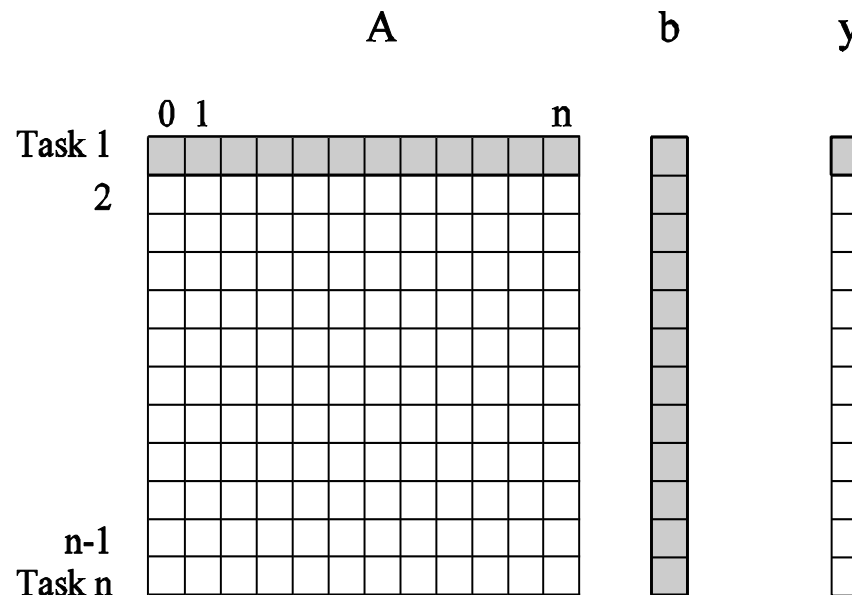


# Fundamentals of parallel programming

# 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 in many different ways.
- A decomposition can be illustrated in the form of a directed 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.

**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 I

Consider the execution of the query:

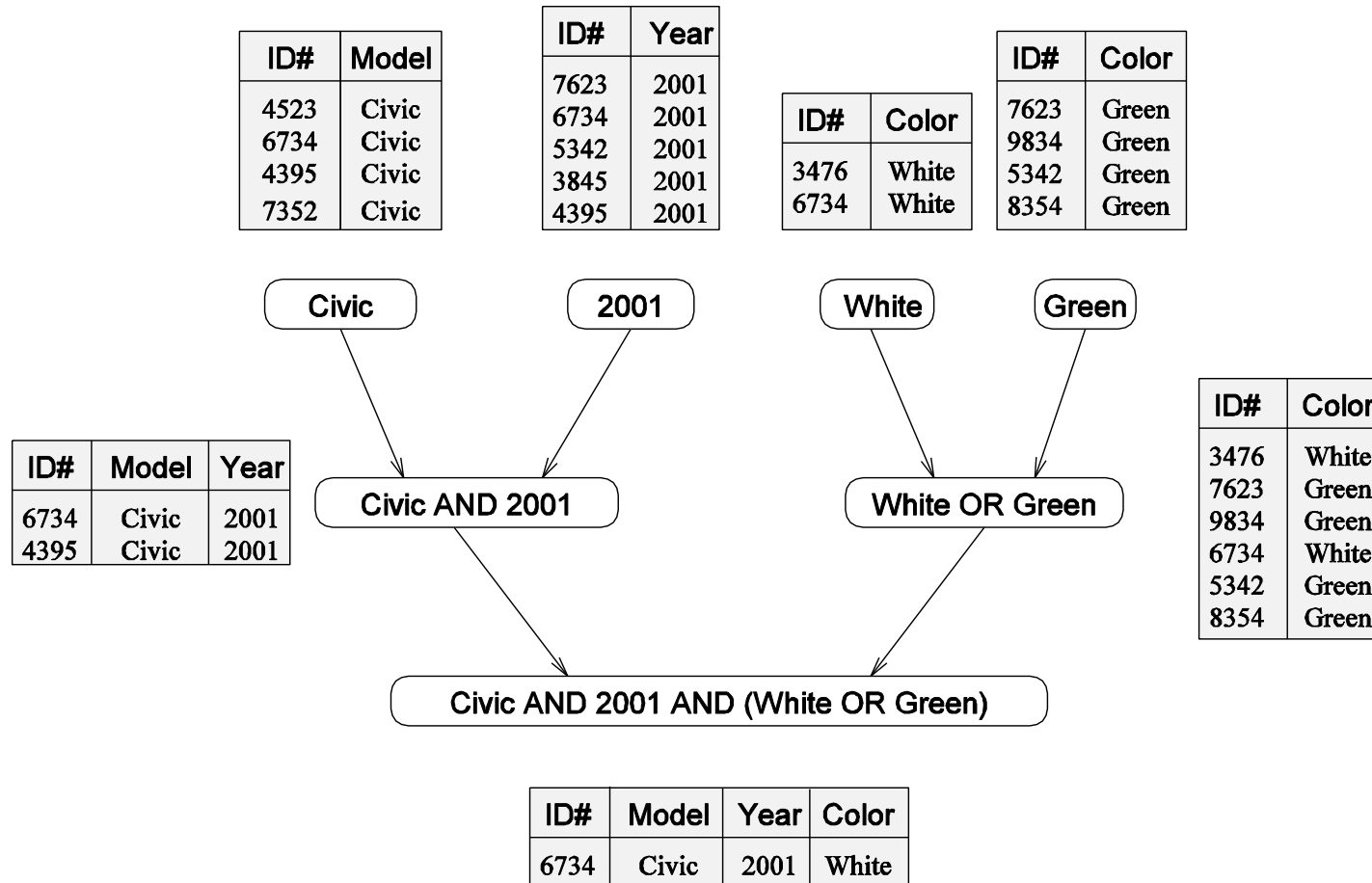
MODEL = "CIVIC" AND YEAR = 2001 AND  
(COLOR = "GREEN" OR COLOR = "WHITE")

