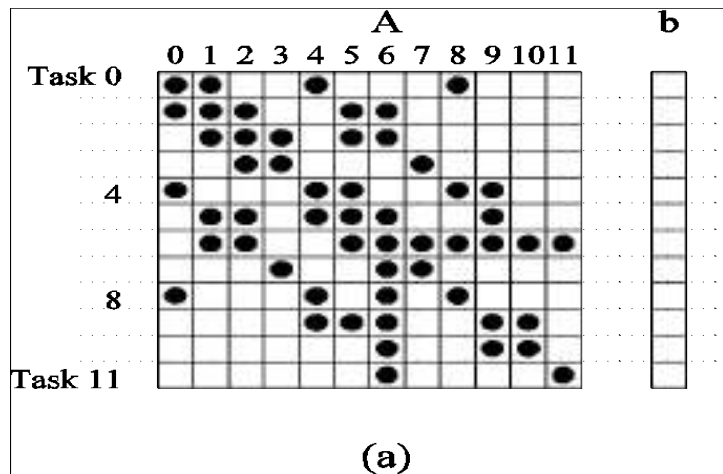
- **Subtasks generally exchange data with others in a decomposition.** For example, even in the trivial decomposition of the dense matrix-vector product, if the vector is not replicated across all tasks, they will have to communicate elements of the vector.
- *The graph of tasks (nodes) and their interactions/data exchange (edges) is referred to as a task interaction graph.*
- Note that *task interaction graphs* represent data dependencies, whereas *task dependency graphs* represent control dependencies.

Task Interaction Graphs: An Example

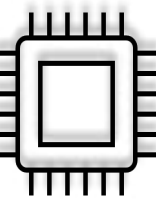


Consider the problem of multiplying a sparse matrix A with a vector b . The following observations can be made:

- As before, the computation of each element of the result vector can be viewed as an independent task.
- Unlike a dense matrix-vector product though, only non-zero elements of matrix A participate in the computation.
- If, for memory optimality, we also partition b across tasks, then one can see that the task interaction graph of the computation is identical to the graph of the matrix A (the graph for which A represents the adjacency structure).
- In general, **if the granularity of a decomposition is finer, the associated overhead** (as a ratio of useful work associated with a task) **increases**.



Processes and Mapping



Process:

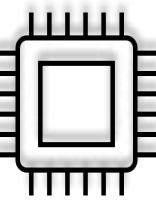
- Refers to a **processing or computing agent** that performs tasks.
- It is an abstract entity that **uses the code and data** corresponding to a task to produce the output of that task within a finite amount of time after the task is activated by the parallel program.

Mapping:

- **The mechanism by which tasks are assigned to processes for execution is called mapping.**
- Even though the degree of concurrency is determined by the decomposition, it is the mapping that determines how much of that concurrency is actually utilized, and how efficiently.
- In general, the number of tasks in a decomposition exceeds the number of processing elements available.
- For this reason, a parallel algorithm must also provide a mapping of tasks to processes.

Note: We refer to the mapping as being from tasks to processes, as opposed to processors. This is because typical programming APIs, as we shall see, do not allow easy binding of tasks to physical processors. Rather, we aggregate tasks into processes and rely on the system to map these processes to physical processors. We use processes, not in the UNIX sense of a process, rather, simply as a collection of tasks and associated data.

Processes and Mapping



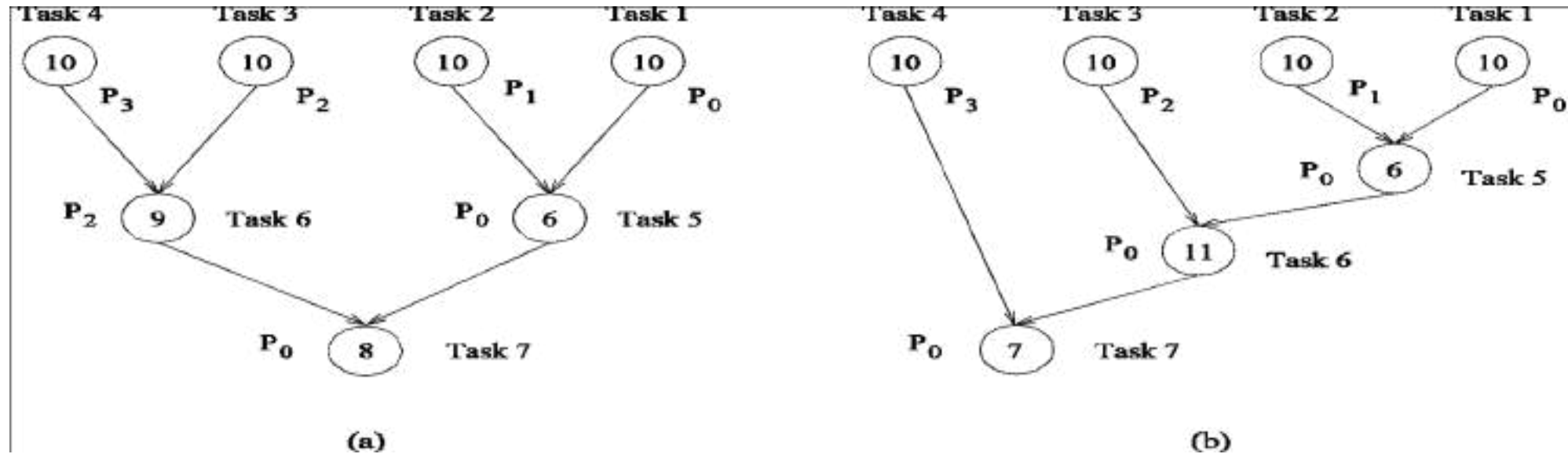
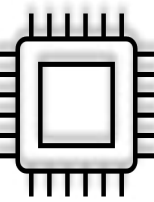
- Appropriate *mapping of tasks to processes* is critical to the parallel performance of an algorithm.
- Mappings are determined by both the **task dependency and task interaction graphs**.
- **Task dependency graphs** can be used to ensure that work is equally spread across all processes at any point (**minimum idling and optimal load balance**).
- **Task interaction graphs** can be used to make sure that processes need minimum interaction with other processes (**minimum communication**).

An appropriate mapping must minimize parallel execution time by:

- Mapping independent tasks to different processes.
- Assigning tasks on critical path to processes as soon as they become available.
- Minimizing interaction between processes by mapping tasks with dense interactions to the same process.

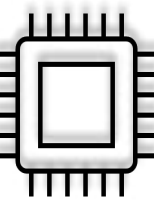
Note: These criteria often conflict with each other. For example, a decomposition into one task (or no decomposition at all) minimizes interaction but does not result in a speedup at all!

Processes and Mapping: Example



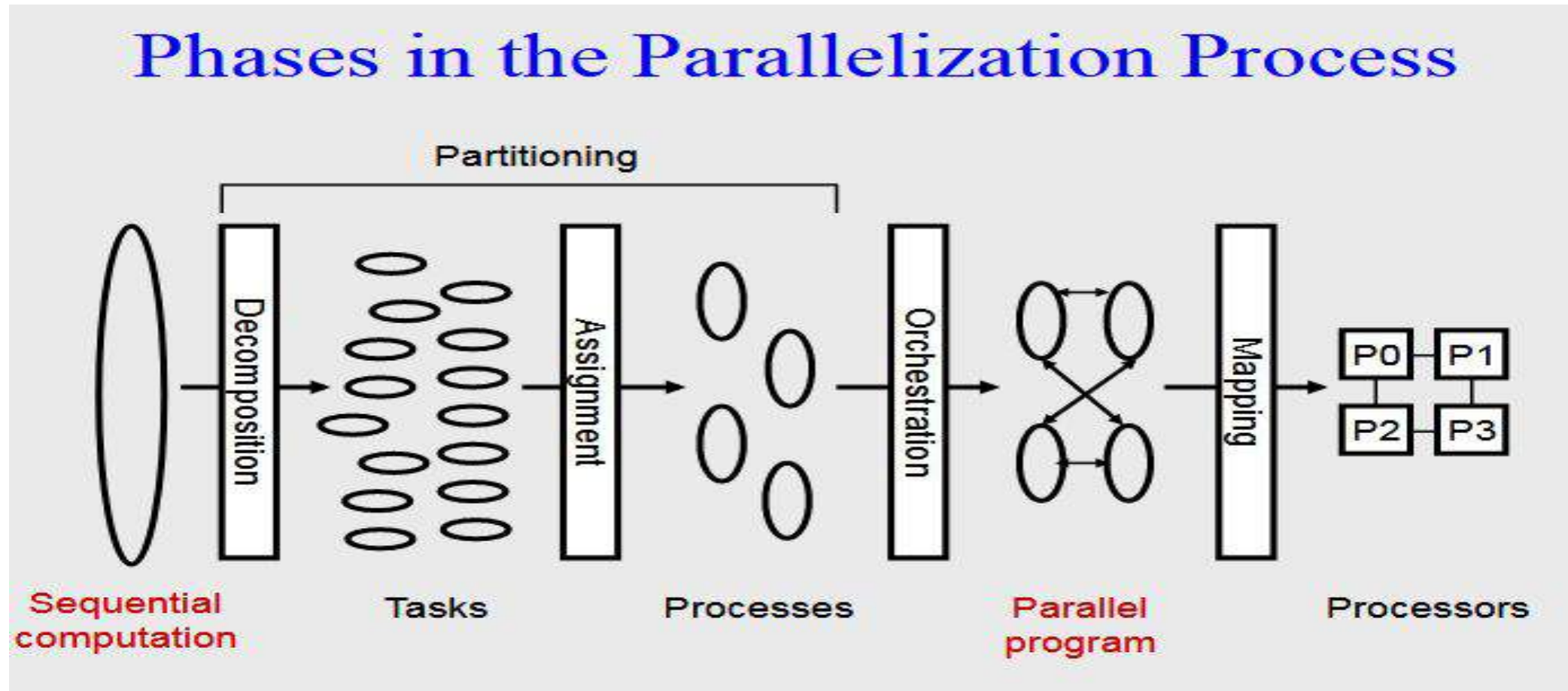
Mapping tasks in the database query decomposition to processes. These mappings were arrived at by viewing the dependency graph in terms of levels (no two nodes in a level have dependencies). Tasks within a single level are then assigned to different processes.

Processes versus Processors

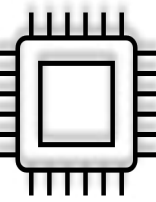


Processes are logical computing agents that perform tasks.

Processors are the hardware units that physically perform computations.



Decomposition Techniques

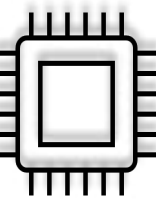


So how does **one decompose a task into various subtasks?**

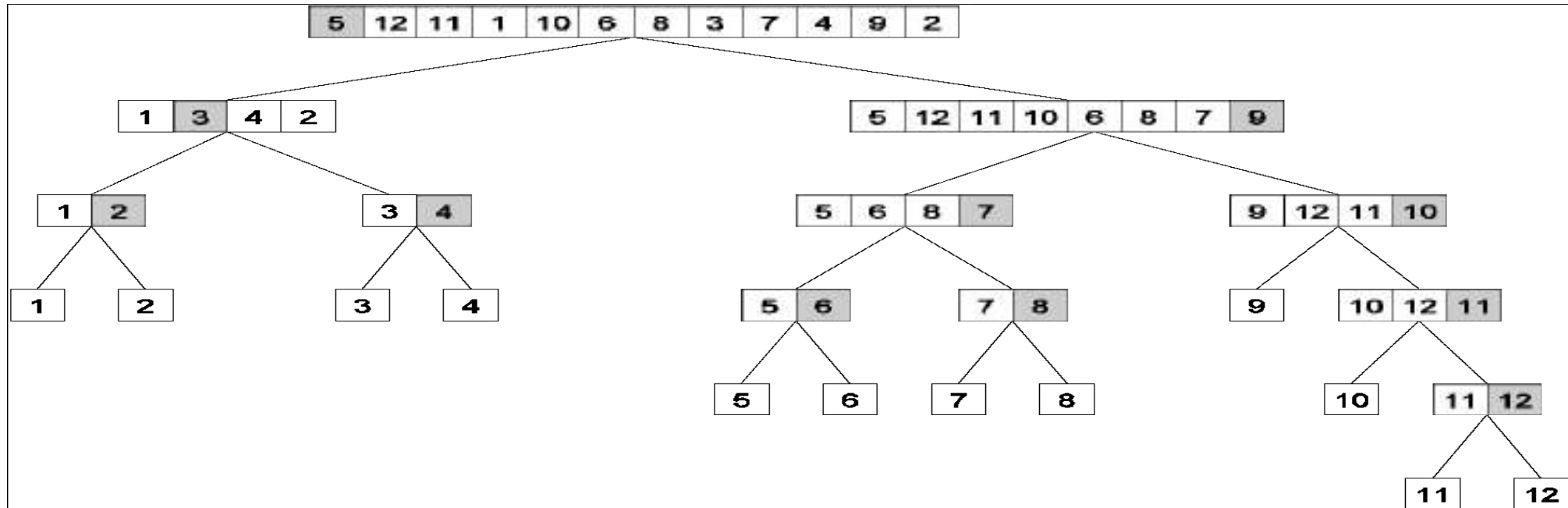
While there is no single recipe that works for all problems, we present a set of commonly used techniques that apply to broad classes of problems. These include:

- Recursive decomposition
- Data decomposition
- Exploratory decomposition
- Speculative decomposition
- Hybrid decomposition

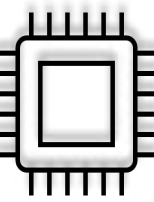
Recursive Decomposition



- Generally suited to problems that are solved using the **divide-and-conquer strategy**.
- A given problem is first decomposed into a set of sub-problems.
- These sub-problems are recursively decomposed further until a desired granularity is reached.
- A classic example of a divide-and-conquer algorithm on which we can apply recursive decomposition is **Quicksort**.
- In this example, once the list has been partitioned around the pivot, each sublist can be processed concurrently (i.e., each sublist represents an independent subtask). This can be repeated recursively.

