# CS3400 - Principles of Software Engineering
## Software Engineering for Multicore Systems
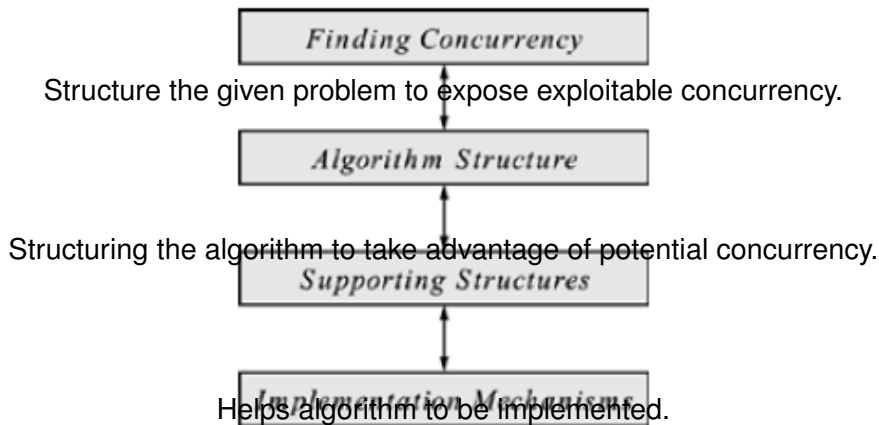
**V. Krishna Nandivada**

IIT Madras

# A pattern language

- **Pattern**: "a careful description of a perennial solution to a recurring problem within a . . . context."
- **Origin** Christopher Alexander, 1977 in the context of design and construction of building and town.
- Patterns in software engineering: Beck and Cunningham (1987), Gamma, Helm, Johnson, Vlissides (1995).
- **Pattern Language**: a structured method of describing good design practices within a field of expertise.
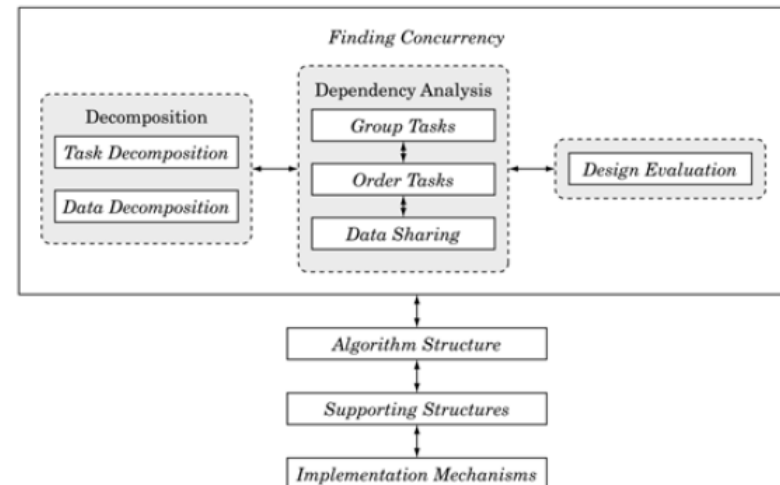
# A pattern language for parallel programs



Structure the given problem to expose exploitable concurrency.

Structuring the algorithm to take advantage of potential concurrency.

Helps algorithm to be implemented.

How the high level specifications are mapped.

**Goal**: Identify patterns in each stage.

# Finding concurrency in a given problem - deep dive

## Decomposition Patterns

- Task decomposition: A program to a sequence of "tasks".
  - Some of the tasks can run in parallel.
  - Independent the tasks the better.

- Data decomposition: Focus on the data used by the program. Decompose the program into tasks based on distinct chunks of data.
  - Efficiency depends on the independence of the chunks.

- Task decomposition may lead to data decomposition and vice versa.

Q: Are they really independent?

## Task decomposition: An approach

- Identify "resource" intensive parts of the problem.
- Identify different tasks that make up the problem. Challenge: write the algorithms and run the tasks concurrently.
- Sometimes the problem will naturally break into a collection of (nearly) independent tasks. Sometimes, not!
- Q: Are there enough tasks to keep the map all the H/W cores?
- Q: Does each task have enough work to keep the individual cores busy?
- Q: Are the number of tasks dependent or independent of the number of H/W core?
- Q: Are these tasks relatively independent?
- Instances of tasks: Independent modules, loop iterations.
- Relation between tasks and ease of programming, debugging and maintenance.

