# Scalable Algorithmic Techniques
## *Decompositions & Mapping*

Alexandre David

1.2.05

adavid@cs.aau.dk

# Introduction

- Focus on data parallelism, scale with size.
  - Task parallelism limited.
- Notion of scalability is fuzzy in the book.
  - More precision later.
  - Idea: You lose efficiency with # of processors, gain efficiency with the size of the problem. Scalability measures the ratio.
    - You can experiment with assignment 2 to see that.

# In Practice...

- Typical tasks:
    - Identify concurrent works.
    - Map them to processors.
    - Distribute inputs, outputs, and other data.
    - Manage shared resources.
    - Synchronize the processors.

# Basic Principles

- Large blocks of independent computations.
  - Rare, seti@home.
    - Better when computations >> size of data.
  - Matrix multiplication too.
- Good performance recipe:
  - minimize interaction (= communication)
  - maximize locality (= blocks of computation)

# Minimizing Interaction Overheads

- **Maximize data locality.**
    - Minimize volume of data-exchange.
    - Minimize frequency of interactions.
- **Minimize contention and hot spots.**
    - Share a link, same memory block, etc...
    - Re-design original algorithm to change the interaction pattern.
    - Use task interaction graph to help.

# Minimizing Interaction Overheads

- **Overlapping computations with interactions – to reduce idling.**
  - Initiate interactions in advance.
  - Non-blocking communications.
  - Multi-threading.
- **Replicating data or computation.**
- **Group communication instead of point to point.**
- **Overlapping interactions.**

# Decomposing Problems

- Decomposition into *concurrent* tasks.
  - No unique solution.
  - Different sizes.
  - Decomposition illustrated as a directed graph:
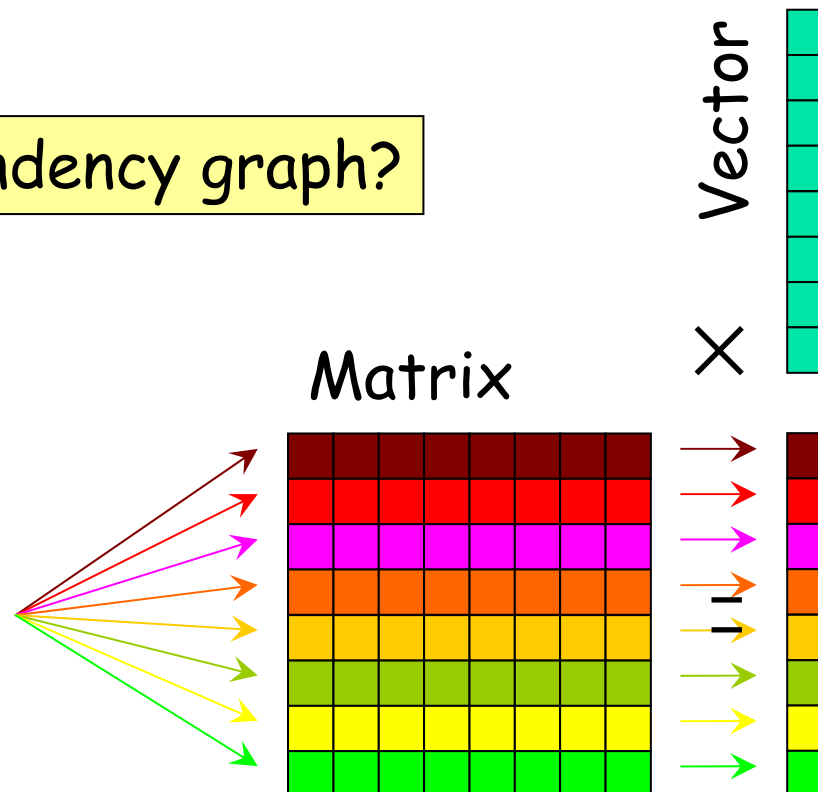    - Nodes = tasks.
    - Edges = dependency.