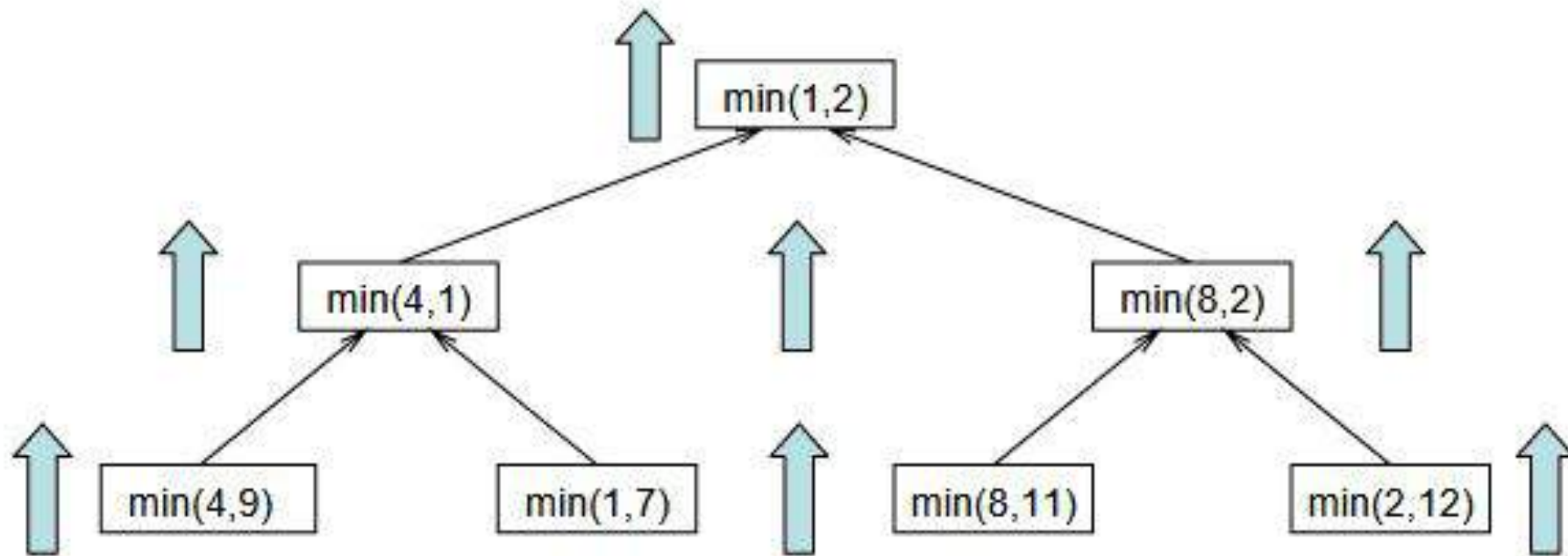


Recursive Decomposition: Example

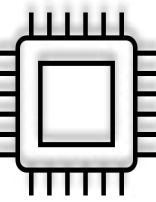


The problem of **finding the minimum number** in a given list (or indeed any other associative operation such as sum, AND, etc.) can be fashioned as a divide-and-conquer algorithm.

Consider the following example of finding the minimum number in the set {4, 9, 1, 7, 8, 11, 2, 12}. The task dependency graph associated with this computation is as follows:



Data Decomposition

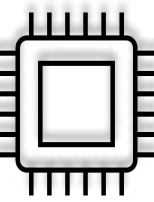


- Identify the data on which computations are performed.
- Partition this data across various tasks.
- This partitioning induces a decomposition of the problem.
- *Data is different with every processor but computation is same.*
- Data can be partitioned in various ways - this critically impacts performance of a parallel algorithm.

Types:

- Partitioning Output Data
- Partitioning Input Data
- Partitioning both Input and Output Data
- Partitioning Intermediate Data

Output Data Decomposition: Example



Often, each element of the output can be computed independently of others (but simply as a function of the input).

A partition of the output across tasks decomposes the problem naturally.

Consider the problem of multiplying two $n \times n$ matrices A and B to yield matrix C . The output matrix C can be partitioned into four tasks as follows:

$$\begin{pmatrix} A_{1,1} & A_{1,2} \\ A_{2,1} & A_{2,2} \end{pmatrix} \cdot \begin{pmatrix} B_{1,1} & B_{1,2} \\ B_{2,1} & B_{2,2} \end{pmatrix} \rightarrow \begin{pmatrix} C_{1,1} & C_{1,2} \\ C_{2,1} & C_{2,2} \end{pmatrix}$$

$$C_{1,1} = A_{1,1}B_{1,1} + A_{1,2}B_{2,1}$$

$$C_{1,2} = A_{1,1}B_{1,2} + A_{1,2}B_{2,2}$$

$$C_{2,1} = A_{2,1}B_{1,1} + A_{2,2}B_{2,1}$$

$$C_{2,2} = A_{2,1}B_{1,2} + A_{2,2}B_{2,2}$$

Task 1: $C_{1,1} = A_{1,1}B_{1,1} + A_{1,2}B_{2,1}$

Task 2: $C_{1,2} = A_{1,1}B_{1,2} + A_{1,2}B_{2,2}$

Task 3: $C_{2,1} = A_{2,1}B_{1,1} + A_{2,2}B_{2,1}$

Task 4: $C_{2,2} = A_{2,1}B_{1,2} + A_{2,2}B_{2,2}$

Decomposition I

Task 1: $C_{1,1} = A_{1,1} B_{1,1}$

Task 2: $C_{1,1} = C_{1,1} + A_{1,2} B_{2,1}$

Task 3: $C_{1,2} = A_{1,1} B_{1,2}$

Task 4: $C_{1,2} = C_{1,2} + A_{1,2} B_{2,2}$

Task 5: $C_{2,1} = A_{2,1} B_{1,1}$

Task 6: $C_{2,1} = C_{2,1} + A_{2,2} B_{2,1}$

Task 7: $C_{2,2} = A_{2,1} B_{1,2}$

Task 8: $C_{2,2} = C_{2,2} + A_{2,2} B_{2,2}$

Decomposition II

Task 1: $C_{1,1} = A_{1,1} B_{1,1}$

Task 2: $C_{1,1} = C_{1,1} + A_{1,2} B_{2,1}$

Task 3: $C_{1,2} = A_{1,2} B_{2,2}$

Task 4: $C_{1,2} = C_{1,2} + A_{1,1} B_{1,2}$

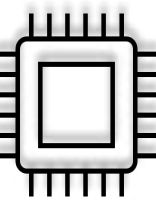
Task 5: $C_{2,1} = A_{2,2} B_{2,1}$

Task 6: $C_{2,1} = C_{2,1} + A_{2,1} B_{1,1}$

Task 7: $C_{2,2} = A_{2,1} B_{1,2}$

Task 8: $C_{2,2} = C_{2,2} + A_{2,2} B_{2,2}$

Output Data Decomposition: Example



Consider the problem of counting the instances of given itemsets in a database of transactions. In this case, the output (itemset frequencies) can be partitioned across tasks.

(a) Transactions (input), itemsets (input), and frequencies (output)

Database Transactions	A, B, C, E, G, H	Itemsets	A, B, C	Itemset Frequency	1
	B, D, E, F, K, L		D, E		3
	A, B, F, H, L		C, F, G		0
	D, E, F, H		A, E		2
	F, G, H, K,		C, D		1
	A, E, F, K, L		D, K		2
	B, C, D, G, H, L		B, C, F		0
	G, H, L		C, D, K		0
	D, E, F, K, L				
	F, G, H, L				

(b) Partitioning the frequencies (and itemsets) among the tasks

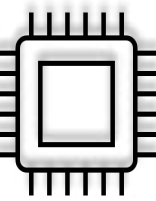
Database Transactions	A, B, C, E, G, H	Itemsets	A, B, C	Itemset Frequency	1
	B, D, E, F, K, L		D, E		3
	A, B, F, H, L		C, F, G		0
	D, E, F, H		A, E		2
	F, G, H, K,				
	A, E, F, K, L				
	B, C, D, G, H, L				
	G, H, L				
	D, E, F, K, L				
	F, G, H, L				

task 1

Database Transactions	A, B, C, E, G, H	Itemsets	C, D	Itemset Frequency	1
	B, D, E, F, K, L		D, K		2
	A, B, F, H, L		B, C, F		0
	D, E, F, H		C, D, K		0
	F, G, H, K,				
	A, E, F, K, L				
	B, C, D, G, H, L				
	G, H, L				
	D, E, F, K, L				
	F, G, H, L				

task 2

Input Data Partitioning



- Generally applicable if each **output can be naturally computed as a function of the input.**
- In many cases, this is the only natural decomposition because the output is not clearly known a-priori (e.g., the problem of finding the minimum in a list, sorting a given list, etc.).
- A task is associated with each input data partition. The task performs as much of the computation with its part of the data. Subsequent processing combines these partial results.
- In the database counting example, the input (i.e., the transaction set) can be partitioned. This induces a task decomposition in which each task generates partial counts for all itemsets. These are combined subsequently for aggregate counts.

Partitioning the transactions among the tasks

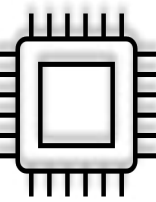
Database Transactions	A, B, C, E, G, H	Itemsets	A, B, C	Itemset Frequency	1
	B, D, E, F, K, L		D, E		2
	A, B, F, H, L		C, F, G		0
	D, E, F, H		A, E		1
	F, G, H, K,		C, D		0
			D, K		1
			B, C, F		0
			C, D, K		0

task 1

Database Transactions	A, E, F, K, L	Itemsets	A, B, C	Itemset Frequency	0
	B, C, D, G, H, L		D, E		1
	G, H, L		C, F, G		0
	D, E, F, K, L		A, E		1
	F, G, H, L		C, D		1
			D, K		1
			B, C, F		0
			C, D, K		0