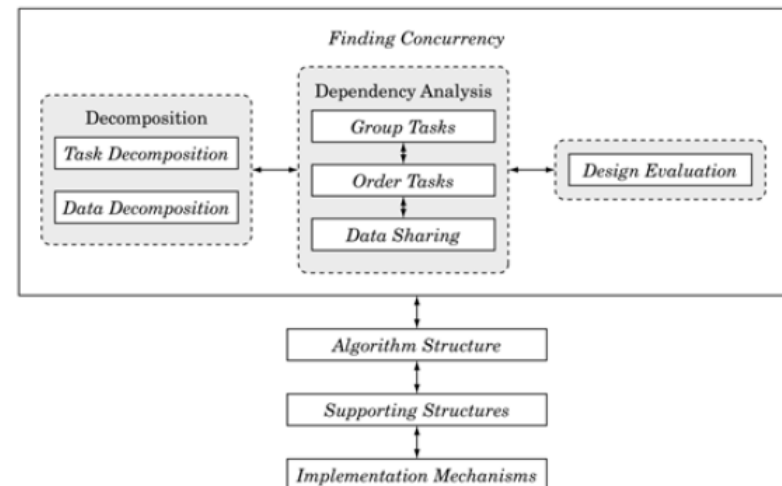## Task decomposition: Matrix multiplication example

$$C = A \times B$$

$$C_{i,j} = \sum_{k=0}^{N-1} A_{i,k} \times B_{k,j}$$

- "Resource" intensive parts?
- Tasks in the problem?
- Are tasks independent? Enough tasks for all the cores? Enough work for each task? Size of tasks and number of cores?
- Each element $C_{i,j}$ is computed in a different task - row major.
- Each element $C_{i,j}$ is computed in a different task - column major.
- Each element $C_{i,j}$ is computed in a different task - diagonals.
- How to reason about Performance? Cache effect?

## Finding concurrency in a given problem

## Data decomposition: Design

- Besides identifying the "resource" intensive parts, identify the key data structures required to solve the problem, and how is the data used during the solution.
- Q: Is the decomposition suitable to a specific system or many systems?
- Q: Does it scale with the size of parallel computer?
- Are similar operations applied to different parts of data, independently?
- Are there different chunks of data that can be distributed?
- Relation between decomposition and ease of programming, debugging and maintenance.
- Examples:
  - Array based computations: concurrency defined in terms of updates of different segments of the array/matrix.
  - Recursive data structures: concurrency by decomposing the parallel updates of a large tree/graph/linked list.

## Data decomposition: Matrix multiplication example

$$C = A \times B$$

$$C_{i,j} = \sum_{k=0}^{N-1} A_{i,k} \times B_{k,j}$$

- "Resource" intensive parts?
- Data chunks in the problem?
- Does it scale with the size of parallel computers?
- Operations (Reads/Writes) applied on independent parts of data?
- Data chunks big enough to deem the thread activity beneficial?
- How to decompose?
- Each row/column of $C_{i,j}$ is computed in a different task.
- Each column of $C_{i,j}$ is computed in a different task.
- Performance? Cache effect?
- Note: Data decomposition also leads to task decomposition as well.