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 II

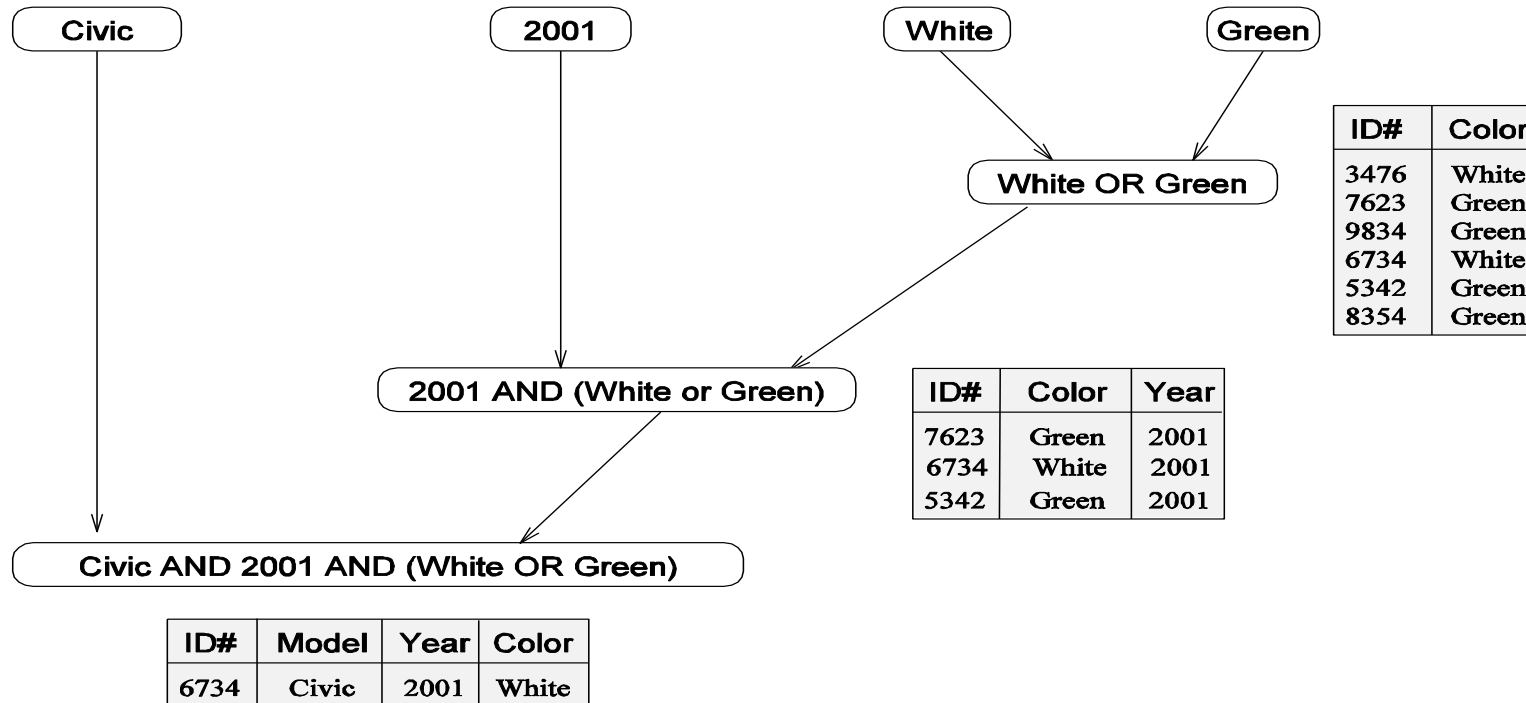
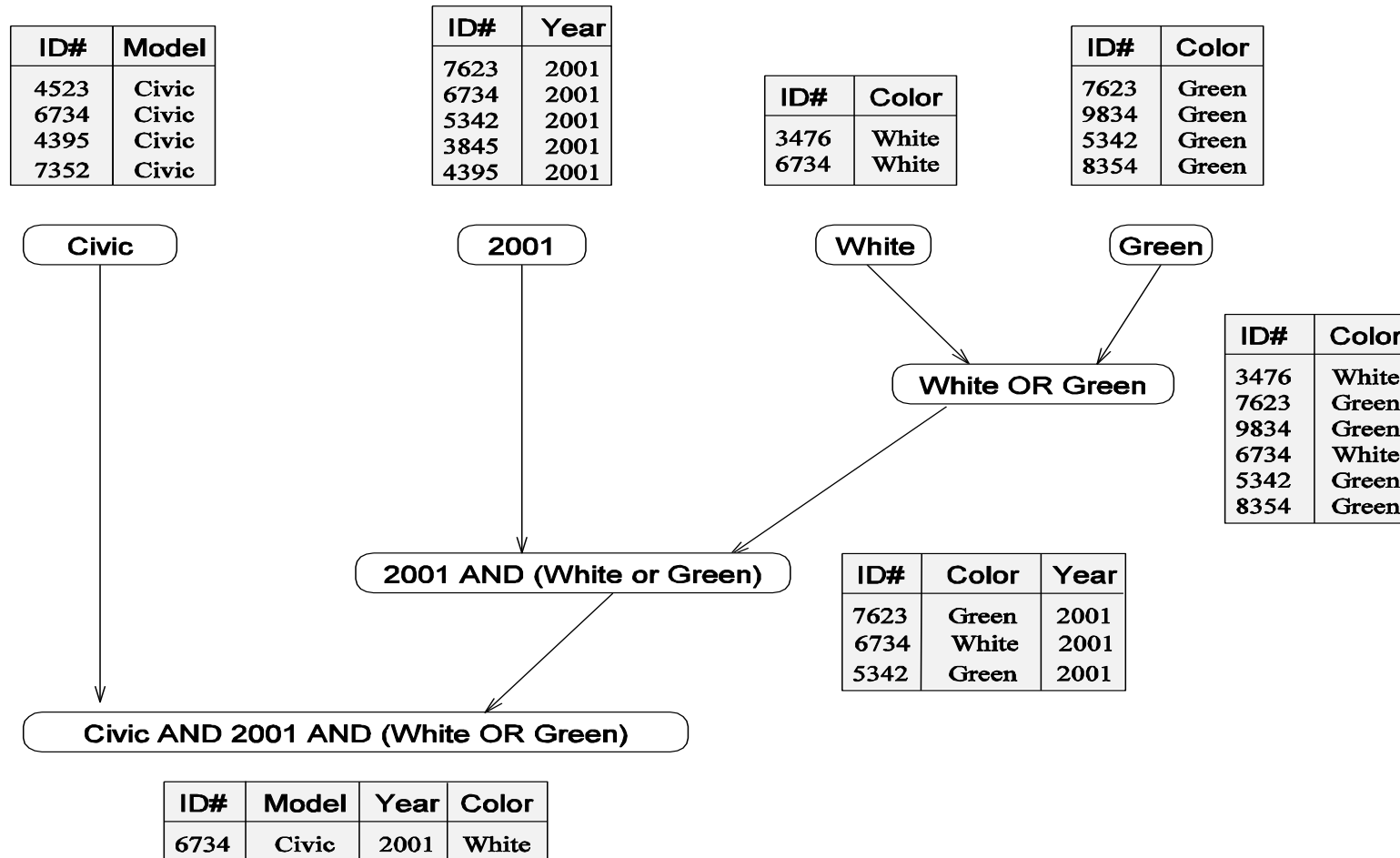
This query can be divided into subtasks in various ways.