⚠ **Task dependency graph**

# Example: Matrix * Vector
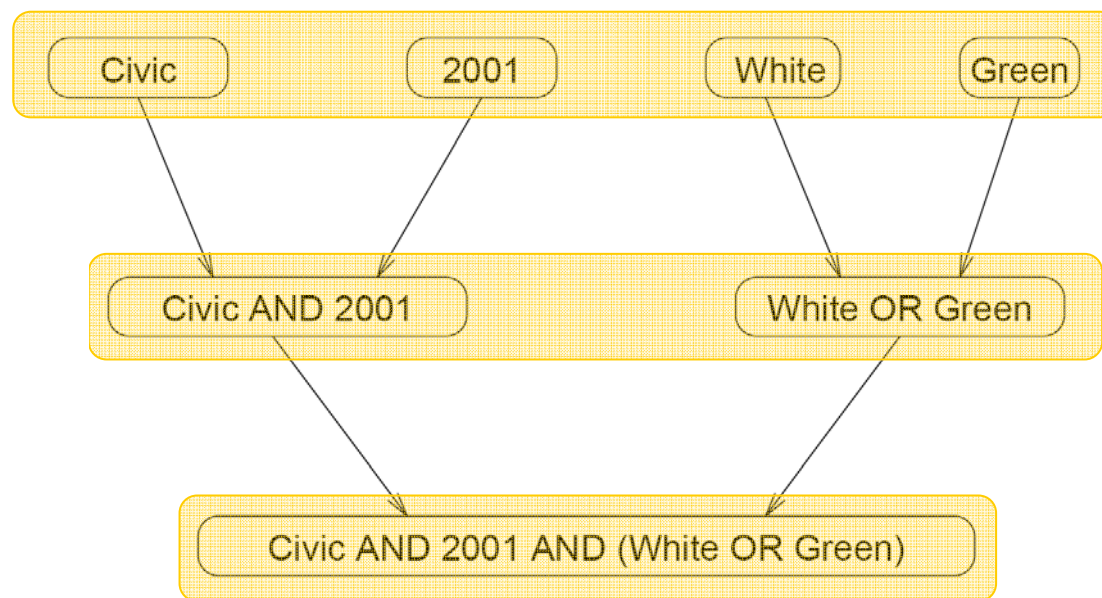


Task dependency graph?

N tasks, 1 task/row:

Matrix $\times$ Vector

# Example: database query processing

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

| 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 |

**Table 3.1**    A database storing information about used vehicles.

# A solution

Measure of concurrency?
Nb. of processors?
Optimal?



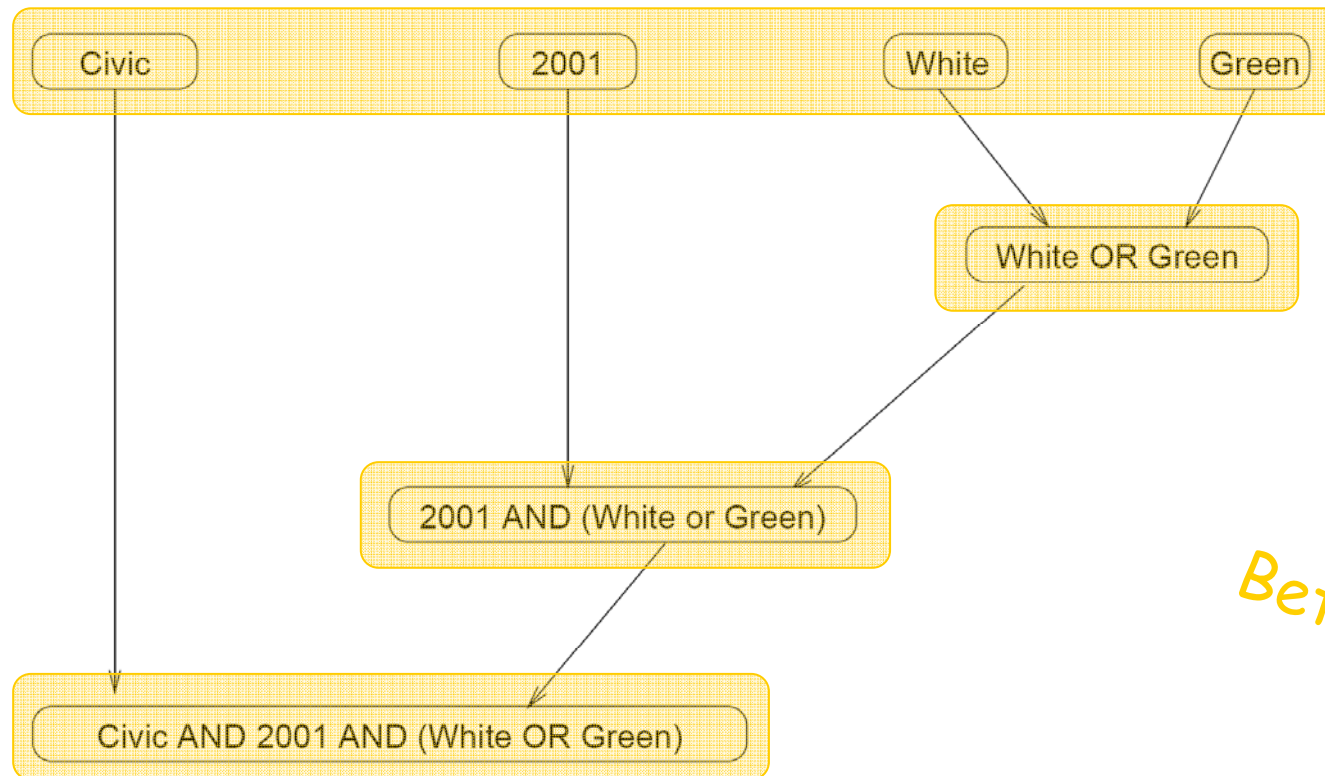**Figure 3.2** The different tables and their dependencies in a query processing operation.