## Matrix multiplication: Data decomposition.

$$C = \begin{pmatrix} A_{1,1} & A_{1,2} \\ A_{2,1} & A_{2,2} \end{pmatrix} \times \begin{pmatrix} B_{1,1} & B_{1,2} \\ B_{2,1} & B_{2,2} \end{pmatrix}$$
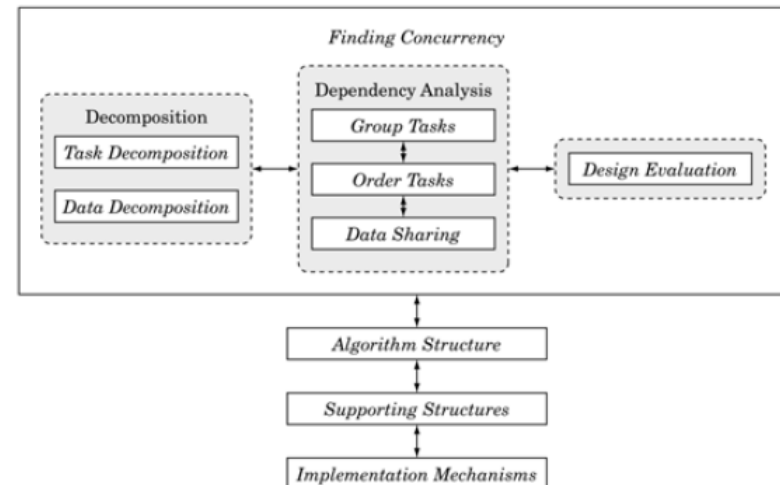
$$= \begin{pmatrix} A_{1,1} \times B_{1,1} + A_{1,2} \times B_{2,1} & A_{1,1} \times B_{1,2} + A_{1,2} \times B_{2,2} \\ A_{2,1} \times B_{1,1} + A_{2,2} \times B_{2,1} & A_{2,1} \times B_{1,2} + A_{2,2} \times B_{2,2} \end{pmatrix}$$

**Advantages**

- Can fit in the blocks into cache.
- Can scale as per the hardware.
- Overlap of communication and computation.

## Finding concurrency in a given problem

## Dependence analysis for managing parallelism: Grouping

- **Background**: Tasks and Data decomposition has been done.
- All the identified tasks may not run in parallel.
- **Q**: How should related tasks be grouped to help manage the dependencies?
- Dependent, related tasks should be (uniquely?) grouped together.
    - Temporal dependency: If task A depends on the result of task *B*, then *A* must wait for the results from *B*. Q: Does *A* have to wait for *B* to terminate?
    - Concurrent dependency: Tasks are expected to run in parallel, and one depends on the updates of the other.
    - Independent tasks: Can run in parallel or in sequence. Is it always better to run them in parallel?
- Advantage of grouping.
    - Grouping enforces partial orders between tasks.
    - Application developer thinks of groups, instead of individual tasks.
- Example: Computing of individual rows.

## Dependence analysis for managing parallelism: Ordering

- **Background**: Tasks and Data decomposition has been done. Dependent tasks have been grouped together.
- Ordering of the tasks and groups not trivial.
- **Q**: How should the groups be ordered to satisfy the constraints among the groups and in turn tasks?
- Dependent groups+tasks should be ordered to preserve the original semantics.
    - Should not be overly restrictive.
    - Ordering is imposed by: Data + Control dependencies.
    - Ordering can also be imposed by external factors: network, i/o and so on.
    - Ordering of independent tasks?
- Importance of grouping.
    - Ensures the program semantics.
    - A key step in program design.

## Dependence analysis for managing parallelism: data sharing

**Background**: Tasks and Data decomposition has been done. Dependent tasks have been grouped together. The ordering between the groups and tasks have been identified.

- Groups and tasks have some level of dependency among each other.
- **Q**: How is data shared among the tasks?
- Identify the data updated/needed by individual tasks - task local data.
- Some data may be updated by multiple tasks - global data.
- Some data may be updated by one data used by multiple tasks - remote data

## Issues in data sharing

- Identify the data being shared - directly follows from the decomposition.
- If sharing is done incorrectly - a task may get invalid data due to race condition.
- A naive way to guarantee correct shared data: synchronize every read with barriers.
- Synchronization of data across different tasks - may require communication. Options:
    - Overlap of communication and computation.
    - Privatization.
    - keep local copies of shared data.

## One special case of sharing

- Accumulation/Reduction: Data being used to accumulate a result; sum, minimum, maximum, variance etc.
  - Each core has a separate copy of data,
  - accumulation happens in these local copies.
  - sub-results are further used to compute the final result.
- Example: Sum elements in an array A[1024]
  - Decompose the array into 32 chunk.
  - Accumulate each chunk separately.
  - Accumulate the sub results into the global "sum".