➔ Edges denote that the output of one task is needed to accomplish the next.

# Example: Database Query Processing III

The query can be decomposed in other ways as well:



ID#	Model	Year	Color
6734	Civic	2001	White

➔ Different task decompositions may lead to significant differences with respect to their eventual parallel performance.

# Data Dependencies

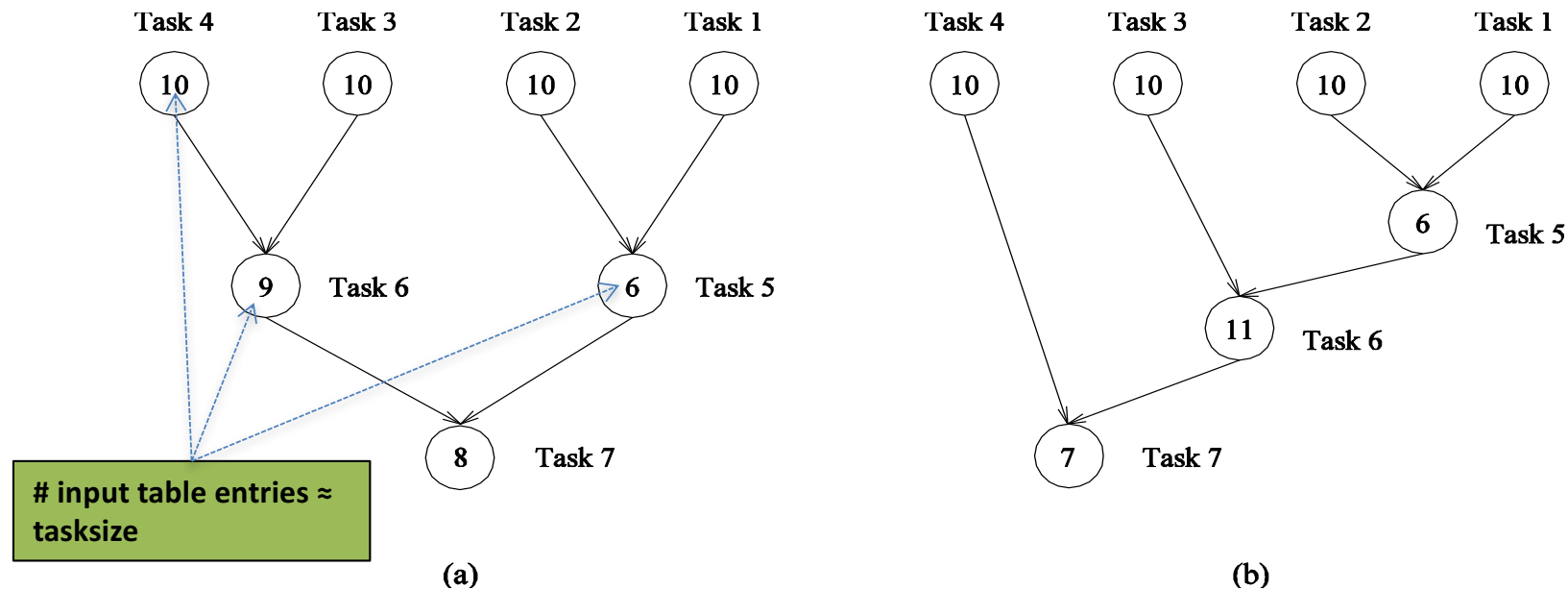
- **True(Flow)-Dependence:** Task 1 computes a value stored at A, and Task 2 retrieves a value stored at A.
- **Anti-Dependence:** Task 1 retrieves a value stored at A, and Task 2 computes a value stored at A
- **Input-Dependence:** Task 1 and Task 2 both retrieve a value stored at A
- **Output-Dependence:** Task 1 and Task 2 both compute a value stored at A

# Critical Path Length

- A **directed path** in the task dependency graph represents a sequence of tasks that must be processed **one after the other**.
- The longest path determines the shortest time in which the program can be executed in parallel.
- The **length of the longest path** in a task dependency graph is called the critical path length.