# Another Solution

| ID# | Model |
|------|-------|
| 4523 | Civic |
| 6734 | Civic |
| 4395 | Civic |
| 7352 | Civic |

| ID# | Year |
|------|------|
| 7623 | 2001 |
| 6734 | 2001 |
| 5342 | 2001 |
| 3845 | 2001 |
| 4395 | 2001 |

| ID# | Color |
|------|-------|
| 3476 | White |
| 6734 | White |

| ID# | Color |
|------|-------|
| 7623 | Green |
| 9834 | Green |
| 5342 | Green |
| 8354 | Green |

Civic    2001    White    Green

White OR Green

2001 AND (White or Green)

Civic AND 2001 AND (White OR Green)

*Better/worse?*

11

# Granularity

- Number and size of tasks.
  - Fine-grained: many small tasks.
  - Coarse-grained: few large tasks.

- Related: *degree of concurrency*.
  (Nb. of tasks executable in parallel).
  - Maximal degree of concurrency.
  - Average degree of concurrency.

# Coarser Matrix * Vector



N tasks, 3 task/row:

Matrix

Vector

# Measures

- Average degree of concurrency if we take into account varying *amount of work?*

- Critical path = longest directed path between any start & finish nodes.

- Critical path length = sum of the weights of nodes along this path.

- Average degree of concurrency = total amount of work / critical path length.

# Database example

Critical path (3).

Critical path length = 27.

Av. deg. of concurrency = 63/27.

Critical path (4).

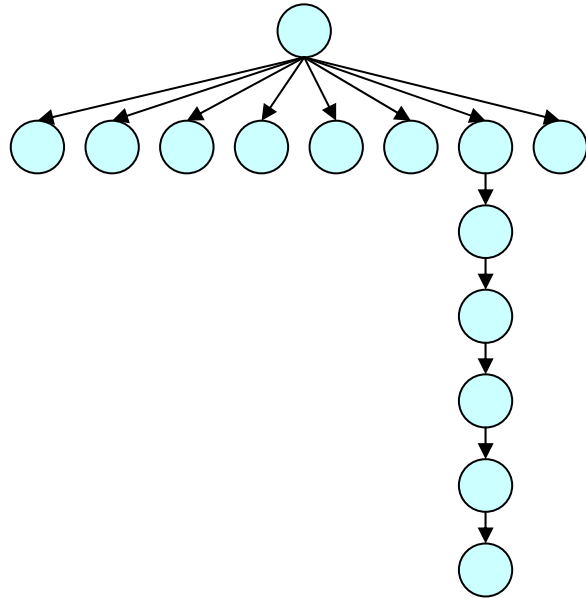Critical path length = 34.

Av. deg. of conc. = 64/34.



2.33

1.88

# Example



Number of tasks: 15.