## Managing parallelism - design evaluation

**Background**: Tasks and Data decomposition has been done. Dependent tasks have been grouped together. The ordering between the groups and tasks have been identified. A scheme for data sharing has also been identified.

- Of the multiple choices present at different points, we have chosen one.
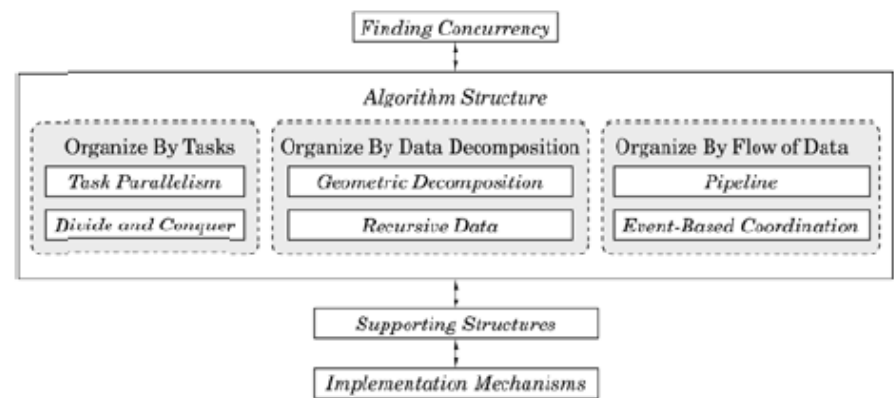- **Q**: Is the chosen path a "good" one?
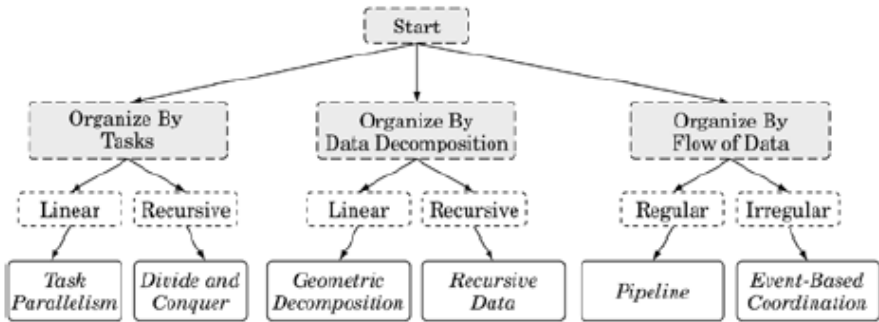
## Design evaluation factors

- Suitability to the target platform (at a high level)
  - Number of cores / HW threads - too few/many tasks?
  - Homogeneous/Heterogeneous multi-cores? And work distribution.
  - Data distribution among the cores - equal/unequal?
  - Cost of communication - fine/coarse grained data sharing.
  - Amount of sharing - shared memory or distributed memory.
- Metrics: simplicity (qualitative) , Efficiency , Flexibility
- Flexibility
  - Flexible/Parametric over the number of tasks?
  - Flexible/Parametric over the number and size of data chunks?
  - Does it handle boundary cases?
- Efficiency.
  - Even load balancing?
  - Minimum overhead? - task creation, synchronization, communication.

## Algorithm Structure - deep dive

# Algorithm Structure design

---

# Task Parallelism

**Q**: A problem is best decomposed into a collection of tasks that can execute concurrently. How to exploit the concurrency efficiently?

- Problem can be decomposed into a collection of concurrent tasks.
- Tasks can be completely independent or can have dependencies.
- Tasks can be known from the beginning (producer/consumer), tasks are created dynamically.
- Solution may or not require all the tasks to finish.

Challenges:

- Assign tasks to cores - to result in a simple, flexible and efficient execution.
- Address the dependencies correctly.

---

# Factors in efficient Task parallel algorithm design

- Tasks:
  1. Enough Tasks to keep the cores busy.
  2. Advantage of creating the tasks should offset the overhead of creating and managing them.
- Dependencies
  1. Ordering constraints.
  2. Dependencies from shared data: synchronization, private data.
  3. Schedule: creation and scheduling.
- Schedule
  1. How are the tasks assigned to cores.
  2. How are the tasks scheduled.