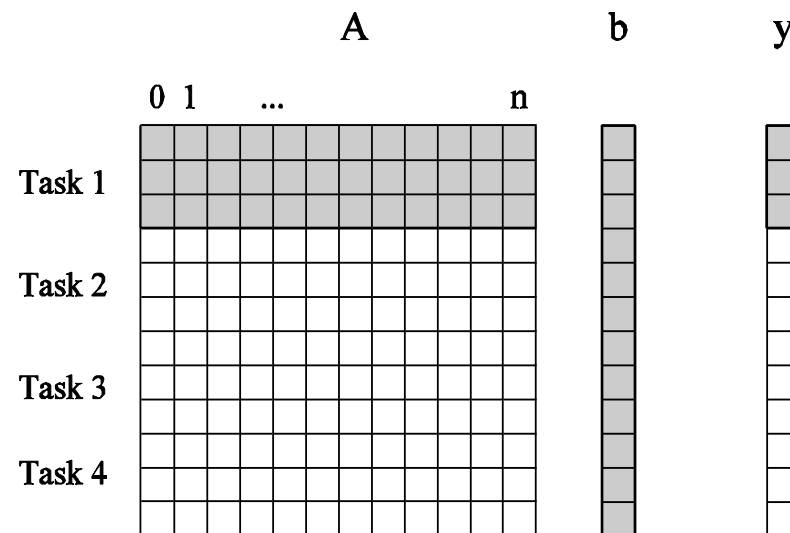
Consider the task dependency graphs of the two database query decompositions:



- ➔ What are the critical path lengths for the two task dependency graphs?
- ➔ If each task takes 10 time units, what is the shortest parallel execution time for each decomposition?
- ➔ How many processors are needed in each case to achieve this minimum parallel execution time?

# 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**.



A coarse grained counterpart to the dense matrix-vector product example.

- ➔ Each task 3 times as big
- ➔ The number of tasks is 3 times less

# 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.

# Limits on Parallel Performance

It would appear that the parallel time can be made arbitrarily small by making the decomposition finer in granularity. However:

- 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  $(n^2)$  concurrent tasks.*
- 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.

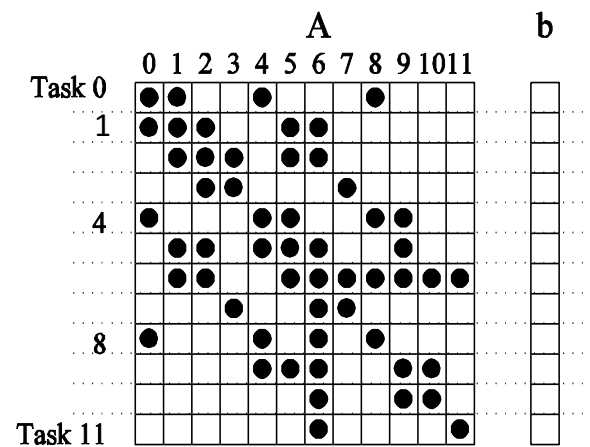
# Task Interaction Graphs

- 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  $b$  is not replicated across all tasks, they will have to communicate elements of the vector  $b$ .
- 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.
- In fact task interaction graphs represent **input- and output-dependencies**, whereas task dependence graphs represent **true-dependencies** or **anti-dependencies**.

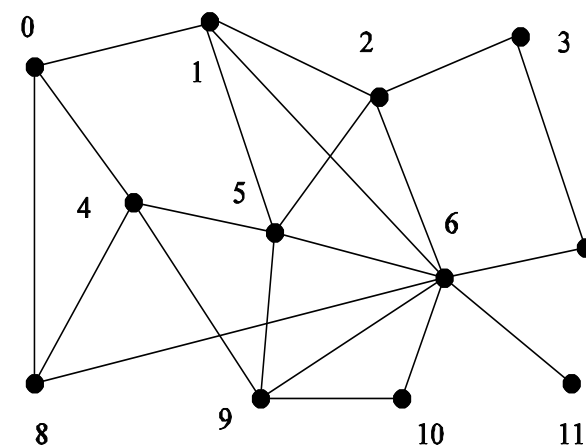
# Task Interaction Graphs: An Example

Consider the problem of multiplying a **sparse** matrix  $\mathbf{A}$  with a vector  $\mathbf{b}$ . The following observations can be made:

- As before, the computation of each element of the result vector can be viewed as an **independent** task. So no true or anti dependencies.
- Unlike a dense matrix-vector product though, only **non-zero elements** of matrix  $\mathbf{A}$  participate in the computation.
- If, for memory optimality, we also partition  $\mathbf{b}$  across tasks, then one can see that the task interaction graph of the computation is identical to the graph of the matrix  $\mathbf{A}$  (the graph for which  $\mathbf{A}$  represents the adjacency structure). **This is not always the case!!**



(a)



(b)

# Task Interaction Graphs, Granularity, and Associated Communication Overhead

In general: if the granularity of a decomposition is finer, then the associated (communication) overhead increases.

**Example:** Consider the sparse matrix-vector product example from previous slide. Assume that each node takes unit time to process and each interaction (edge) causes an overhead of a unit time.

Viewing node 0 as an independent task involves a useful computation of **one time unit** and overhead (communication) of **three time units**.

Now, if we consider nodes 0, 4, and 5 as one task, then the task has useful computation totaling to **three time units** and communication corresponding to **five time units** (five edges to 1,2,6,8 and 9). Clearly, this is a more favorable ratio than the former case.

# Decomposition Techniques

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.

- Recursive Decomposition
- Data Decomposition
- Exploratory Decomposition
- Speculative Decomposition