- Maximum degree of concurrency: 8.
- Critical path length: 4.
- Maximum possible speedup: 15/4.
- Minimum number of processes to reach this speedup: 8.
- Maximum speedup if we limit the processes to 2,4, and 8: 15/8, 3, and 15/4.

# Example
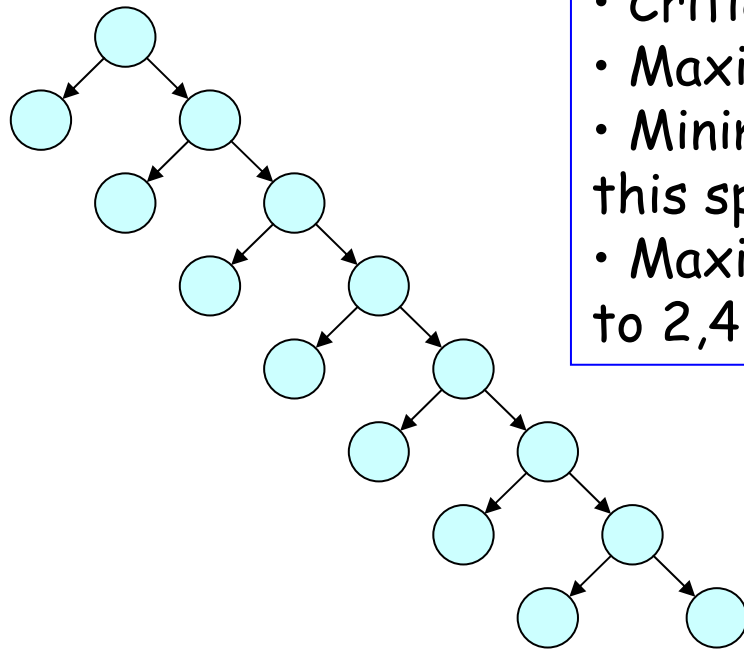


Number of tasks: 15.

- Maximum degree of concurrency: 8.
- Critical path length: 4.
- Maximum possible speedup: 15/4.
- Minimum number of processes to reach this speedup: 8.
- Maximum speedup if we limit the processes to 2,4, and 8: 15/8, 3, and 15/4.

# Example



- Maximum degree of concurrency: 8.
- Critical path length: 7.
- Maximum possible speedup: 14/7.
- Minimum number of processes to reach this speedup: 3.
- Maximum speedup if we limit the processes to 2,4, and 8: 14/8, 14/7, and 14/7.

Number of tasks: 14.

# Example



- Maximum degree of concurrency: 2.
- Critical path length: 8.
- Maximum possible speedup: 15/8.
- Minimum number of processes to reach this speedup: 2.
- Maximum speedup if we limit the processes to 2,4, and 8: 15/8.

Number of tasks: 15.

# Interaction Between Tasks

- Tasks often share data.

- Task interaction graph:
  - Nodes = tasks.
  - Edges = interaction.
  - Optional weights.

- Task dependency graph is a sub-graph of the task interaction graph.

# Characteristics of Task Interactions

- **One-way interactions.**
  - Only one task initiates and completes the communication **without** interrupting the other one.

- **Two-way interactions.**
  - Producer – consumer model.

# Processes and Mapping

- Tasks run on processors.

- Process: processing agent executing the tasks. *Not exactly like in your OS course.* Processes ~ threads here.

- Mapping = assignment of tasks to processes.

- API exposes processes and binding to processors not always controlled.

  - Scheduling of threads is not controlled.
  - What makes a good mapping?

# Mapping example



**Figure 3.7** Mappings of the task graphs of Figure 3.5 onto four processes.

# Processes vs. processors

- Processes = logical computing agent.
- Processor = hardware computational unit.
- In general 1-1 correspondence but this model gives better abstraction.
- Useful for hardware supporting multiple programming paradigms.

How do you decompose?

# Decomposition Techniques

- **Recursive decomposition.**
  - Divide-and-conquer.

- **Data decomposition.**
  - Large data structure.

- Exploratory decomposition.
  - Search algorithms.

- Speculative decomposition.
  - Dependent choices in computations.