---

# Example: Task parallel algorithm

| Machine | Job1 | Job2 | Job3 | Job4 |
|---------|------|------|------|------|
| M1      | 4    | 4    | 3    | 5    |
| M2      | 2    | 3    | 4    | 4    |



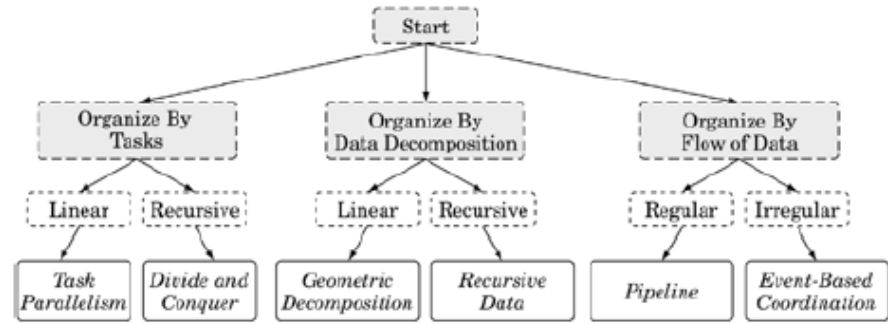- Say Job1 to M1, Job2 to M2, Job3 to M1, Job4 to M2 = 7.

# Solution to Branch and Bound ILP

- Maintain a list of tasks.
- Remove a solution from the list.
- **Examine the solution. Either discard it or declare it a solution, or add a sub-problem to task list.**
- The tasks depend depend on each other through the task-list.
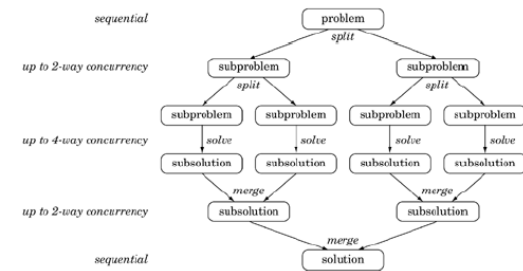
# Algorithm Structure design

# Divide and conquer

Q: Tasks are created recursively to solve a problem in a divide conquer strategy. How to exploit the concurrency?

- Divide and Conquer: Problem is solved by splitting it into a number of smaller subproblems. Examples?
- Each subproblems can be solved "fairly" independently. Directly or further divide and conquer.
- Solutions of the smaller problems is merged to compute the final solution.
- Each divide doubles the concurrency.
- Each merge halves the concurrency.

# Divide and Conquer pattern: features



- The amount of exploitable concurrency varies.
- At the beginning and end very little exploitable concurrency.
- Note: "split" and "merge" are serial parts.
- Amdahl's law - speed up constrained by the serial part. Impact?
- Too many parallel threads?
- What if cores are distributed? - data movement?
- Tasks are created dynamically - load balancing?
- What if the sub-problems are not equal-sized?
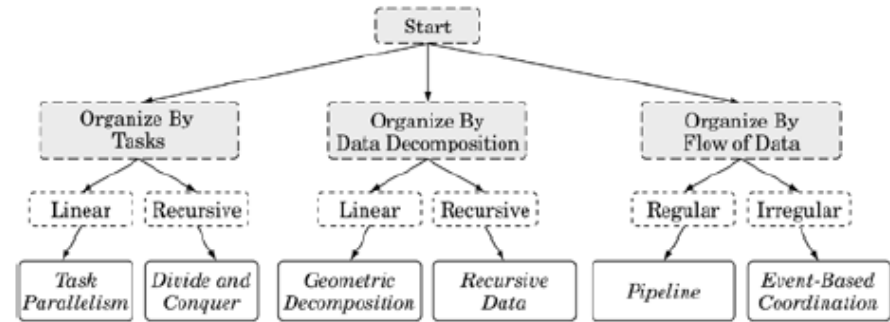
## Divide and conquer - example Mergesort

```
int[] mergesort(int[]A,int L,int H){
    if (H - L <= 1) return;
    if (H-L <= T) {quickSort(A, L, H); return;}
    int m = (L+H)/2;
    A1 = mergesort(A, L, m);
    A2 = mergesort(A, m+1, H);
    return merge(A1, A2);
    // returns a merged sorted array.
}
```

- split cost?
- merge cost?
- Value of threshold $T$?