task 2

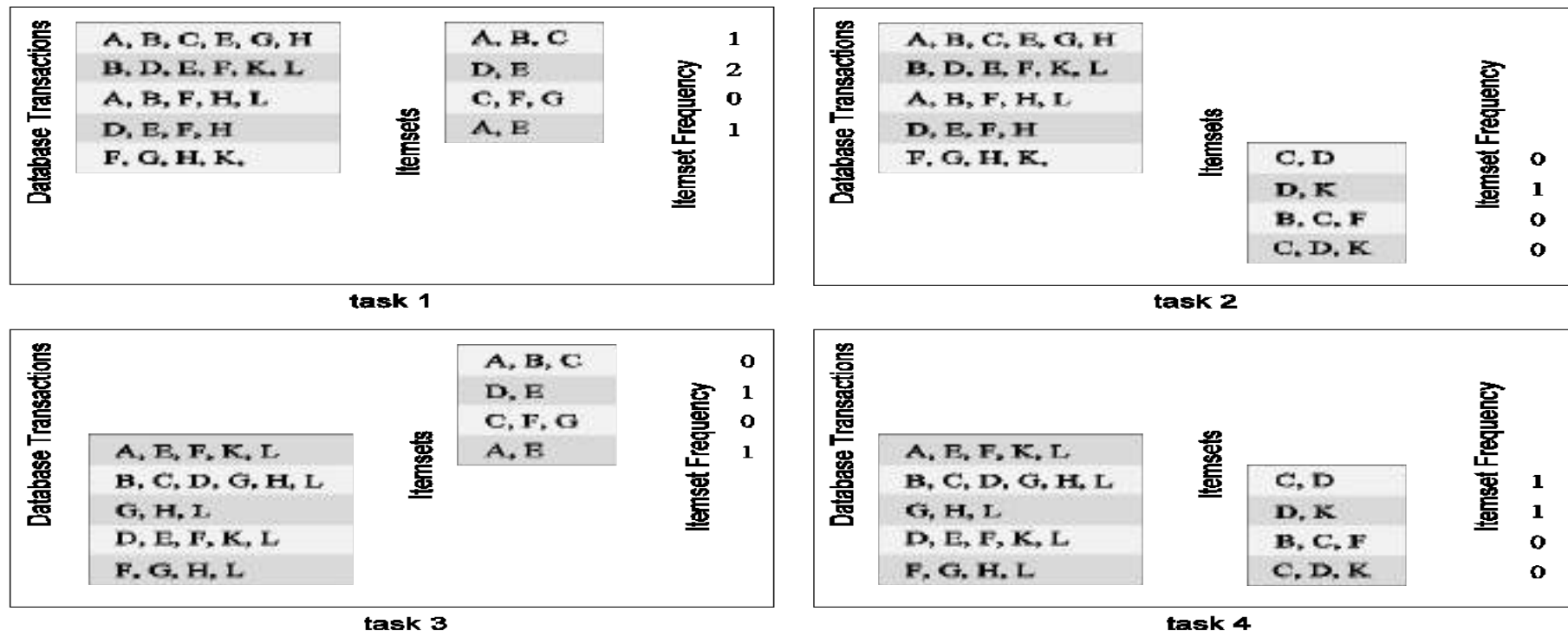
Partitioning Input *and* Output Data



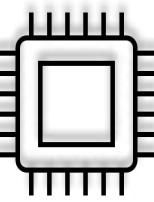
From the previous example, the following observations can be made:

- If the database of transactions is replicated across the processes, **each task can be independently accomplished with no communication.**
- If the database is partitioned across processes as well (**for reasons of memory utilization**), each task first computes partial counts. These counts are then aggregated at the appropriate task.
- Often input and output data decomposition can be combined for a higher degree of concurrency.

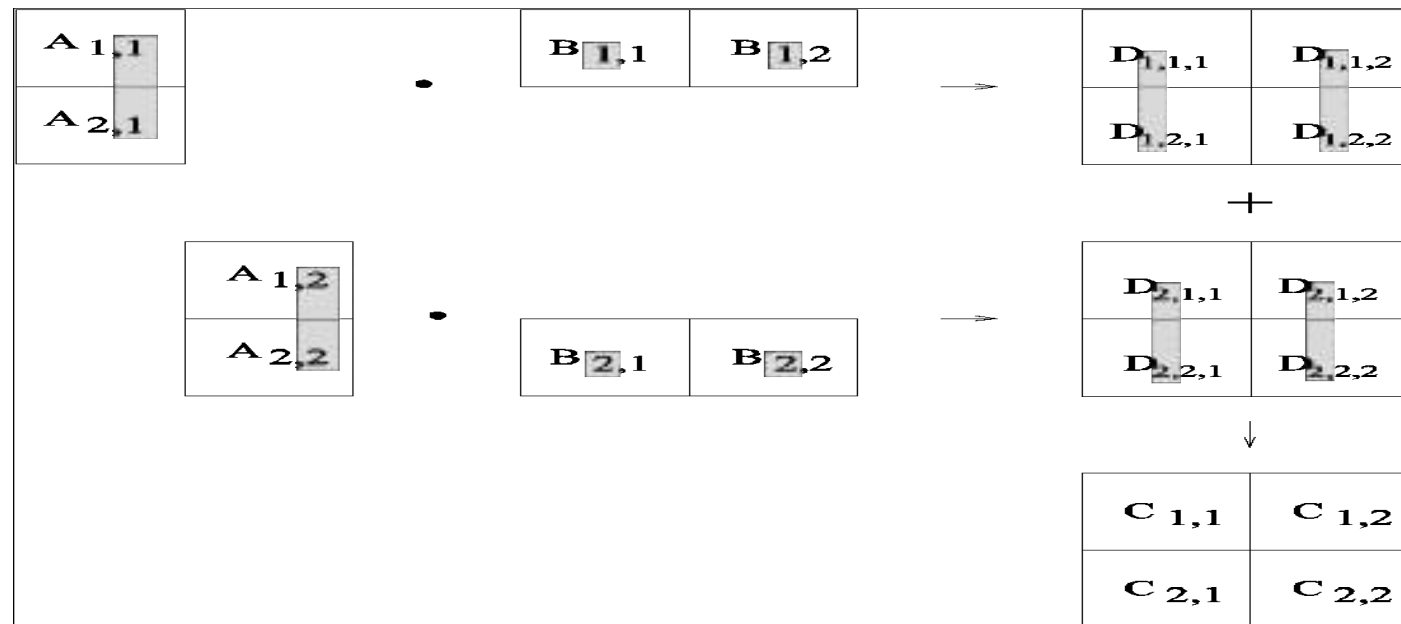
Partitioning both transactions and frequencies among the tasks

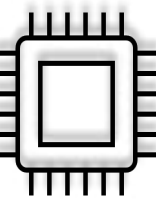


Intermediate Data Partitioning



- Computation can often be viewed as a sequence of transformation from the input to the output data.
- In these cases, it is often beneficial to use one of the intermediate stages as a basis for decomposition.
- Let us revisit the example of dense matrix multiplication. We first show how we can visualize this computation in terms of intermediate matrices D .





Intermediate Data Partitioning: Example

A decomposition of intermediate data structure leads to the following decomposition into 8 + 4 tasks:

Stage I

$$\begin{pmatrix} A_{1,1} & A_{1,2} \\ A_{2,1} & A_{2,2} \end{pmatrix} \cdot \begin{pmatrix} B_{1,1} & B_{1,2} \\ B_{2,1} & B_{2,2} \end{pmatrix} \rightarrow \left(\begin{pmatrix} D_{1,1,1} & D_{1,1,2} \\ D_{1,2,2} & D_{1,2,2} \end{pmatrix} \right)$$

Stage II

$$\begin{pmatrix} D_{1,1,1} & D_{1,1,2} \\ D_{1,2,2} & D_{1,2,2} \end{pmatrix} + \begin{pmatrix} D_{2,1,1} & D_{2,1,2} \\ D_{2,2,2} & D_{2,2,2} \end{pmatrix} \rightarrow \begin{pmatrix} C_{1,1} & C_{1,2} \\ C_{2,1} & C_{2,2} \end{pmatrix}$$

Task 01: $D_{1,1,1} = A_{1,1} B_{1,1}$

Task 03: $D_{1,1,2} = A_{1,1} B_{1,2}$

Task 05: $D_{1,2,1} = A_{2,1} B_{1,1}$

Task 07: $D_{1,2,2} = A_{2,1} B_{1,2}$

Task 09: $C_{1,1} = D_{1,1,1} + D_{2,1,1}$

Task 11: $C_{2,1} = D_{1,2,1} + D_{2,2,1}$

Task 02: $D_{2,1,1} = A_{1,2} B_{2,1}$

Task 04: $D_{2,1,2} = A_{1,2} B_{2,2}$

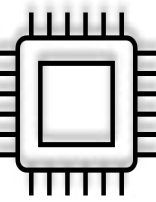
Task 06: $D_{2,2,1} = A_{2,2} B_{2,1}$

Task 08: $D_{2,2,2} = A_{2,2} B_{2,2}$

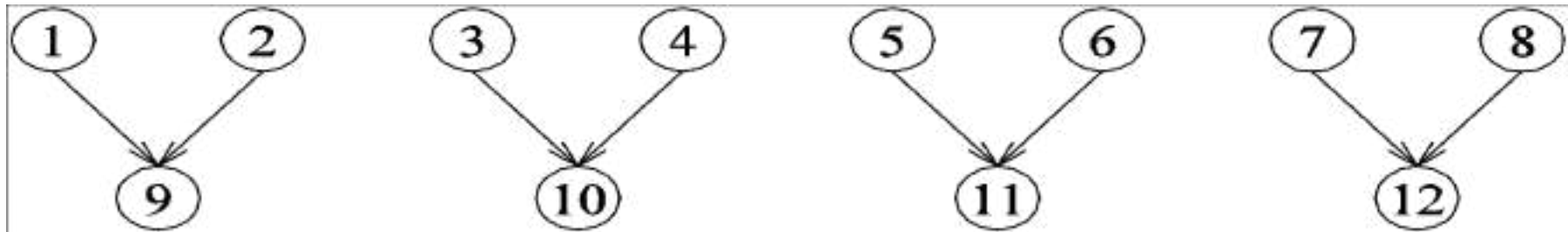
Task 10: $C_{1,2} = D_{1,1,2} + D_{2,1,2}$

Task 12: $C_{2,2} = D_{1,2,2} + D_{2,2,2}$

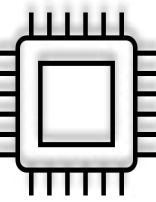
Intermediate Data Partitioning: Example



The task dependency graph for the decomposition (shown in previous foil) into 12 tasks is as follows:

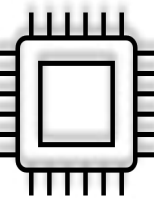


The Owner Computes Rule



- The *Owner Computes Rule* generally states that the process assigned a particular data item is responsible for all computation associated with it.
- In the case of *input data decomposition*, the owner computes rule implies that all computations that use the input data are performed by the process.
- In the case of *output data decomposition*, the owner computes rule implies that the output is computed by the process to which the output data is assigned.

Exploratory Decomposition



- In many cases, the **decomposition** of the problem goes **hand-in-hand with its execution**.
- These problems typically involve the exploration **(search) of a state space of solutions**.
- Problems in this class include a variety of discrete optimization problems (0/1 integer programming), theorem proving, game playing, etc.
- A simple application of exploratory decomposition is in the solution to a 15 puzzle (a tile puzzle). We show a sequence of three moves that transform a given initial state (a) to desired final state (d).
- Of-course, the problem of computing the solution, in general, is much more difficult than in this simple example

1	2	3	4
5	6	7	8
9	10	11	12
13	14	15	16

(a)

1	2	3	4
5	6	7	8
9	10	11	12
13	14	15	16

(b)

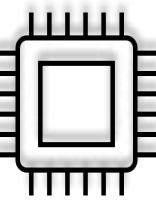
1	2	3	4
5	6	7	8
9	10	11	12
13	14	15	16

(c)

1	2	3	4
5	6	7	8
9	10	11	12
13	14	15	16

(d)

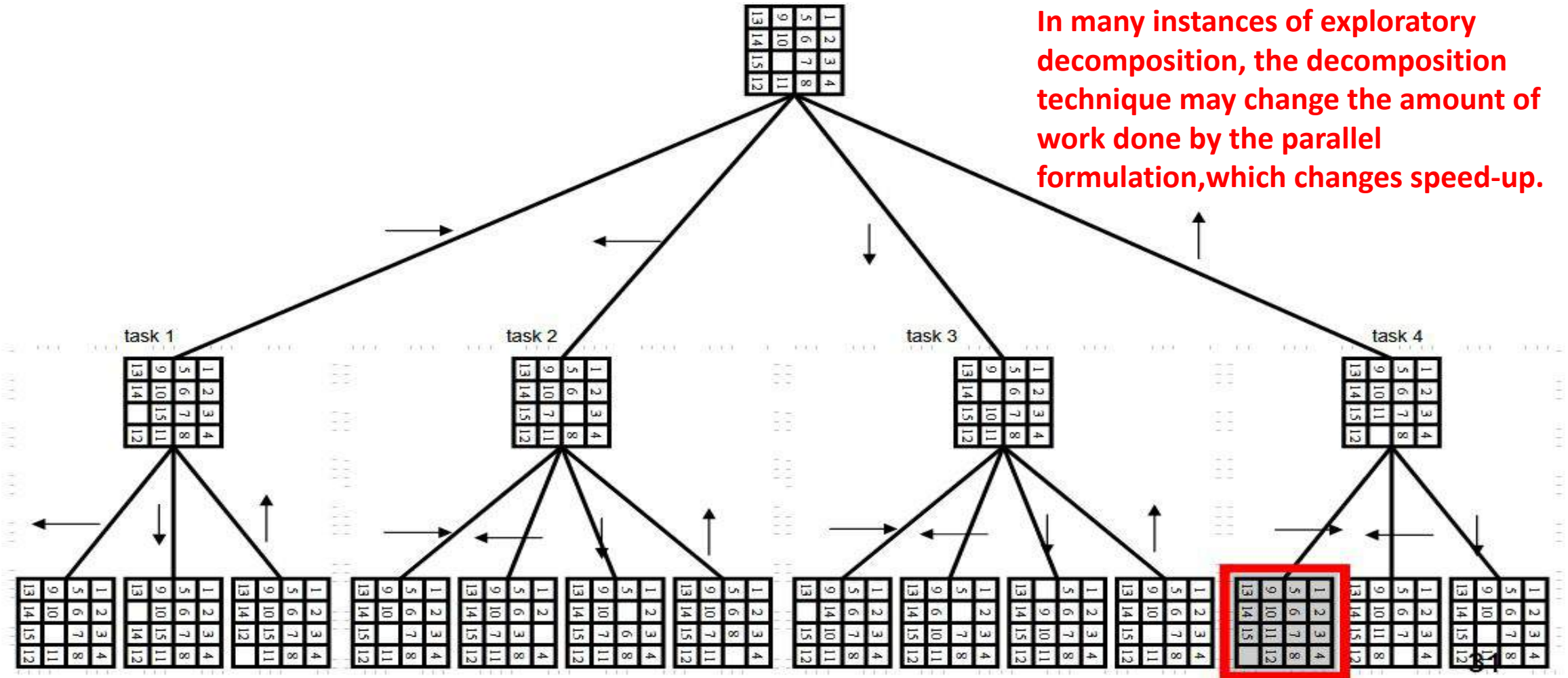
Exploratory Decomposition: Example



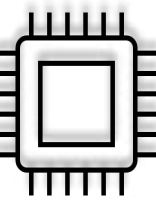
The state space can be explored by generating various successor states of the current state and to view them as independent tasks.