# Recursive decomposition

- **Problem solvable by divide-and-conquer:**
  - **Decompose** into sub-problems.
    - Do it recursively.
  - **Combine** the sub-solutions.
    - Do it recursively.
- **Concurrency: The sub-problems are solved in parallel.**

# Schwartz's Algorithm
## *+ reduce by block*

- **Reduce with maximal concurrency:**
  - one thread per pair (n/2)
  - combine results in a tree structure
- **Schwartz:**
  - one thread per n/p block of numbers
  - local sums
  - combine results in a tree structure
  - *follows recipe*

**Recursive Decomposition**

# Combining Results

# Peril-L

full/empty variable

intermediate value – constant amount of extra space



```
int nodeval'[P];
…
forall(index in (0..P-1))
{
  int tally;
  stride=1;
  …
  while(stride<P)
  {
    if (index%(2*stride)==0)
    {
      tally += nodeval'[index+stride];
      stride *= 2;
    }
    else
    {
      nodeval'[index]=tally;
      break;
    }
  }
}
```

# Reduce & Scan Abstractions

- Reduce: combine values to a single one.
  - Almost always needed.

- Scan: prefix computation.
  - Logic that performs sequential operations and carries along intermediate results.

- Lesson: Try to use them as much as possible.
  - Abstract them as functions.
    - high-level, contain information
    - may customize implementation (e.g. BlueGene).

# Small typo p130
$A=\{0,2,4\} \Rightarrow A=\{0,2,6\}$

# Basic Structure

- Idea:
    - Assume block allocation,
    - use Schwartz's like algorithm,
    - local variable – tally – stores intermediate results.
- Primitives:
    - `init()` – init tally
    - `accum()` – local accumulation
    - `combine()` – combines tally results
    - **x**`-gen()` – final answer

# Example: + reduce

- `init():` tally=0
- `accum(tally,val):` tally+=value
- `combine(left,right):` left+right sent to parent
- `reduce-gen(root):` return

# Reduce Basic Structure

```
1    int nodeval'[P];                                  Global full/empty variables
2    int result;
3    forall(index in(0..P-1))
4    {
5      int myData[size]=localize(dataarray[]);         Local portion of global data values
6      int tally;
7      int stride=1;
8      tally=init ()                                    Initialize tally
9      for(i=0; i<size; i++)
10     {
11       tally=accum (tally, myData[i]);                Local accumulation
12     }
13     nodeval'[index]=tally;                           Send initially to parent
14     while(stride < P)                                Begin logic for tree
15     {
16       if(index%(2*stride)==0)
17       {                                              Combine values globally
18         nodeval'[index]=combine(nodeval'[index],
19                                 nodeval'[index+stride]);
20         stride=2*stride;
21       }
22       else
23       {
24         break;
25       }
26     }
27     if(index==0)
28     {
29       result=reduceGen (nodeval'[0]);                Generate reduced value
30     }
31   }
```

```
1   struct tally
2   {
3      float smallest1;                              Smallest element
4      float smallest2;                              Second smallest
5   };
6
7   tally init()                                     Initialize tally
8   {
9      tally t;
10     t.smallest1=MAX_FLOAT;
11     t.smallest2=MAX_FLOAT;
12     return t;
13  }
14
15  tally accum(tally t, float elem)                 Local accumulation
16  {
17     if(t.small1>elem)                             Is this a new smallest?
18     {
19        t.smallest2=t.smallest1;
20        t.smallest1=elem;
21     }
22     else
23     {
24        if(t.smallest2>elem)                       Is it a new second smallest?
25        {
26           t.smallest2=elem;
27        }
28        return t;
29     }
```

```
24      if(t.smallest2>elem)                    Is it a new second smallest?
25      {
26         t.smallest2=elem;
27      }
28      return t;
29    }
30  }
31
32  tally combine(tally left, tally right)       Combine into "left" by
33  {                                            accumulating right values
34    tally t;
35    t=accum(left, right.smallest1);
36    t=accum(t, right.smallest2);
37    return t;
38  }
39
40  float reduceGen(tally t)
41  {
42    return t.smallest2;
43  }
```

# General Scan

- **Difference with reduce:**
  - need to pass intermediate results too.
  - Propagate tally down the tree:
    value from a parent = tally from the left sub-tree of the parent.
  - root has no parent – fix that
- **Idea:**
  - up-sweep with reduce
  - down-sweep to propagate tallys