## Algorithm Structure design

## Geometric decomposition

Q: How can an algorithm be organized around a data structure that has been decomposed into concurrently updatable "chunks"?

- Similar to decomposing a geometric region into subregions.
- Linear Data structures (such as arrays) - can be often decomposed into contiguous sub-structures.
- These individual tasks are processed in different concurrent tasks.
- Note: Sometimes all the required data for a task is present "locally" (embarrassingly parallel - Task parallelism pattern). And sometimes share data with "neighboring" chunks.
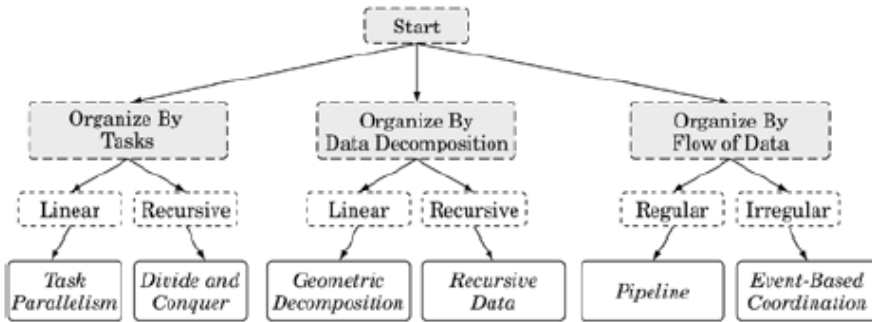
Challenges
- Ensure that each task has access to all data it needs.
- Mapping of chunks to cores giving good performance. Q: Why is it a challenge?
- Granularity of decomposition (coarse or fine-grain) - effect on efficiency? Parametric? Tweaked at compile time or runtime?
- Shape of the chunk: Regular/irregular?

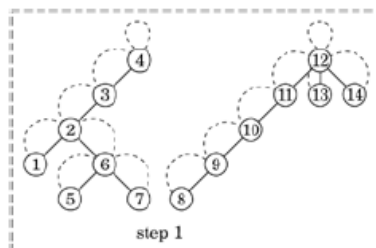## Geometric decomposition: Matrix multiplication

$$
\begin{aligned}
C &= A \times B \\
&= \left( \begin{array}{cc} A_{1,1} & A_{1,2} \\ A_{2,1} & A_{2,2} \end{array} \right) \times \left( \begin{array}{cc} B_{1,1} & B_{1,2} \\ B_{2,1} & B_{2,2} \end{array} \right) \\
&= \left( \begin{array}{cc} A_{1,1} \times B_{1,1} + A_{1,2} \times B_{2,1} & A_{1,1} \times B_{1,2} + A_{1,2} \times B_{2,2} \\ A_{2,1} \times B_{1,1} + A_{2,2} \times B_{2,1} & A_{2,1} \times B_{1,2} + A_{2,2} \times B_{2,2} \end{array} \right)
\end{aligned}
$$

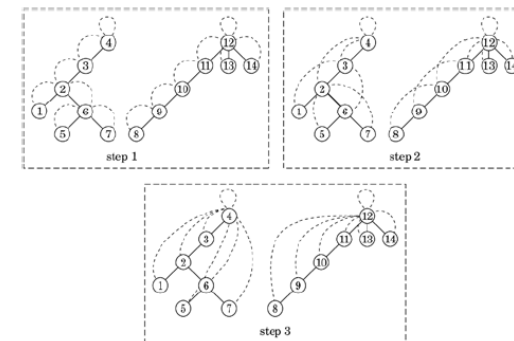## Algorithm Structure design

## Recursive Data Pattern

Q: How can recursive data structures be partitioned so as that operations on them are performed in parallel?