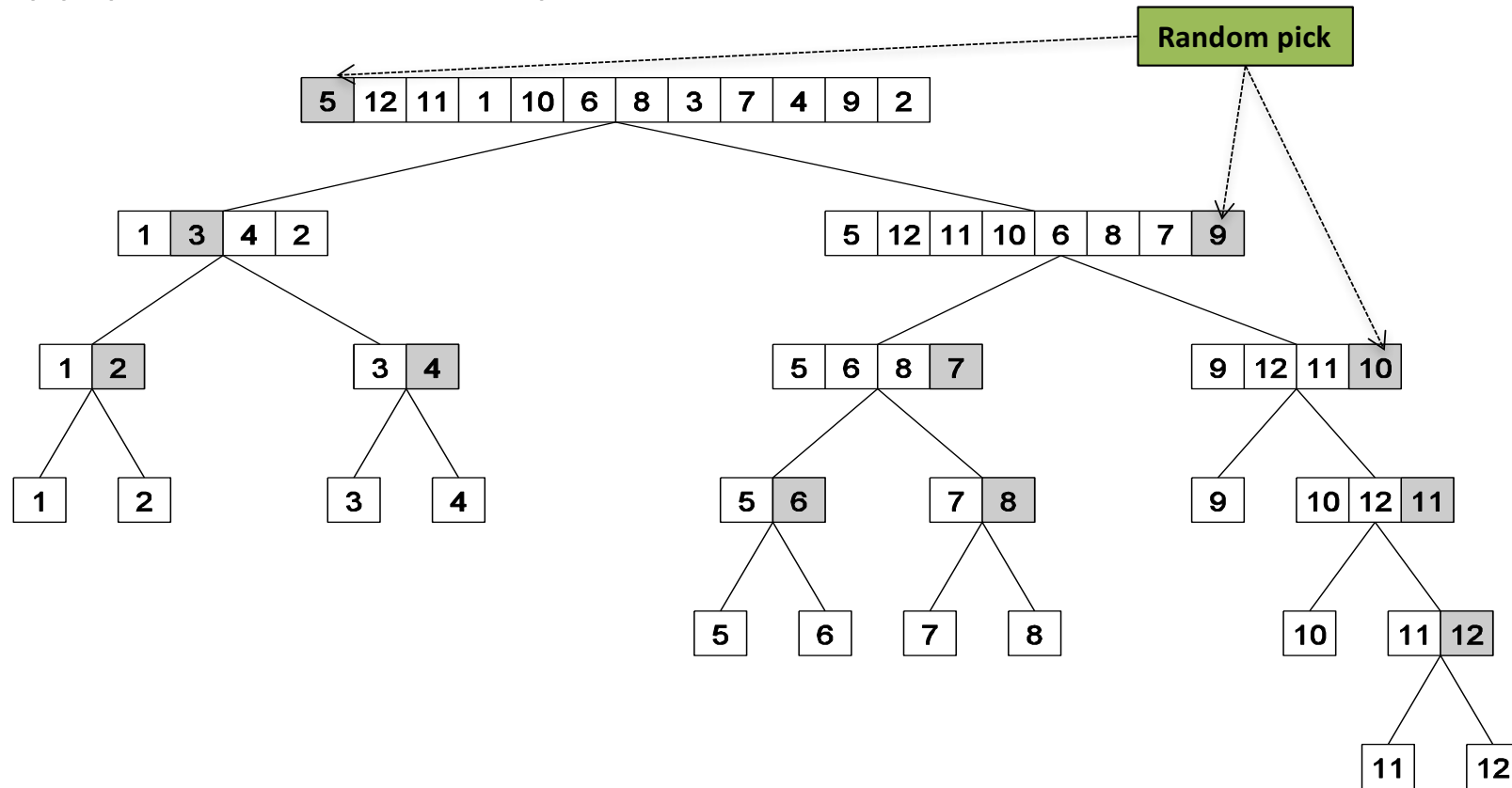


# Recursive Decomposition

- 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.

# Recursive Decomposition: Example

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 sub-list can be processed concurrently (i.e., each sub-list represents an independent subtask). This can be repeated **recursively**.

# Data Decomposition

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

# Output Data Decomposition: Example

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

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

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

# Output Data Decomposition: Example

A partitioning of output data does **not** result in a **unique** decomposition into tasks. For example, for the same problem as in previous slide, with identical output data distribution, we can derive the following two (other) decompositions:

<b>Decomposition I</b>	<b>Decomposition II</b>
Task 1: $\mathbf{C}_{1,1} = \mathbf{A}_{1,1} \mathbf{B}_{1,1}$	Task 1: $\mathbf{C}_{1,1} = \mathbf{A}_{1,1} \mathbf{B}_{1,1}$
Task 2: $\mathbf{C}_{1,1} = \mathbf{C}_{1,1} + \mathbf{A}_{1,2} \mathbf{B}_{2,1}$	Task 2: $\mathbf{C}_{1,1} = \mathbf{C}_{1,1} + \mathbf{A}_{1,2} \mathbf{B}_{2,1}$
Task 3: $\mathbf{C}_{1,2} = \mathbf{A}_{1,1} \mathbf{B}_{1,2}$	Task 3: $\mathbf{C}_{1,2} = \mathbf{A}_{1,2} \mathbf{B}_{2,2}$
Task 4: $\mathbf{C}_{1,2} = \mathbf{C}_{1,2} + \mathbf{A}_{1,2} \mathbf{B}_{2,2}$	Task 4: $\mathbf{C}_{1,2} = \mathbf{C}_{1,2} + \mathbf{A}_{1,1} \mathbf{B}_{1,2}$
Task 5: $\mathbf{C}_{2,1} = \mathbf{A}_{2,1} \mathbf{B}_{1,1}$	Task 5: $\mathbf{C}_{2,1} = \mathbf{A}_{2,2} \mathbf{B}_{2,1}$
Task 6: $\mathbf{C}_{2,1} = \mathbf{C}_{2,1} + \mathbf{A}_{2,2} \mathbf{B}_{2,1}$	Task 6: $\mathbf{C}_{2,1} = \mathbf{C}_{2,1} + \mathbf{A}_{2,1} \mathbf{B}_{1,1}$
Task 7: $\mathbf{C}_{2,2} = \mathbf{A}_{2,1} \mathbf{B}_{1,2}$	Task 7: $\mathbf{C}_{2,2} = \mathbf{A}_{2,1} \mathbf{B}_{1,2}$
Task 8: $\mathbf{C}_{2,2} = \mathbf{C}_{2,2} + \mathbf{A}_{2,2} \mathbf{B}_{2,2}$	Task 8: $\mathbf{C}_{2,2} = \mathbf{C}_{2,2} + \mathbf{A}_{2,2} \mathbf{B}_{2,2}$