# Prefix Sum - sum

# Down-sweep

# Scan in Peril-L

```
 1   int nodeval'[P];                                Global full/empty memory
 2   int ltally[P];                                  Store left operand of combine
 3   forall(index in(0..P-1))
 4   {
 5     int myData[size]=localize(operandArray[]);    Local data values
 6     int tally;                                     Tally
 7     int ptally;                                    Tally from parent
 8     int stride=1;
 9     tally=init ();                                 Initialize
10     for(i=0; i<size; i++)
11     {
12       tally=accum (tally, myData[i]);             Accumulate
13     }
14     nodeval'[index]=tally;                         Send initially to parent
15     while(stride<P)                                Begin logic for tree
```

```
16    {
17      if(index%(2*stride)==0)
18      {                                              Combine
19        ltally[index+stride]=nodeval'[index];
20        nodeval'[index]=combine (ltally[index+stride],
21                                 nodeval'[index+stride]);
22        stride=2*stride;
23      }
24      else
25      {
26        break;
27      }
28    }
29    stride=P/2;
30    if(index==0)
31    {
32      ptally=nodeval'[0];                            Clear existing up sweep value
33      nodeval'[0]=init ();                           Set init() as parent input
34    }
35    while(stride>1)                                  Begin logic for tree descent
36    {
37      ptally=nodeval'[index];                        Grab parent value
38      nodeval'[index]=ptally;                        Send it down to left
39      nodeval'[index+stride]=                        Send parent + left child right
          combine (ptally, ltally[index+stride]);
40      stride=stride/2;                               Go down to next level
41    }
42    for(i=0; i<size; i++)
43    {
44      myResult[i]=scanGen (ptally, myData[i]);       Generate Scan
45    }
46  }
```

# Lesson

- Structure the algorithm with reduce & scan.
- Use efficient implementations of reduce & scan.

# Data Decomposition

- 2 steps:
  - Partition the data.
  - Induce partition into tasks.
- How to partition data?
- Partition output data:
  - Independent "sub-outputs".
- Partition input data:
  - Local computations, followed by combination.
- 1-D, 2-D, 3-D block decomposition.

# Static Allocation of Work to Processes

- # of threads fixed but unknown.
  - Allocate data to threads.
  - Owner compute rule.
- Block allocation – maximize locality
  - 1-D or 2-D depending on the communication pattern – minimize communication
    *surface area to volume in favour of 2-D*

# Owner-Compute Rule

- **Process assigned to some data**
  - is responsible for all computations associated with it.

- **Input data decomposition:**
  - All computations done on the (partitioned) input data are done by the process.

- **Output data decomposition:**
  - All computations for the (partitioned) output data are done by the process.

# 1-D & 2-D Block Allocations



(a)

(b)

# Overlap Regions

- Obtain data from neighbors.

- Compute locally.

- Avoid false sharing.

- Use local matrix
  - no special edge cases
  - uniform indices
  - batch communication cheaper

# Overlap Regions

# Cyclic & Block Cyclic

- Cyclic = round-Robin. Idea:
    - Partition an array into many *more blocks than available processes*.
    - Assign partitions (tasks) to processes in a round-robin manner.
    - $\rightarrow$ each process gets several non adjacent blocks.
- Useful when computations are not proportional to the data.
    - ex: assignment 2
    - otherwise poor load balance
- Good: load balance.
- Bad: more communication, break large blocks.

# Example: LU factorization

- Non singular square matrix A (invertible).
- A = L*U.
- Useful for solving linear equations.

$$A = L \times U$$

# LU factorization

In practice we work on A.



N steps

# Load Imbalance: LU-Decomposition



Matrix inversion – similar.

# 8x8 Array on 5 Processes

# Block Cyclic Distribution

# Julia Sets
## *Assignment: Mandelbrot*

# Irregular Allocations

# Randomized Distributions



(a)