- Linked list, tree, graphs . . .
- Inherently operations on recursive data structures are serial - as one has to sequentially move through the data structure.
- For example linked list traversal or traversing a binary tree.
- Sometimes it is possible to reshape operations to derive and exploit concurrency.

## Recursive Data Pattern - example Find roots



- Given a forest of rooted trees: compute the root of each node.
- Serial version: Do a depth-first or breadth first traversal from root to the leaf nodes.
- For each visited node - set the root. Total running time?

Q: Is there concurrency?

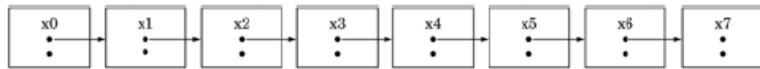## Recursive Data structures: Parallel find roots



- Transformed the original serial computation to one where we compute partial result and repeatedly combine partial results. Total Cost = ?
- Total cost = $O(N \log N)$
- However, if we exploit the parallelism - running time will come down to $O(\log N)$.
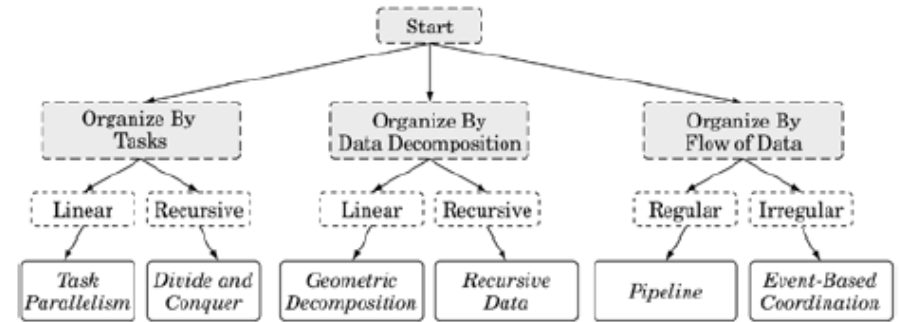
## Parallelizing recursive data structures

- Recasting the problem increases the cost. Find a way to get it back.
- Effective exploitation of the derived concurrency depends on factors such as - amount of work available for each task, amount of serial code ...
- Restructuring may make the solution complex.
- Requirement of synchronization - Why?
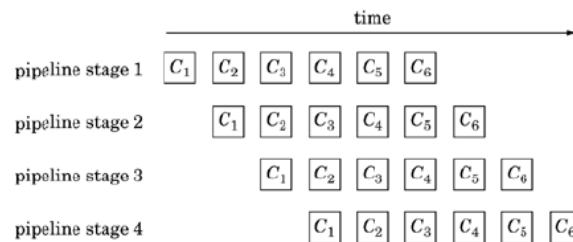- Another example: Find partial sums in a linked list.

## Algorithm Structure design

## Pipieline pattern

Q: The computation may involve performing similar sets of operations on many sets of data. Is there concurrency? How to exploit it?
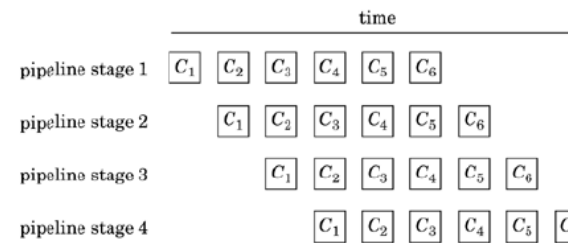
- Factory assembly line, Network Packet processing, Instruction processing in CPUs etc.



- There are ordering constraints on each operation on any one set of data: Operation $C_2$ can be undertaken only after $C_1$.
- Key requirement: Number of operations $> 1$.