Note:

In many instances of exploratory decomposition, the decomposition technique may change the amount of work done by the parallel formulation, which changes speed-up.

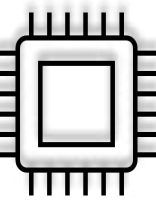


Speculative Decomposition



- **Speculative decomposition** is used when a program may take one of many possible computationally significant branches depending on the output of other computations that precede it.
- In this situation, 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.
- In some applications, dependencies between tasks are not known a-priori.
- For such applications, it is impossible to identify independent tasks.
- There are generally two approaches to dealing with such applications:
 - **Conservative approaches**, which identify independent tasks only when they are guaranteed to not have dependencies,
 - **Optimistic approaches**, which schedule tasks even when they may potentially be erroneous.
- Conservative approaches may yield little concurrency and optimistic approaches may require roll-back mechanism in the case of an error.

Speculative Decomposition: Example



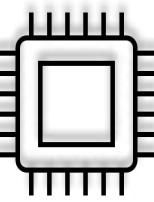
A classic example of speculative decomposition is in **discrete event simulation**.

- The central data structure in a discrete event simulation is a **time-ordered event list**.
- Events are extracted precisely in time order, processed, and if required, resulting events are inserted back into the event list.
- Consider your day today as a discrete event system - you get up, get ready, drive to work, work, eat lunch, work some more, drive back, eat dinner, and sleep.
- Each of these events may be processed independently, however, in driving to work, you might meet with an unfortunate accident and not get to work at all.
- Therefore, an optimistic scheduling of other events will have to be rolled back.

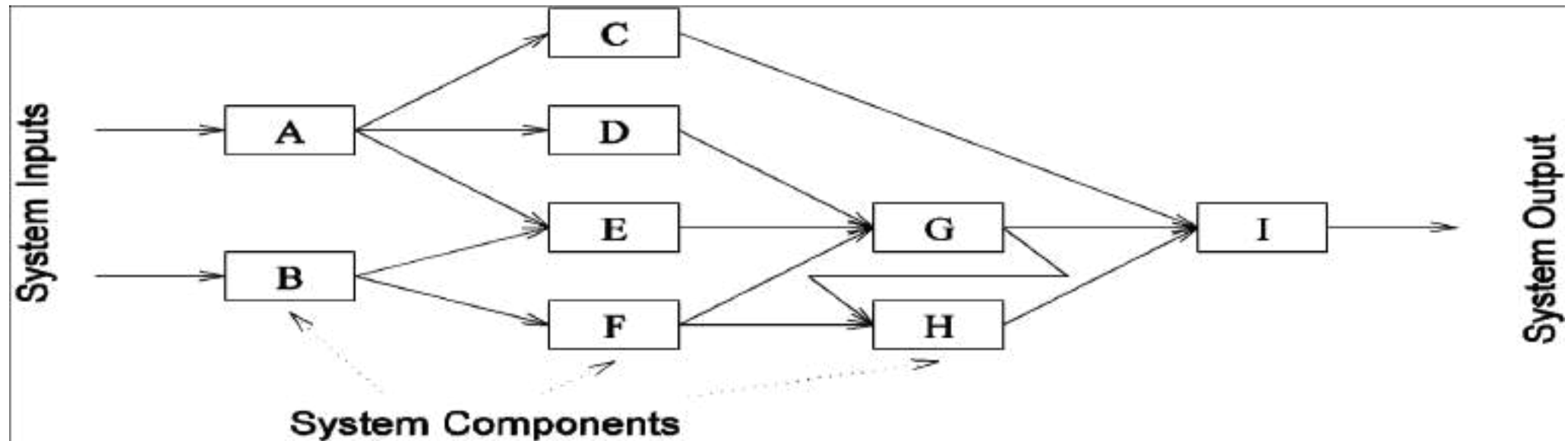
Example :Switch Case

Among multiple cases available which one is to be considered is based on the input(expression) which has come from its preceding part.

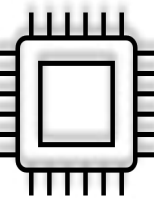
Speculative Decomposition: Example



Another example is the simulation of a **network of nodes** (for instance, an assembly line or a computer network through which packets pass). The task is to simulate the behavior of this network for various inputs and node delay parameters (note that networks may become unstable for certain values of service rates, queue sizes, etc.).

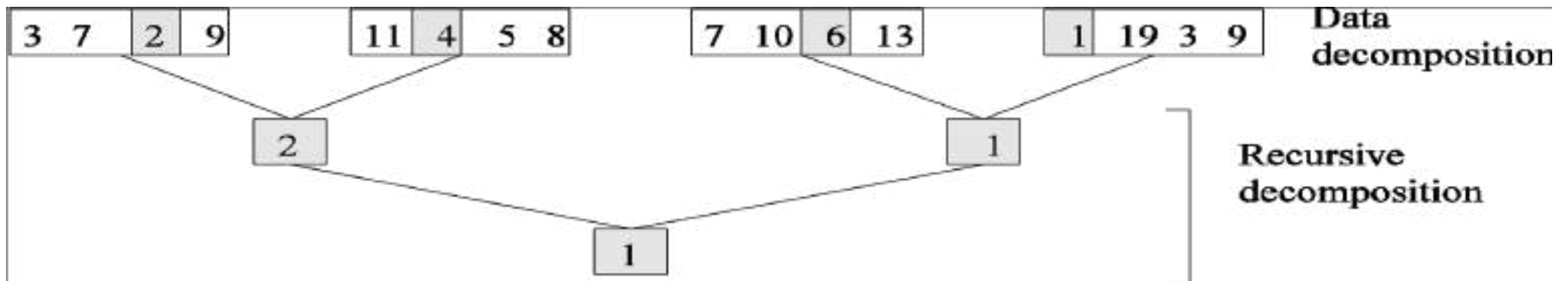


Hybrid Decompositions

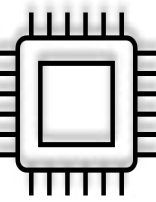


Often, a mix of decomposition techniques is necessary for decomposing a problem. Consider the following examples:

- In quicksort, recursive decomposition alone limits concurrency (Why?). A mix of data and recursive decompositions is more desirable.
- In discrete event simulation, there might be concurrency in task processing. A mix of speculative decomposition and data decomposition may work well.
- Even for simple problems like finding a minimum of a list of numbers, a mix of data and recursive decomposition works well.



Characteristics of Tasks



Once a problem has been decomposed into independent tasks, the characteristics of these tasks critically impact choice and performance of parallel algorithms. Relevant task characteristics include:

a) Task generation:

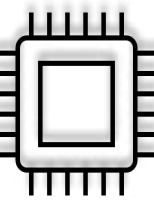
Static task generation:

Concurrent tasks can be identified a-priori. Typical matrix operations, graph algorithms, image processing applications, and other regularly structured problems fall in this class. These can typically be decomposed using data or recursive decomposition techniques.

Dynamic task generation:

Tasks are generated as we perform computation. A classic example of this is in game playing - each 15 puzzle board is generated from the previous one. These applications are typically decomposed using exploratory or speculative decompositions.

Characteristics of Tasks



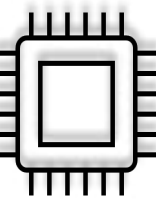
b) Task sizes:

- Task sizes may be **uniform** (i.e., all tasks are the same size) or **non-uniform**.
- Non-uniform task sizes may be such that they can be determined (or estimated) a-priori or not.
- Examples in this class include discrete optimization problems, in which it is difficult to estimate the effective size of a state space.

c) Size of data associated with tasks:

- The size of data associated with a task may be small or large when viewed in the context of the size of the task.
- A **small context** of a task implies that an algorithm can **easily communicate this task to other processes dynamically** (e.g., the 15 puzzle).
- A **large context** ties the task to a process, or alternately, an algorithm may attempt to **reconstruct the context at another processes as opposed to communicating the context of the task** (e.g., 0/1 integer programming).

Characteristics of Task Interactions



Tasks may communicate with each other in various ways. The associated dichotomy is:

Static versus dynamic interactions:

Static interactions:

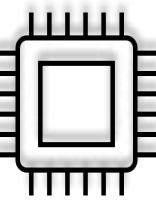
- The tasks and their interactions are known a-priori.
- These are relatively simpler to code into programs.
- Example: Matrix Multiplication.

Dynamic interactions:

- The timing or interacting tasks cannot be determined a-priori.
- These interactions are harder to code.

Ex. Puzzle game. The task has exhausted its work can pick up an unexplored state from the queue of another busy task and start exploring it.

Characteristics of Task Interactions



Tasks may communicate with each other in various ways. The associated dichotomy is:

Regular versus Irregular interactions:

Regular interactions:

- There is a definite pattern (in the graph sense) to the interactions.
- These patterns can be exploited for efficient implementation.

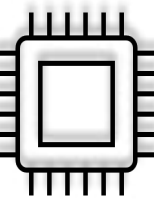
Example: Explicit finite difference for solving PDEs. Image dithering.

Irregular interactions:

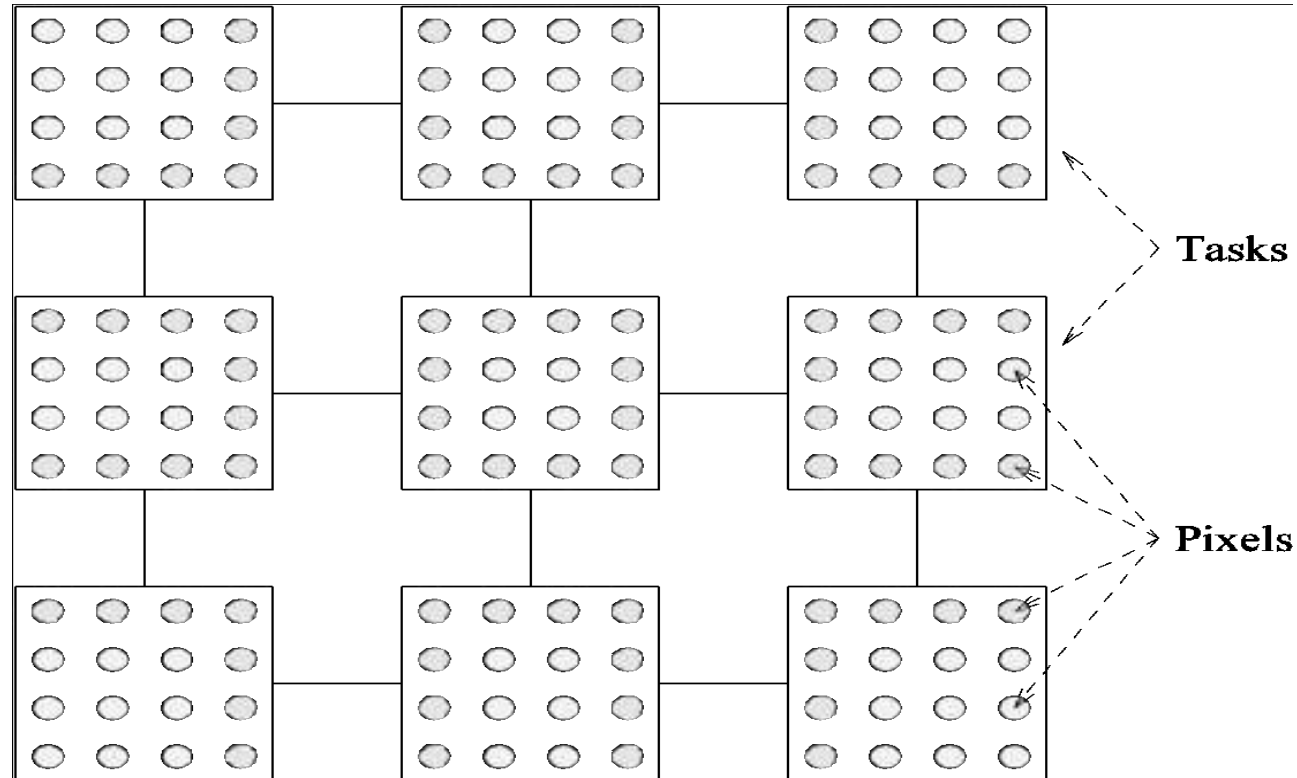
- Interactions lack well-defined topologies.