| $P_0$ | $P_1$ | $P_2$ | $P_3$ | $P_0$ | $P_1$ | $P_2$ | $P_3$ |
|---|---|---|---|---|---|---|---|
| $P_4$ | $P_5$ | $P_6$ | $P_7$ | $P_4$ | $P_5$ | $P_6$ | $P_7$ |
| $P_8$ | $P_9$ | $P_{10}$ | $P_{11}$ | $P_8$ | $P_9$ | $P_{10}$ | $P_{11}$ |
| $P_{12}$ | $P_{13}$ | $P_{14}$ | $P_{15}$ | $P_{12}$ | $P_{13}$ | $P_{14}$ | $P_{15}$ |
| $P_0$ | $P_1$ | $P_2$ | $P_3$ | $P_0$ | $P_1$ | $P_2$ | $P_3$ |
| $P_4$ | $P_5$ | $P_6$ | $P_7$ | $P_4$ | $P_5$ | $P_6$ | $P_7$ |
| $P_8$ | $P_9$ | $P_{10}$ | $P_{11}$ | $P_8$ | $P_9$ | $P_{10}$ | $P_{11}$ |
| $P_{12}$ | $P_{13}$ | $P_{14}$ | $P_{15}$ | $P_{12}$ | $P_{13}$ | $P_{14}$ | $P_{15}$ |

(b)

Irregular distribution with regular mapping!
Not good.

# 1-D Randomized Distribution

$$V = [0, 1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11]$$

Permutation

$$random(V) = [8, 2, 6, 0, 3, 7, 11, 1, 9, 5, 4, 10]$$

$$mapping = 8 \ 2 \ 6 \quad 0 \ 3 \ 7 \quad 11 \ 1 \ 9 \quad 5 \ 4 \ 10$$

$$P_0 \qquad P_1 \qquad P_2 \qquad P_3$$

**Figure 3.32** A one-dimensional randomized block mapping of 12 blocks onto four process (i.e., $\alpha = 3$).
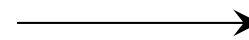
# 2-D Randomized Distribution



(a)                  (b)                  (c)

2-D block random distribution.

Block mapping.

# Irregular Allocations

- Same idea as overlap regions:
  - get data local – *inspector*
  - local computations – *executor*

# Dynamic Allocations

- Keys
  - dynamic load balancing
  - dynamic interactions
  - choose right granularity of tasks
- Work queues
  - e.g. producer/consumer
  - centralized schemes with master/slave
  - different queue orderings
  - multiple queues – issues with load balancing

# Graph Partitioning

- For sparse data structures and data dependent interaction patterns.
  - Numerical simulations. Discretize the problem and represent it as a mesh.
- Sparse matrix: assign equal number of nodes to processes & minimize interaction.
- Example: simulation of dispersion of a water contaminant in Lake Superior.

# Discretization



**Figure 3.34**   A mesh used to model Lake Superior.
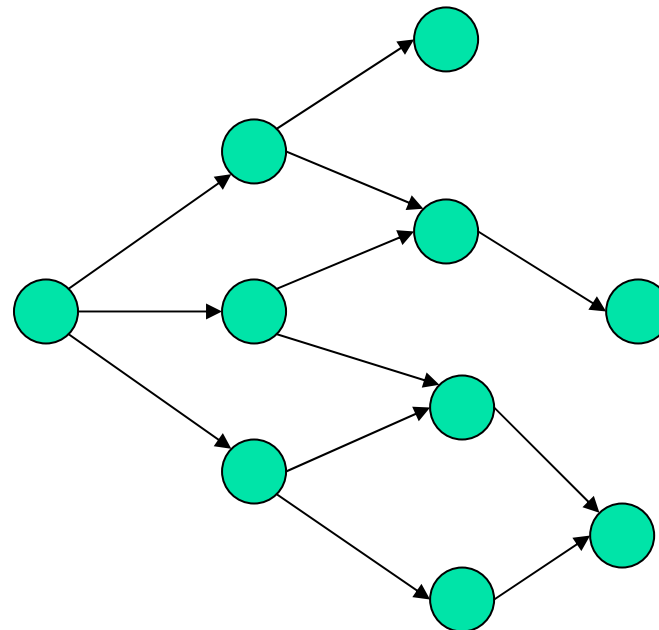
# Partitioning Lake Superior



Random partitioning.

Partitioning with minimum edge cut.

Finding an exact optimal partitioning is an NP-complete problem.