## Pipeline pattern features



- Once the pipeline is full maximum parallelism is observed.
- Number of stages should be small compared to the number of items processed.
- Efficiency improves if time taken in each stage is roughly the same. Else?
- Amount of concurrency depends on the number of stages.
- Too many stages, disadvantage?
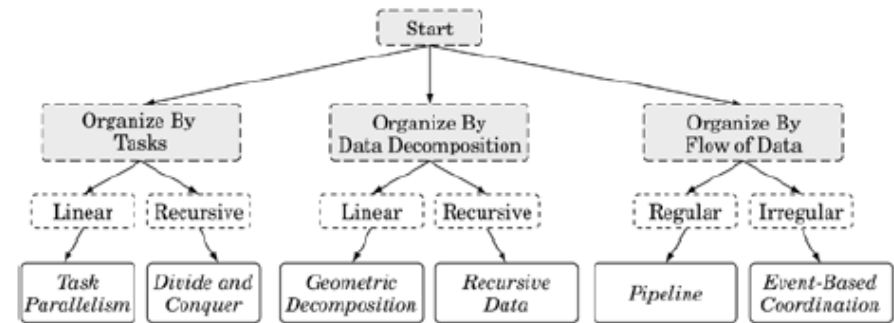- Communication across stages?

## Pipieline pattern. Issues

- Error handling.
  - Create a separate task for error handling - which will run exception routines.
- Processor allocation, load balancing
- Throughput and Latency.

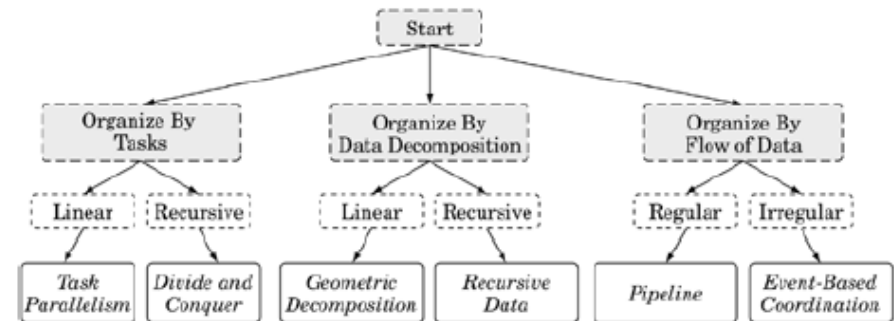## Overall big picture

## Event based coordination

**Challenges**

- Identifying the tasks.
- Identifying the events flow.
- Enforcing the events ordering.
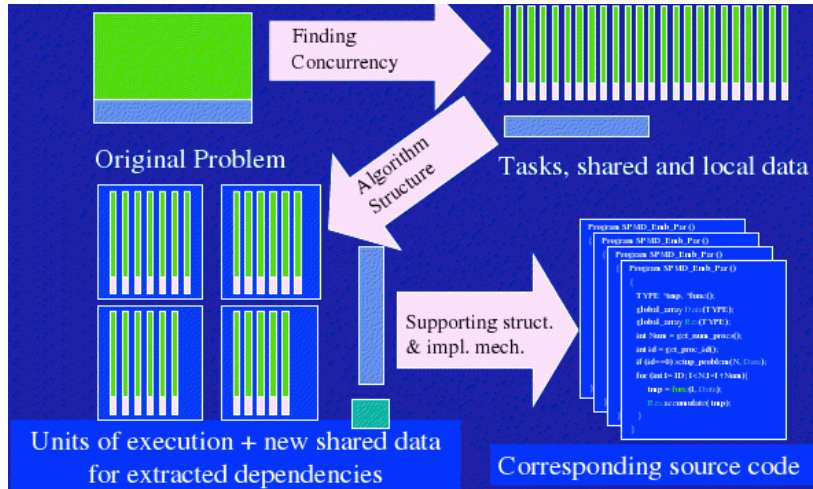- Avoiding deadlock.
- Efficient communication of events.

Left for self reading.

## Algorithm Structure design

# Overall big picture



- Original Problem
- Finding Concurrency
- Tasks, shared and local data
- Algorithm Structure
- Supporting struct. & impl. mech.
- Units of execution + new shared data for extracted dependencies
- Corresponding source code

# Sources

- Patterns for Parallel Programming: Sandors, Massingills.
- multicoreinfo.com
- Wikipedia
- fixstars.com
- Jernej Barbic slides.
- Loop Chunking in the presence of synchronization.
- Java Memory Model JSR-133: "Java Memory Model and Thread Specification Revision"