Example: Sparse matrix-vector multiplication

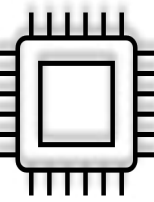
Characteristics of Task Interactions: Example



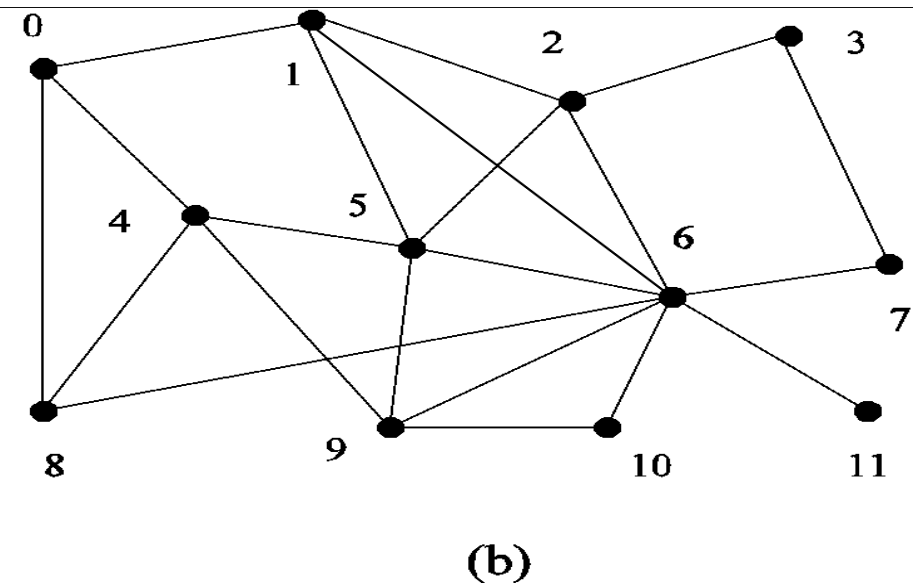
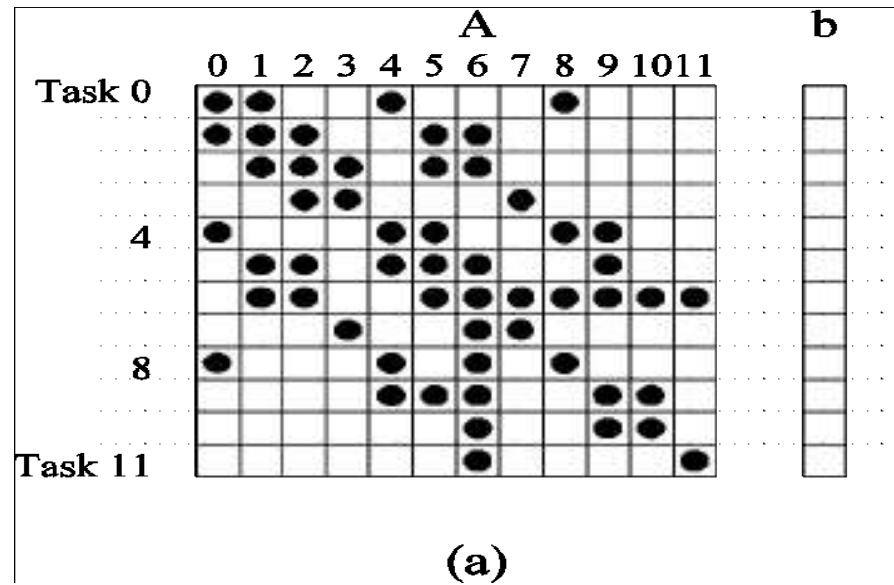
A simple example of a regular static interaction pattern is in image dithering. The underlying communication pattern is a structured (2-D mesh) one as shown here:



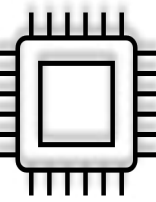
Characteristics of Task Interactions: Example



The multiplication of a sparse matrix with a vector is a good example of a static irregular interaction pattern. Here is an example of a sparse matrix and its associated interaction pattern.



Characteristics of Task Interactions



Tasks may communicate with each other in various ways. The associated dichotomy is:

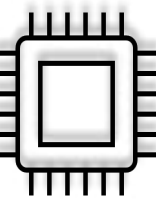
Read-Only versus Read-Write interactions:

- In **read-only interactions**, tasks just read data items associated with other tasks.
- In **read-write interactions** tasks read, as well as modify data items associated with other tasks.
- In general, read-write interactions are harder to code, since they require additional synchronization primitives.

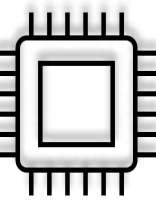
One-way versus two way interactions:

- A **one-way interaction** can be initiated and accomplished by one of the two interacting tasks.
- A **two-way interaction** requires participation from both tasks involved in an interaction.
- One way interactions are somewhat harder to code in message passing APIs.

Mapping Techniques

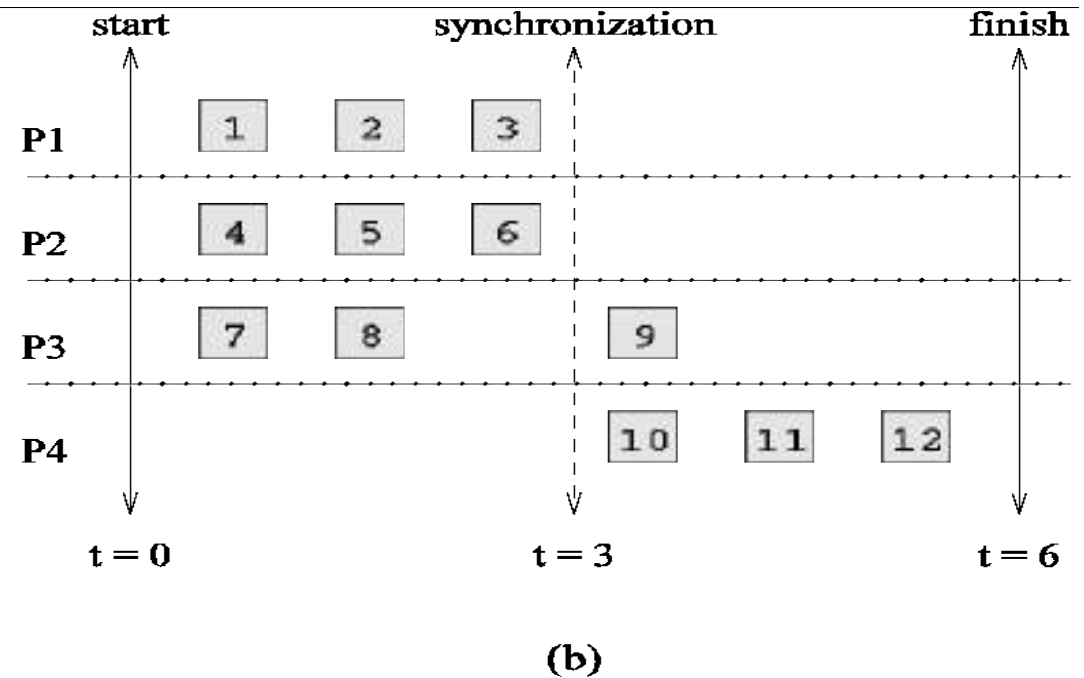
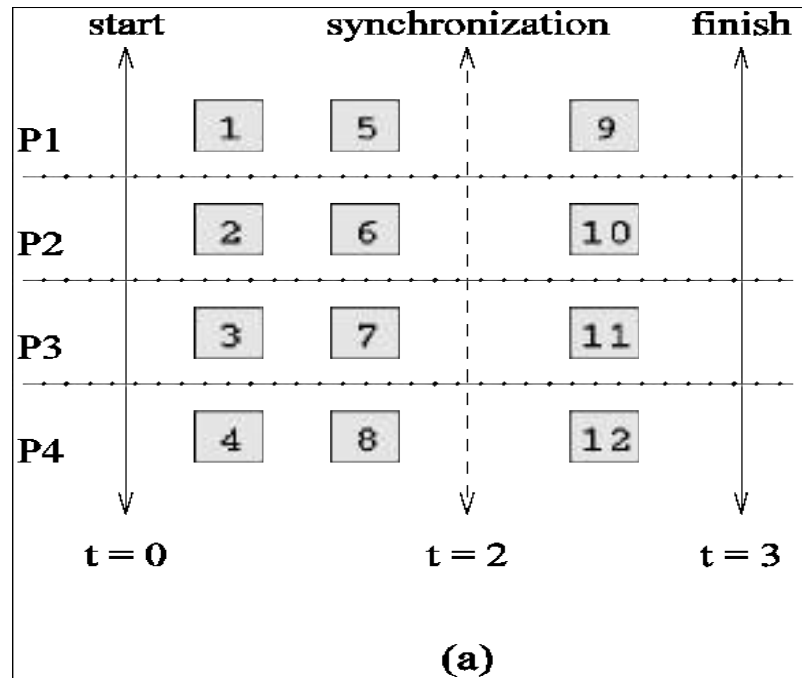


- Once a problem has been decomposed into concurrent tasks, these must be mapped to processes (that can be executed on a parallel platform).
- Mappings must minimize overheads.
- Primary overheads are communication and idling.
- Minimizing these overheads often represents contradicting objectives.
- Assigning all work to one processor trivially minimizes communication at the expense of significant idling.

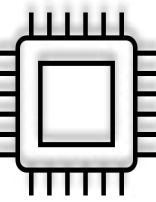


Mapping Techniques for Minimum Idling

Mapping must simultaneously minimize idling and load balance. Merely balancing load does not minimize idling.



Mapping Techniques for Minimum Idling

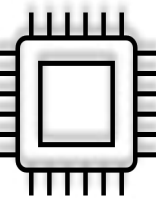


Mapping techniques can be static or dynamic.

- **Static Mapping:** Tasks are mapped to processes a-priori. For this to work, we must have a good estimate of the size of each task. Even in these cases, the problem may be NP complete.
- **Dynamic Mapping:** Tasks are mapped to processes at runtime. This may be because the tasks are generated at runtime, or that their sizes are not known.

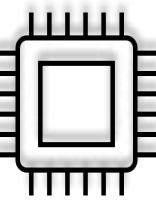
Other factors that determine the choice of techniques include the size of data associated with a task and the nature of underlying domain.

Schemes for Static Mapping



- Mappings based on data partitioning
- Mappings based on task graph partitioning
- Mapping based on task partitioning
- Hierarchical mapping

Mappings Based on Data Partitioning



We can combine data partitioning with the "owner-computes" rule to partition the computation into subtasks. The simplest data decomposition schemes for dense matrices are 1-D block distribution schemes.