# Input Data Partitioning

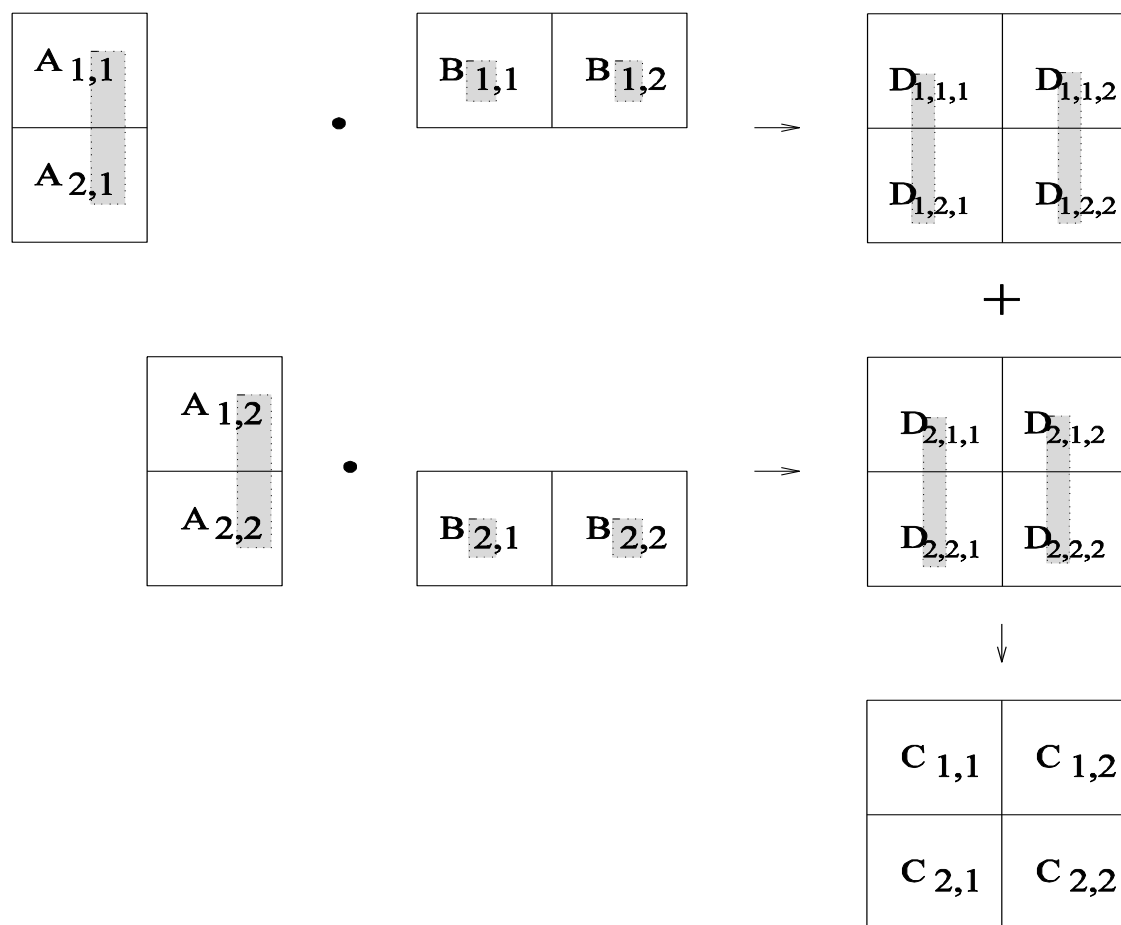
- Generally applicable if each output can be naturally computed as a function of the input.
- In many cases, this is the only natural decomposition because the output is not clearly known a-priori (e.g., the problem of finding the minimum in a list, sorting a given list, etc.).
- A task is associated 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.

# 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.

# Intermediate Data Partitioning: Example

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

Stage I

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

Stage II

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

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

# 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.

# Exploratory Decomposition: Example

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).

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

(a)

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

(b)

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

(c)

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

(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
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	15	11
13	14	12	

task 2

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

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

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

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

task 4

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

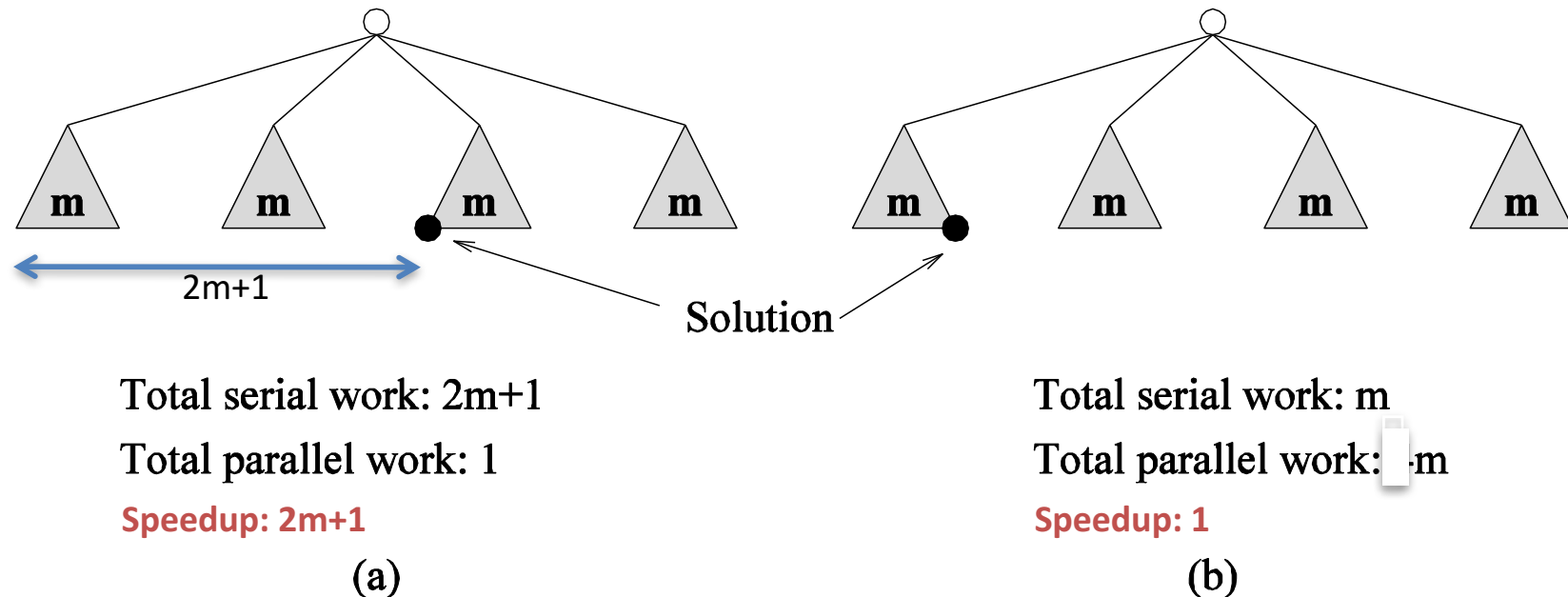
1	2	3	4
5	6	7	8
9	10	11	
13	14	15	12

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

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

# Exploratory Decomposition: Anomalous Computations

- In many instances of exploratory decomposition, the decomposition technique may change the amount of work done by the parallel formulation.
- This change results in **super-** or **sub-linear speedups**.