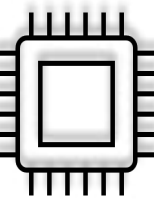
row-wise distribution

P_0
P_1
P_2
P_3
P_4
P_5
P_6
P_7

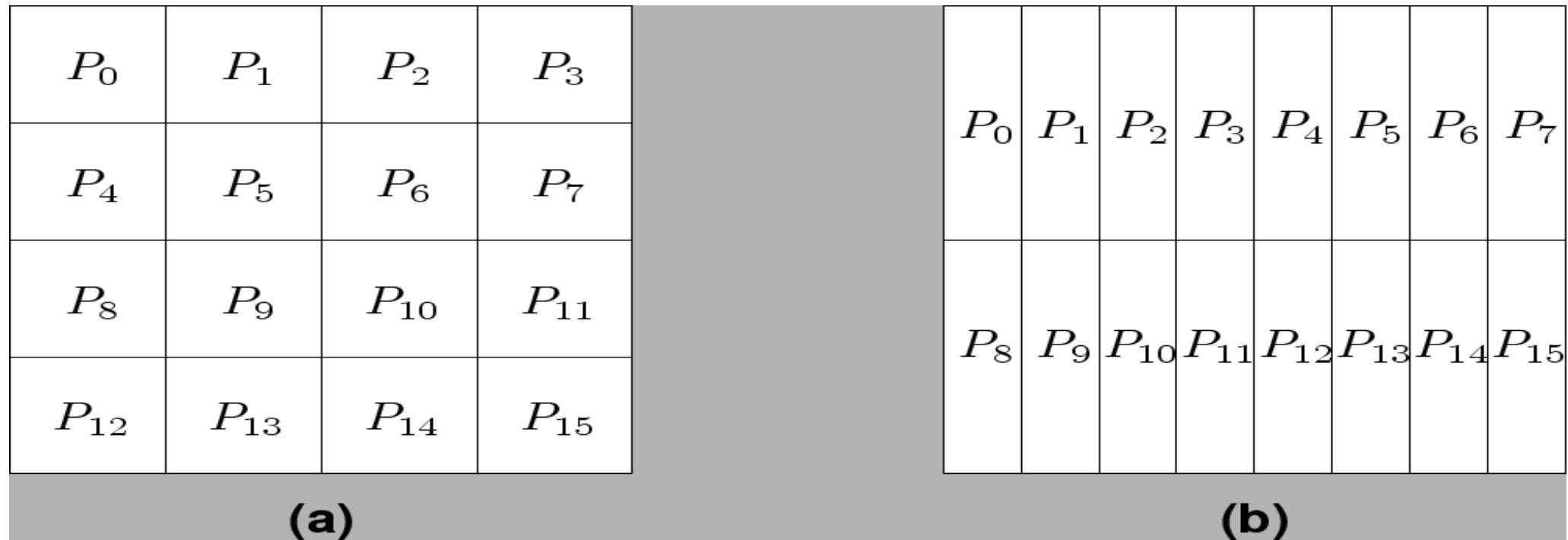
column-wise distribution

P_0	P_1	P_2	P_3	P_4	P_5	P_6	P_7
-------	-------	-------	-------	-------	-------	-------	-------

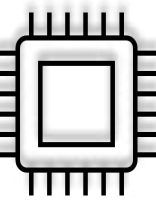
Block Array Distribution Schemes



Block distribution schemes can be generalized to higher dimensions as well.

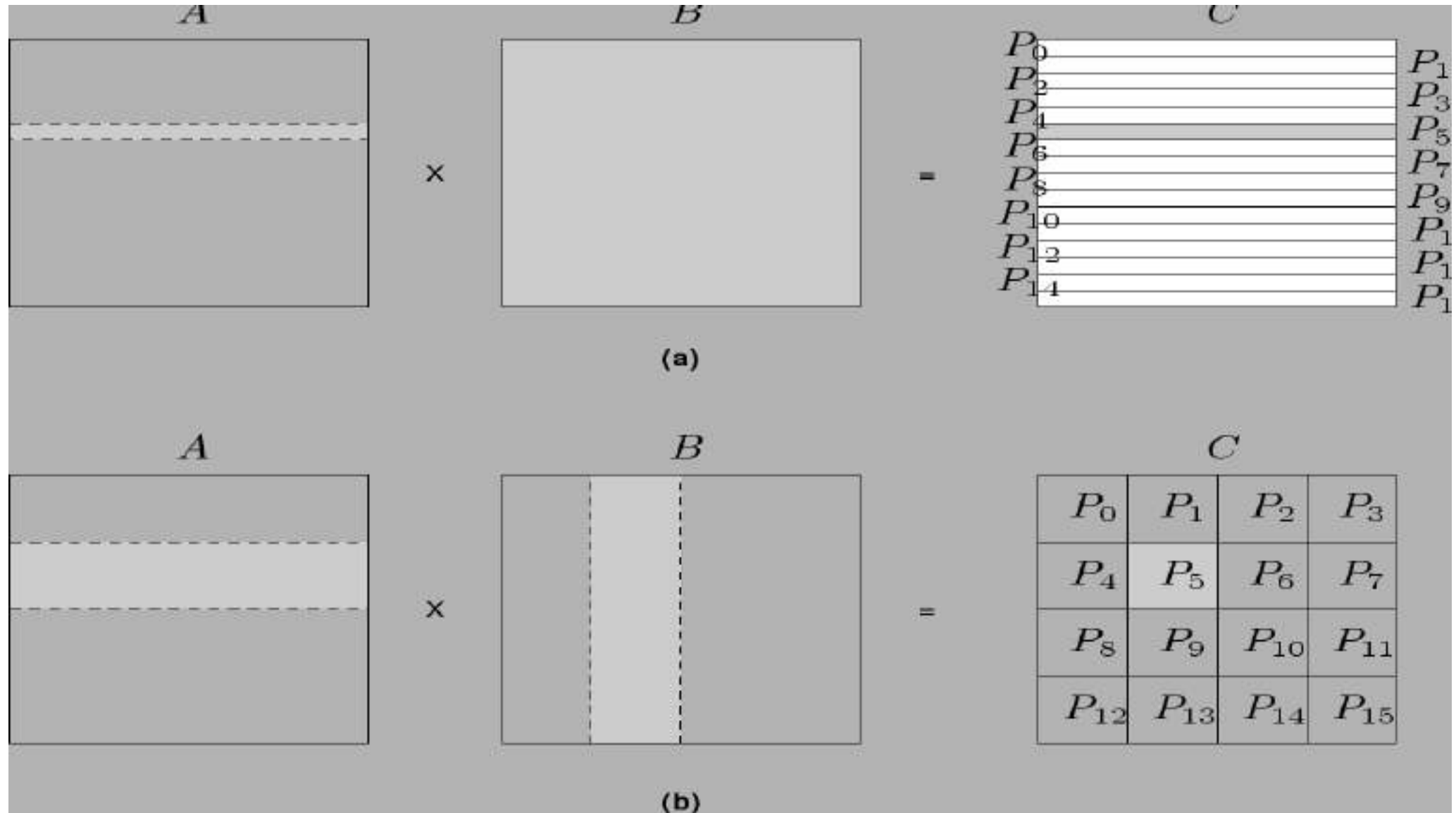
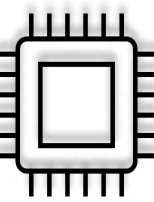


Block Array Distribution Schemes: Examples

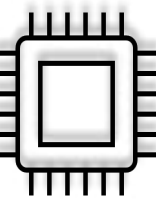


- For multiplying two dense matrices \mathbf{A} and \mathbf{B} , we can partition the output matrix \mathbf{C} using a block decomposition.
- For load balance, we give each task the same number of elements of \mathbf{C} . (Note that each element of \mathbf{C} corresponds to a single dot product.)
- The choice of precise decomposition (1-D or 2-D) is determined by the associated communication overhead.
- In general, higher dimension decomposition allows the use of larger number of processes.

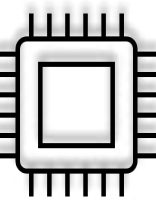
Data Sharing in Dense Matrix Multiplication



Cyclic and Block Cyclic Distributions



- If the amount of computation associated with data items varies, a block decomposition may lead to significant load imbalances.
- A simple example of this is in LU decomposition (or Gaussian Elimination) of dense matrices.



LU Factorization of a Dense Matrix

A decomposition of LU factorization into 14 tasks - notice the significant load imbalance.

$$\begin{pmatrix} A_{1,1} & A_{1,2} & A_{1,3} \\ A_{2,1} & A_{2,2} & A_{2,3} \\ A_{3,1} & A_{3,2} & A_{3,3} \end{pmatrix} \rightarrow \begin{pmatrix} L_{1,1} & 0 & 0 \\ L_{2,1} & L_{2,2} & 0 \\ L_{3,1} & L_{3,2} & L_{3,3} \end{pmatrix} \cdot \begin{pmatrix} U_{1,1} & U_{1,2} & U_{1,3} \\ 0 & U_{2,2} & U_{2,3} \\ 0 & 0 & U_{3,3} \end{pmatrix}$$

1: $A_{1,1} \rightarrow L_{1,1}U_{1,1}$

2: $L_{2,1} = A_{2,1}U_{1,1}^{-1}$

3: $L_{3,1} = A_{3,1}U_{1,1}^{-1}$

4: $U_{1,2} = L_{1,1}^{-1}A_{1,2}$

5: $U_{1,3} = L_{1,1}^{-1}A_{1,3}$

6: $A_{2,2} = A_{2,2} - L_{2,1}U_{1,2}$

7: $A_{3,2} = A_{3,2} - L_{3,1}U_{1,2}$

8: $A_{2,3} = A_{2,3} - L_{2,1}U_{1,3}$

9: $A_{3,3} = A_{3,3} - L_{3,1}U_{1,3}$

10: $A_{2,2} \rightarrow L_{2,2}U_{2,2}$

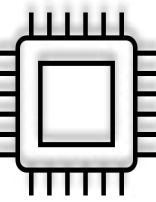
11: $L_{3,2} = A_{3,2}U_{2,2}^{-1}$

12: $U_{2,3} = L_{2,2}^{-1}A_{2,3}$

13: $A_{3,3} = A_{3,3} - L_{3,2}U_{2,3}$

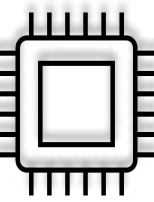
14: $A_{3,3} \rightarrow L_{3,3}U_{3,3}$

Block Cyclic Distributions

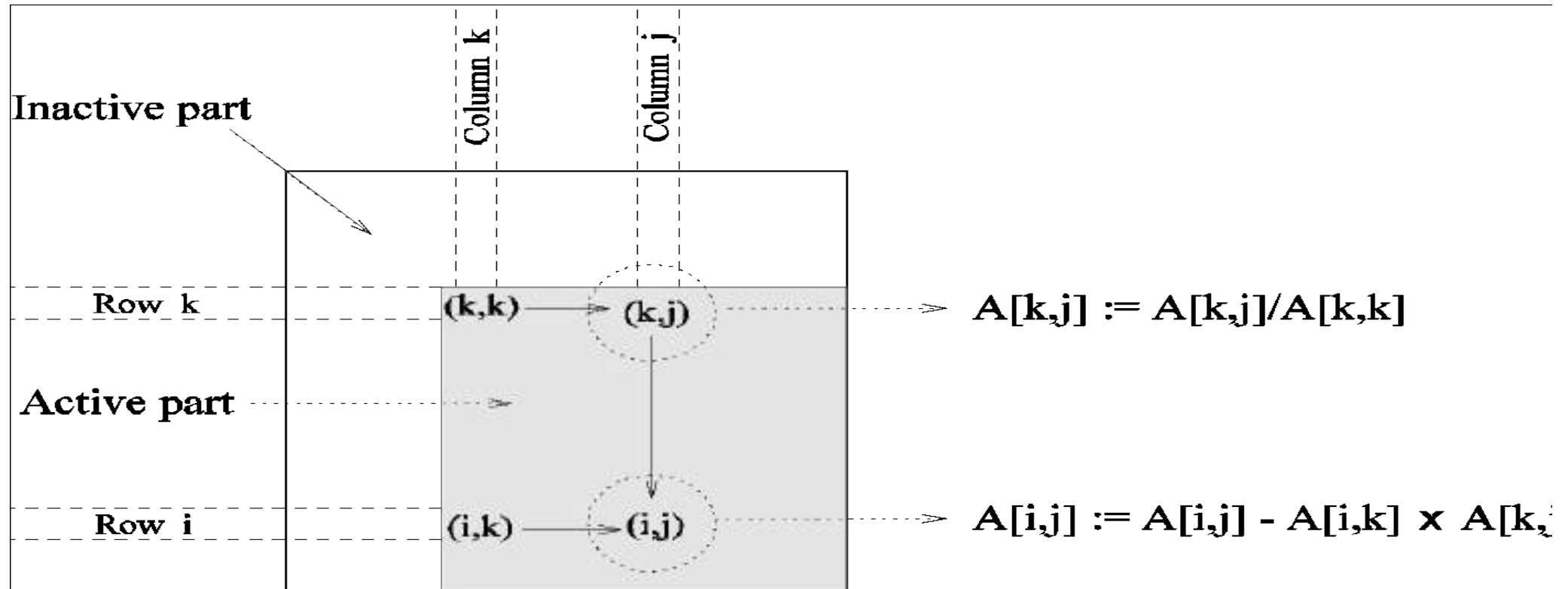


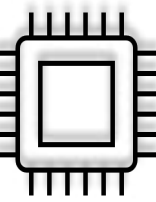
- Variation of the block distribution scheme that can be used to alleviate the load-imbalance and idling problems.
- Partition an array into many more blocks than the number of available processes.
- Blocks are assigned to processes in a round-robin manner so that each process gets several non-adjacent blocks.

Block-Cyclic Distribution for Gaussian Elimination



The active part of the matrix in Gaussian Elimination changes. By assigning blocks in a block-cyclic fashion, each processor receives blocks from different parts of the matrix.

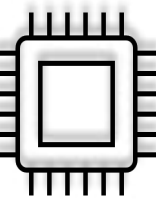




Block-Cyclic Distribution: Examples

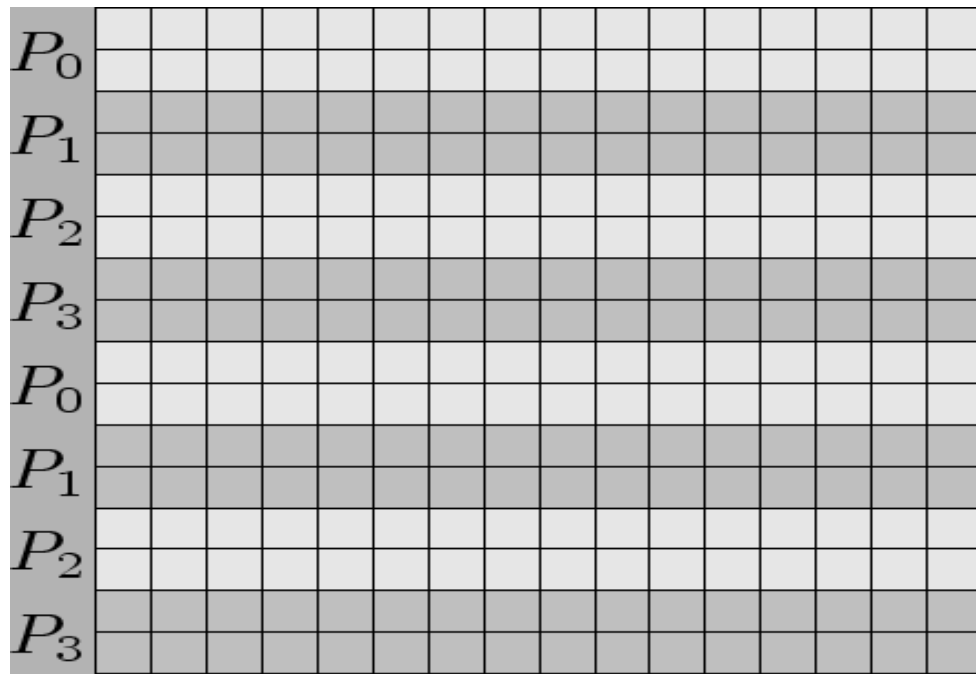
One- and two-dimensional block-cyclic distributions among 4 processes.

P₀ T ₁	P₃ T ₄	P₆ T ₅
P₁ T ₂	P₄ T ₆ T ₁₀	P₇ T ₈ T ₁₂
P₂ T ₃	P₅ T ₇ T ₁₁	P₈ T ₉ T ₁₃ T ₁₄

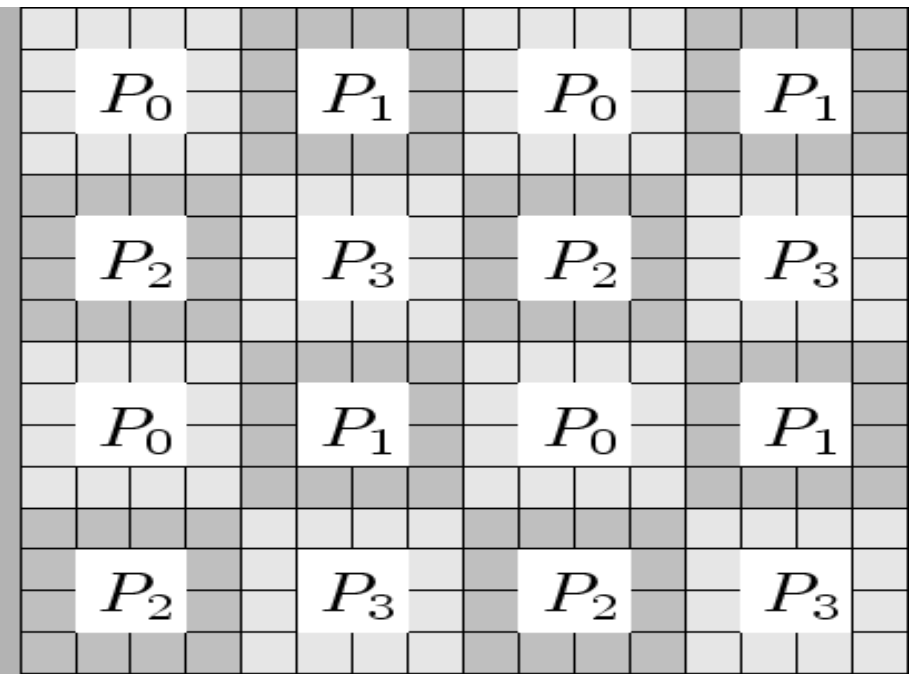


Block-Cyclic Distribution

- A cyclic distribution is a special case in which block size is one.
- A block distribution is a special case in which block size is n/p , where n is the dimension of the matrix and p is the number of processes.

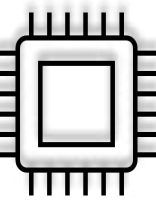


(a)



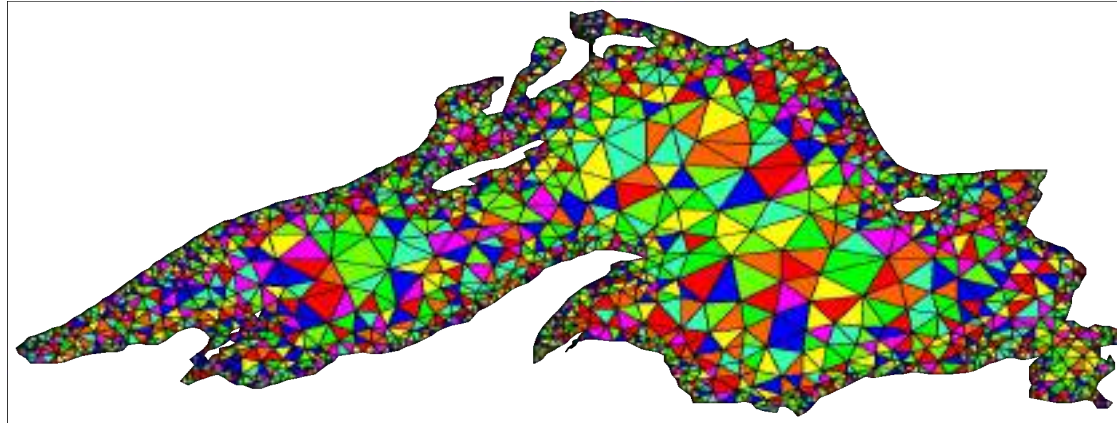
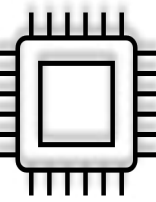
(b)

Graph Partitioning based Data Decomposition

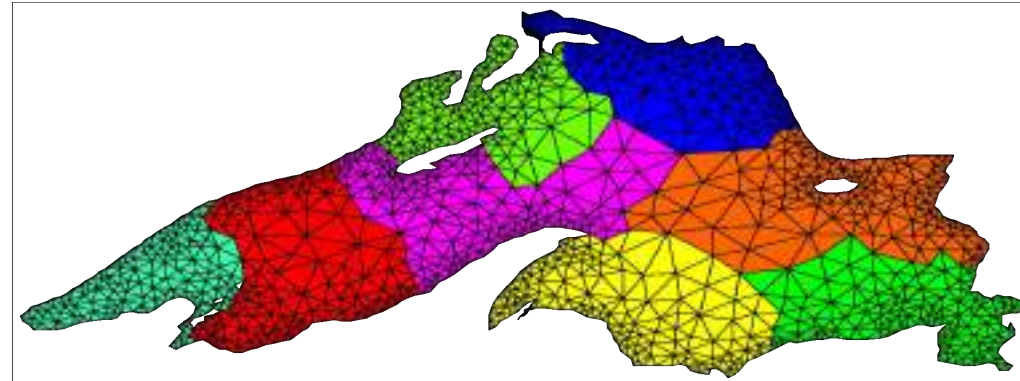


- In case of sparse matrices, block decompositions are more complex.
- Consider the problem of multiplying a sparse matrix with a vector.
- The graph of the matrix is a useful indicator of the work (number of nodes) and communication (the degree of each node).
- In this case, we would like to partition the graph so as to assign equal number of nodes to each process, while minimizing edge count of the graph partition.

Partitioning the Graph of Lake Superior

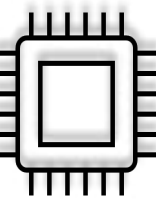


Random Partitioning



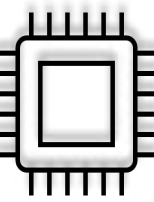
Partitioning for minimum edge-cut.

Mappings Based on Task Partitioning

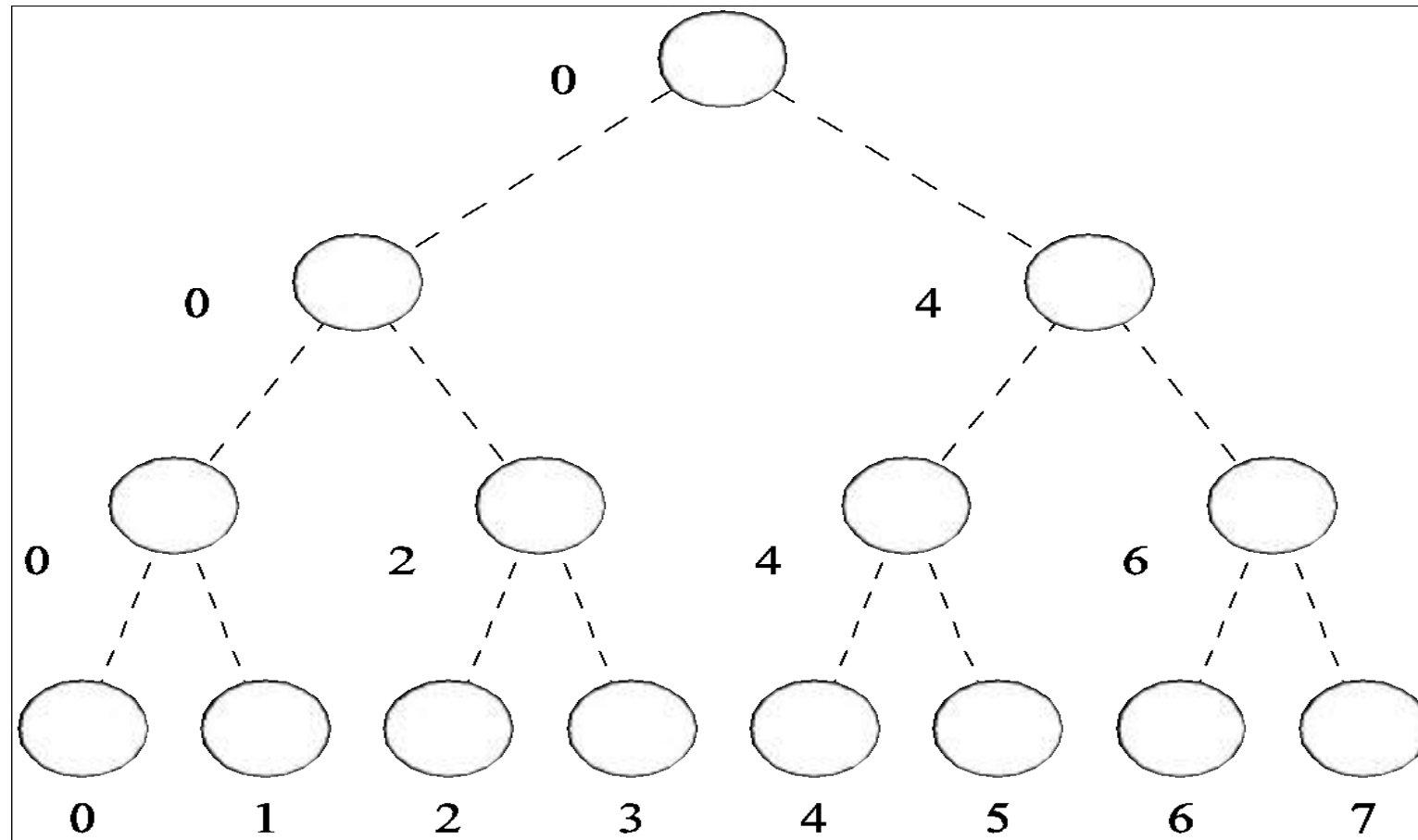


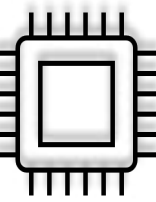
- Partitioning a given task-dependency graph across processes.
- Determining an optimal mapping for a general task-dependency graph is an NP-complete problem.
- Excellent heuristics exist for structured graphs.