(b) In case the serial search would have searched the solution space from right to left, with the subtrees interchanged accordingly then Speedup would have been  $1/m$

# Speculative Decomposition

- 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, and, **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.

# Mapping Tasks onto Processes

- In general, the number of tasks in a decomposition exceeds the number of processing elements (Processors, CPU's) available.
- For this reason, a parallel algorithm must also provide a mapping of tasks to **processes**.

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.



# Mapping Tasks onto Processes

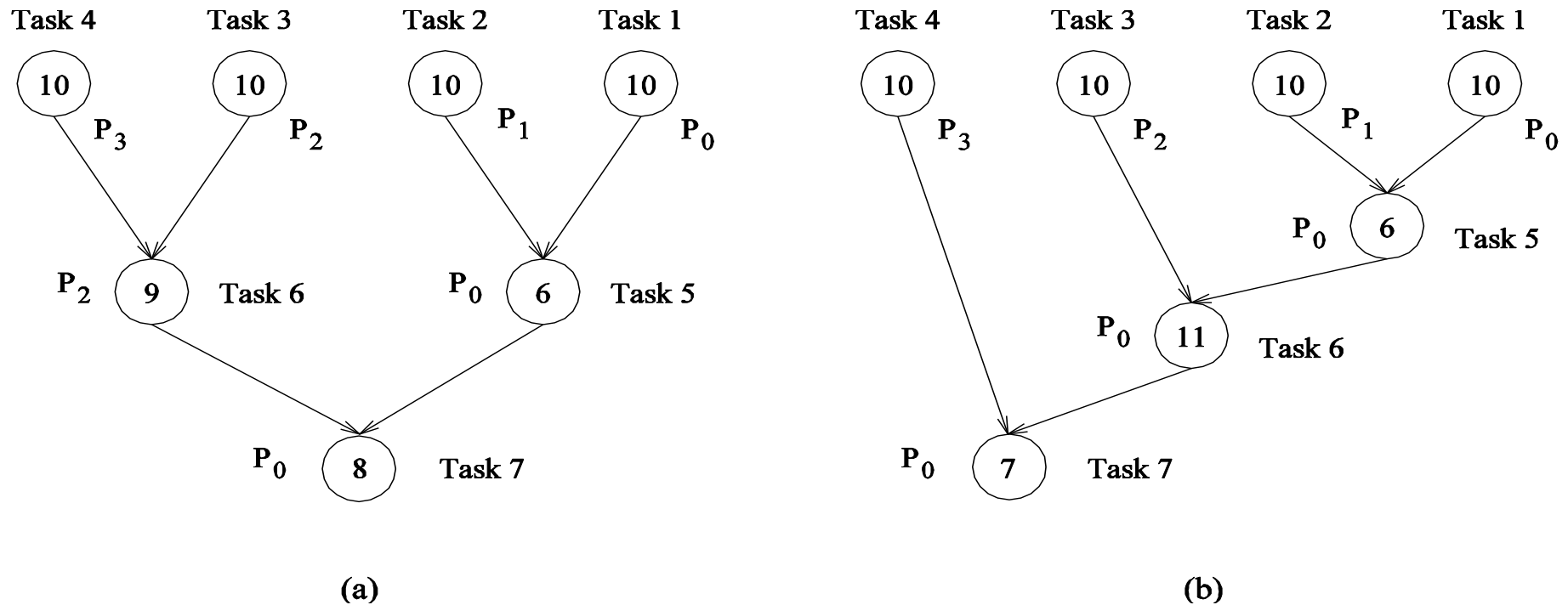
- Appropriate mapping of tasks to processes is critical to the parallel performance of an algorithm.
- Mappings are determined by both the task dependency and task interaction graphs.
- Task dependency graphs can be used to ensure that work is equally spread across all processes at any point (minimum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

# Mapping Tasks onto Processes

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 or to nearby processes.

# Example



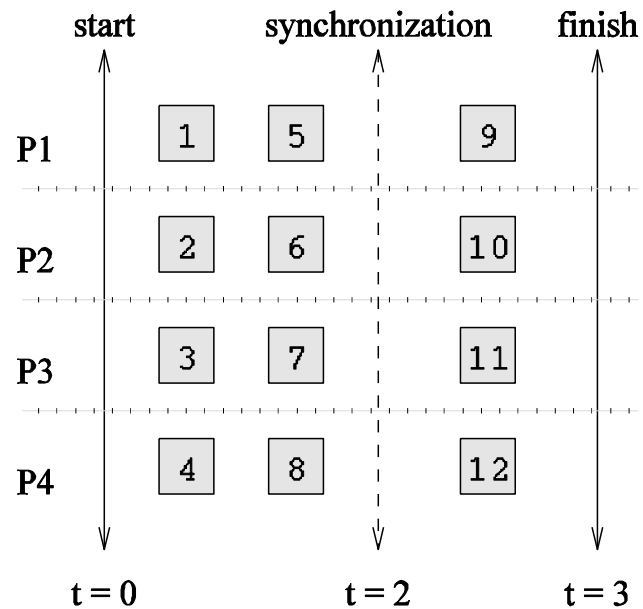
These mappings were arrived at by viewing the dependency graph in terms of levels (**no two nodes in a level have dependencies**). Tasks within a **single level** are then assigned to **different processes**.

# Mapping Trade-offs

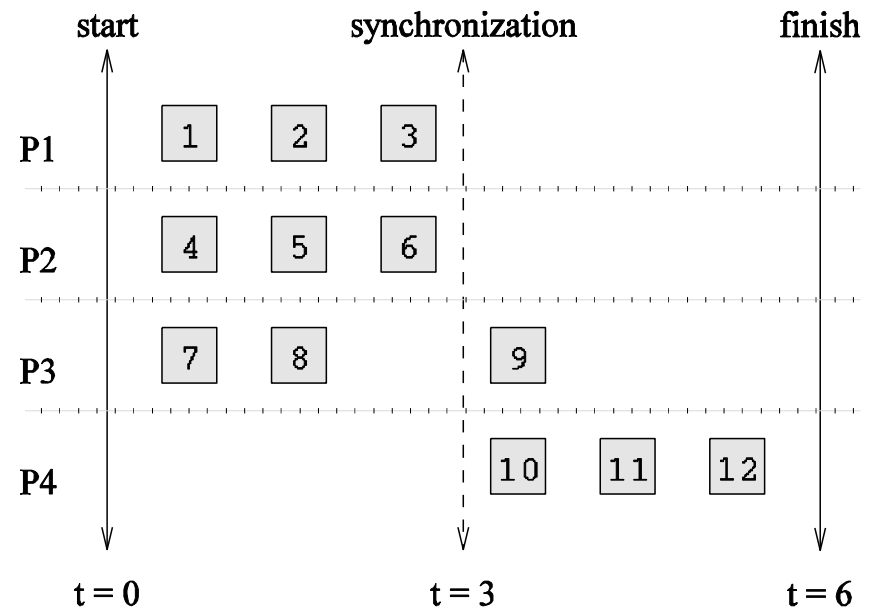
- Mappings must **minimize overheads**.
- Primary overheads are **communication** and **idling**.
- Minimizing these overheads often represents **contradicting objectives**.
- Assigning all work to one processor trivially minimizes communication at the expense of significant idling.

# Mapping Techniques for Minimum Idling

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



(a)



(b)

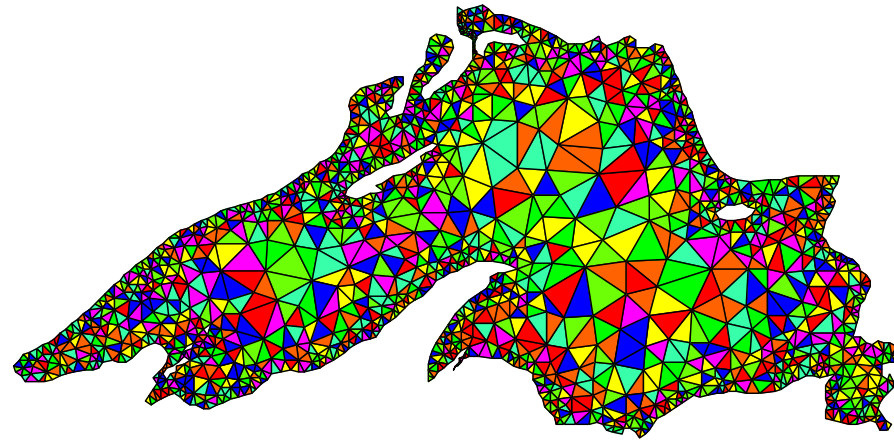
Mapping techniques can be **static** or **dynamic**.

- **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.

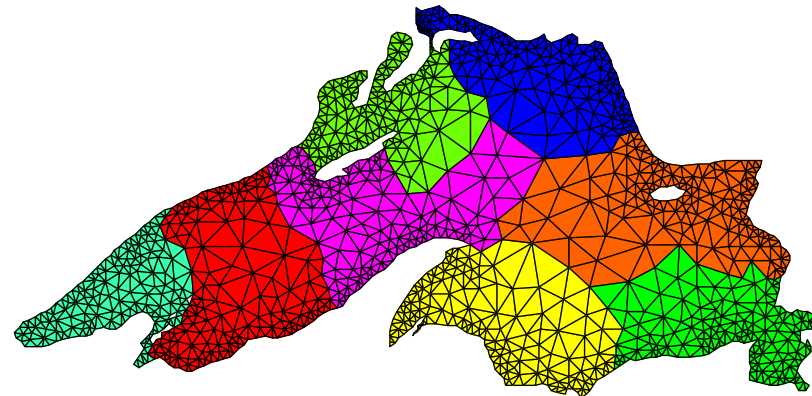
# Mappings Based on Task Partitioning

- 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.

# Partitioning the Graph of Lake Superior



Random Partitioning



Partitioning for minimum edge-cut.



# 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.
- 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 sending 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**.

# Minimizing Interaction Overheads

- **Maximize data locality:** Where possible, reuse intermediate data. Restructure computation so that data can be reused in smaller time windows.
- **Minimize volume of data exchange:** There is a cost associated with each word that is communicated. For this reason, we must minimize the volume of data communicated.
- **Minimize frequency of interactions:** There is a startup cost associated with each interaction. Therefore, try to merge multiple interactions to one, where possible.
- **Minimize contention and hot-spots:** Use decentralized techniques, replicate data where necessary.
- **Overlapping computations with interactions:** Use non-blocking communications, multithreading, and prefetching to hide latencies.
- **Replicating data or computations.**
- **Using group communications** instead of point-to-point primitives.
- **Overlap interactions** with other interactions.

# Parallel Algorithm Models

- **Data Parallel Model:** Tasks are statically (or semi-statically) mapped to processes and each task performs similar operations on different data.
- **Task Graph Model:** Starting from a task dependency graph, the interrelationships among the tasks are utilized to promote locality or to reduce interaction costs.
- **Master-Slave Model:** One or more processes generate work and allocate it to worker processes. This allocation may be static or dynamic.
- **Pipeline / Producer-Consumer Model:** A stream of data is passed through a succession of processes, each of which perform some task on it.
- **Hybrid Models:** A hybrid model may be composed either of multiple models applied hierarchically or multiple models applied sequentially to different phases of a parallel algorithm.