Task Partitioning: Mapping a Binary Tree Dependency Graph



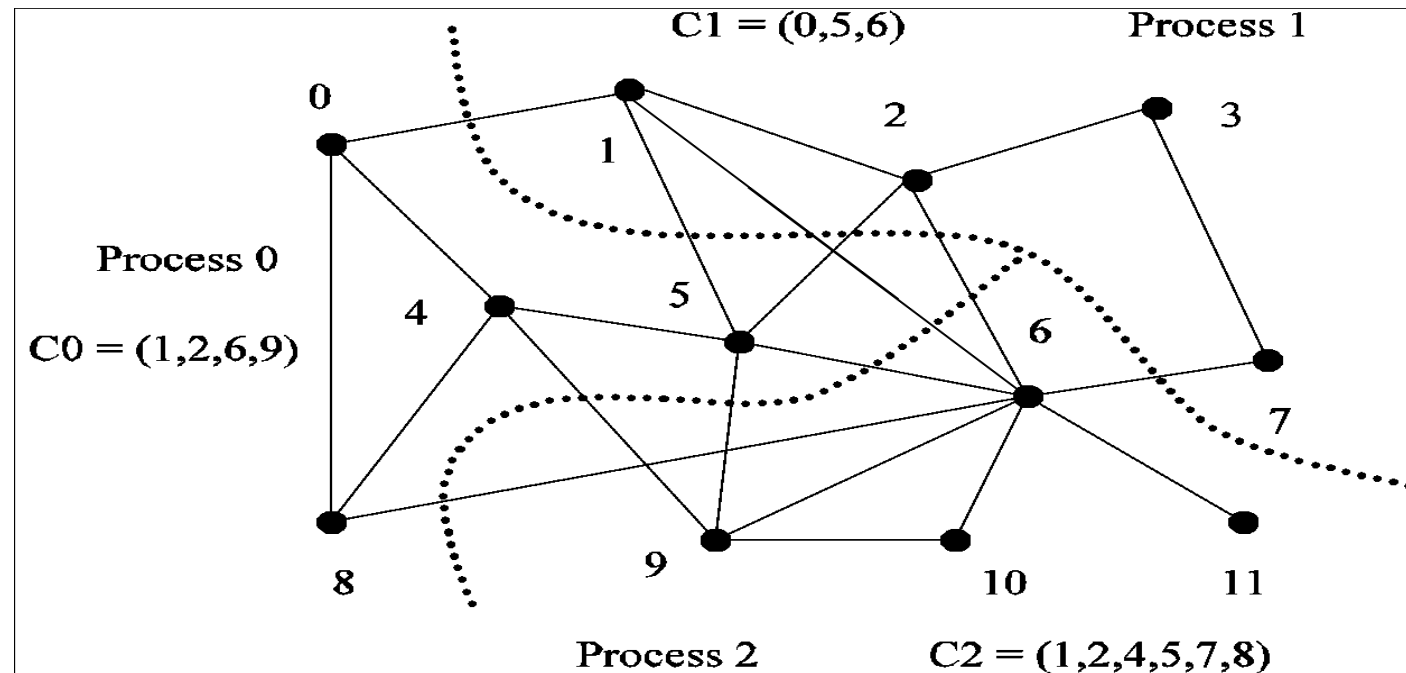
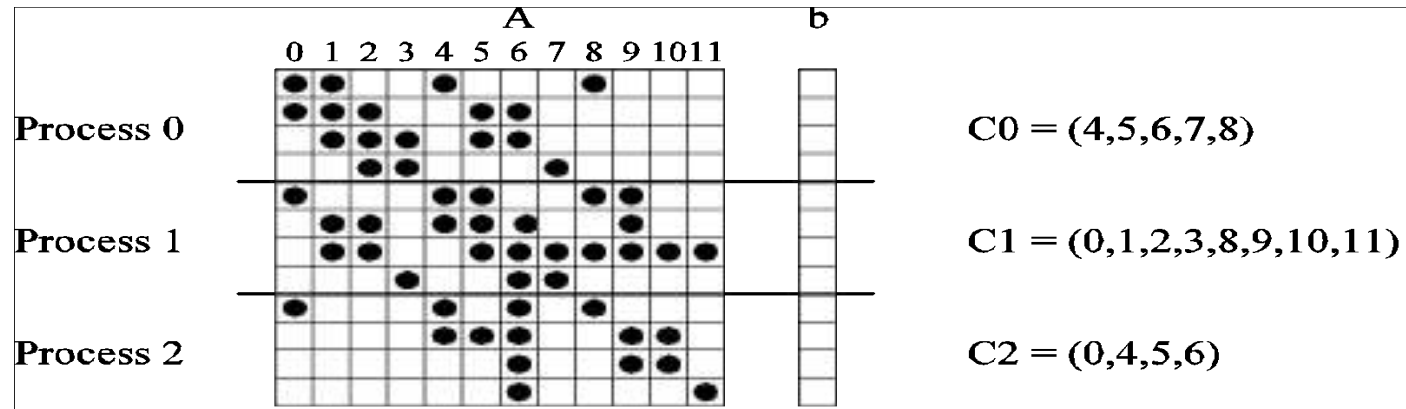
Example illustrates the dependency graph of one view of quick-sort and how it can be assigned to processes in a hypercube.



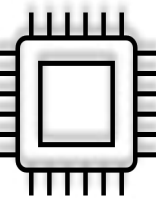


Task Partitioning: Mapping a Sparse Graph

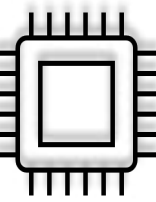
Sparse graph for computing a sparse matrix-vector product and its mapping.



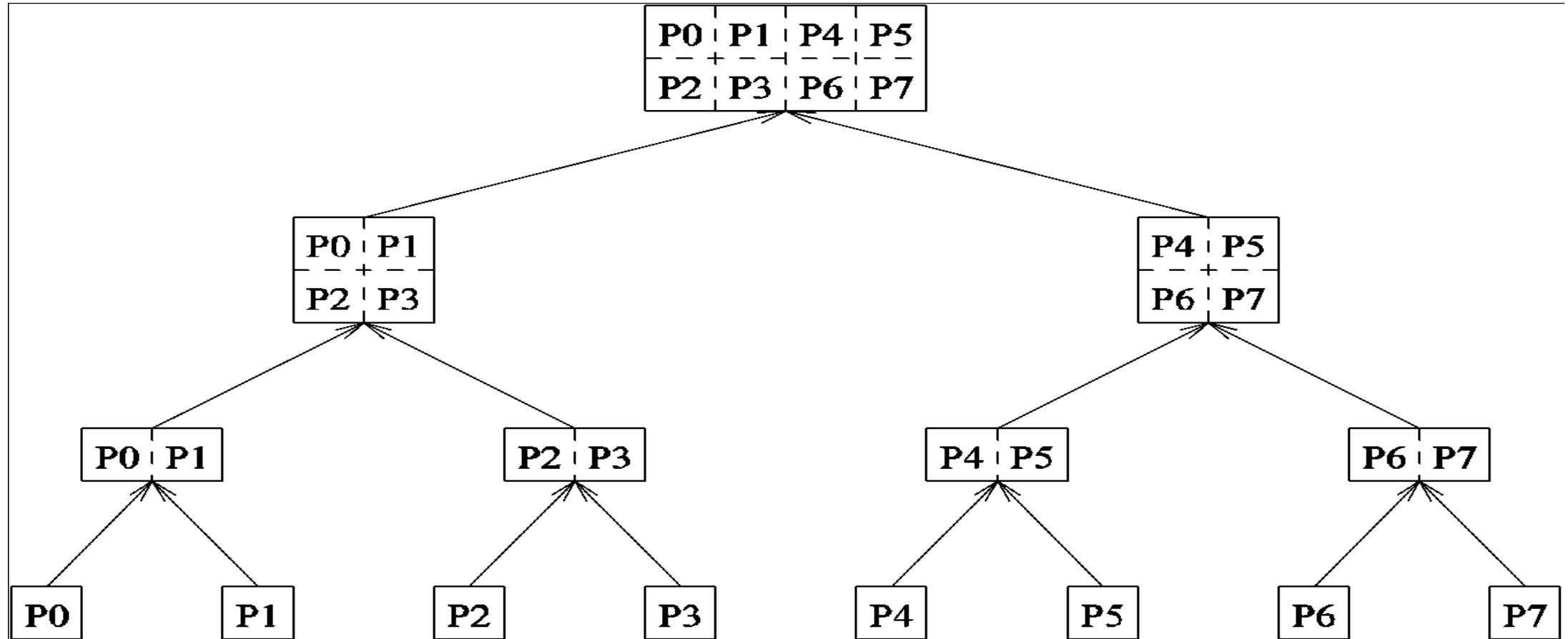
Hierarchical Mappings



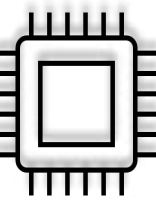
- Sometimes a single mapping technique is inadequate.
- For example, the task mapping of the binary tree (quicksort) cannot use a large number of processors.
- For this reason, task mapping can be used at the top level and data partitioning within each level.



An example of task partitioning at top level with data partitioning at the lower level.

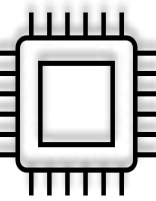


Schemes for Dynamic Mapping



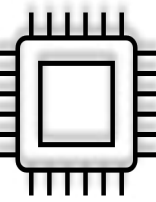
- **Dynamic mapping is sometimes also referred to as dynamic load balancing, since load balancing is the primary motivation for dynamic mapping.**
- **Dynamic mapping schemes can be centralized or distributed.**

Centralized Dynamic Mapping



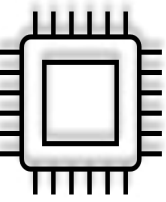
- Processes are designated as masters or slaves.
- When a process runs out of work, it requests the master for more work.
- In an example of sorting rows of matrix, there can be a central pool of indices containing the rows to be sorted. When a process is idle it can pick up a row by deleting its indice from the pool and sort it. This method of scheduling independent iterations of a loop among processes is known as **Self scheduling**.
- When the number of processes increases, the master may become the bottleneck.
- To alleviate this, a process may pick up a number of tasks (a chunk) at one time. This is called **Chunk scheduling**.
- Selecting large chunk sizes may lead to significant load imbalances as well.
- A number of schemes have been used to gradually decrease chunk size as the computation progresses.

Distributed Dynamic Mapping

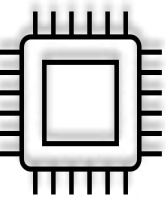


- Each process can send or receive work from other processes.
- This alleviates the bottleneck in centralized schemes.
- There are four critical questions: how are sensing and receiving processes paired together, who initiates work transfer, how much work is transferred, and when is a transfer triggered?
- Answers to these questions are generally application specific.

You have completed this topic, you should be able to:



- Draw task dependency graph and determine the average degree of concurrency?
- Describe the terms:
 - a)Task
 - b)Decomposition
 - c)Task Dependency graph
 - d)Task Interaction Graph
 - e)Process
 - f)Mapping
 - g)Critical path
 - h)degree of concurrence
- Explain different decomposition techniques?
- Explain characteristics of task and task interactions?
- Explain the static and dynamic mapping load balancing techniques?



References Used

- **Ananth Grama, Anshul Gupta, George Karypis, Vipin Kumar ,
—Introduction to Parallel Computing, Pearson Education, Second
Edition, 2007.**
- https://www.tutorialspoint.com/parallel_algorithm/parallel_algorithm_models.html