# Exploratory Decomposition - Trees
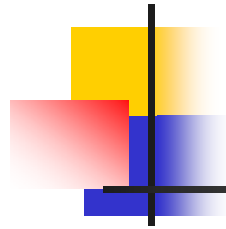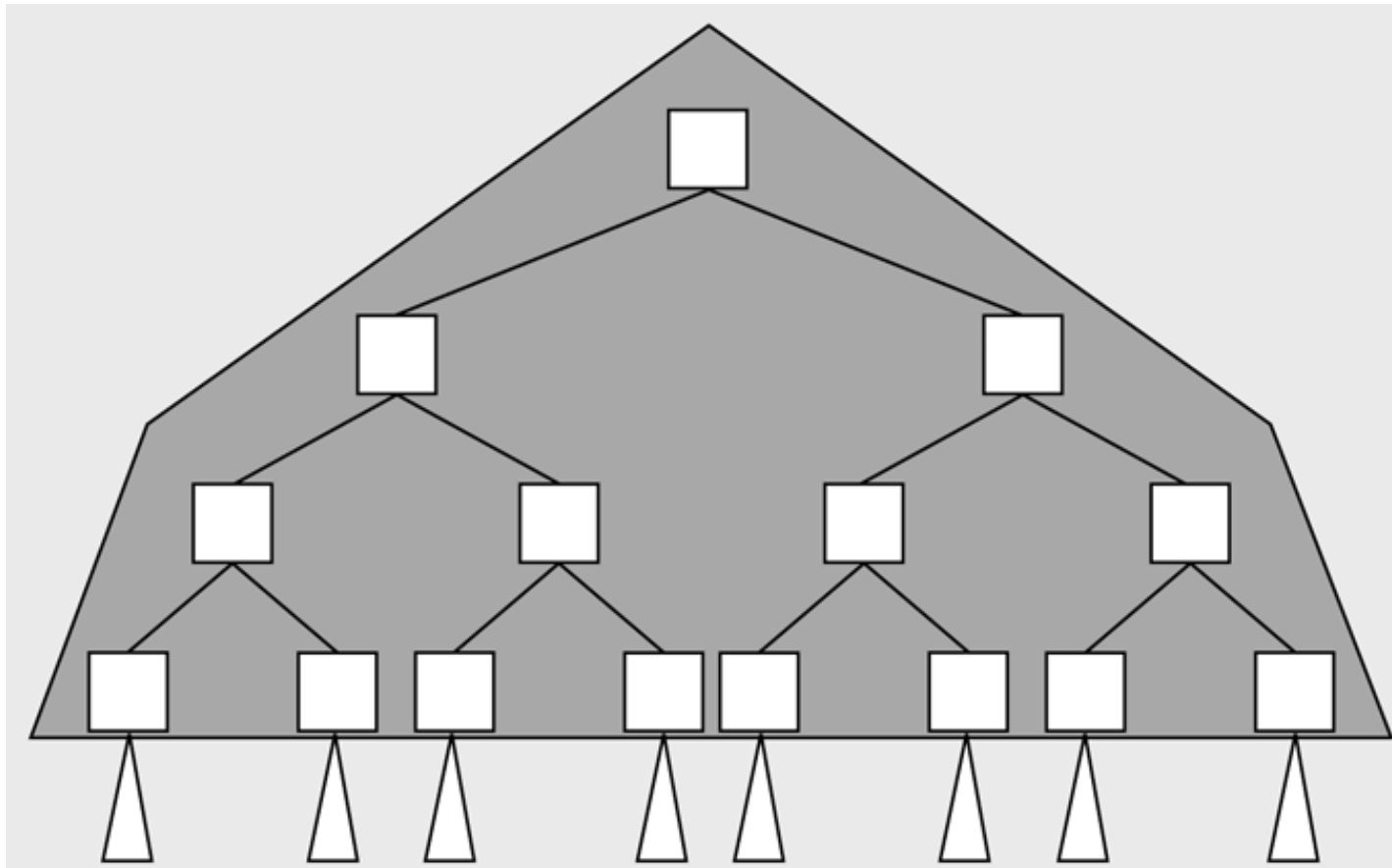
Exploration of states.

# Trees

- Useful data-structures.

- Usually constructed with pointers.

- Challenges for
  - communication
  - load balance on irregular trees

- If little communication among sub-trees:
  - Allocate sub-trees to processes, copy the "cap".
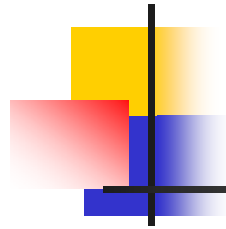  - All processes know the structure.
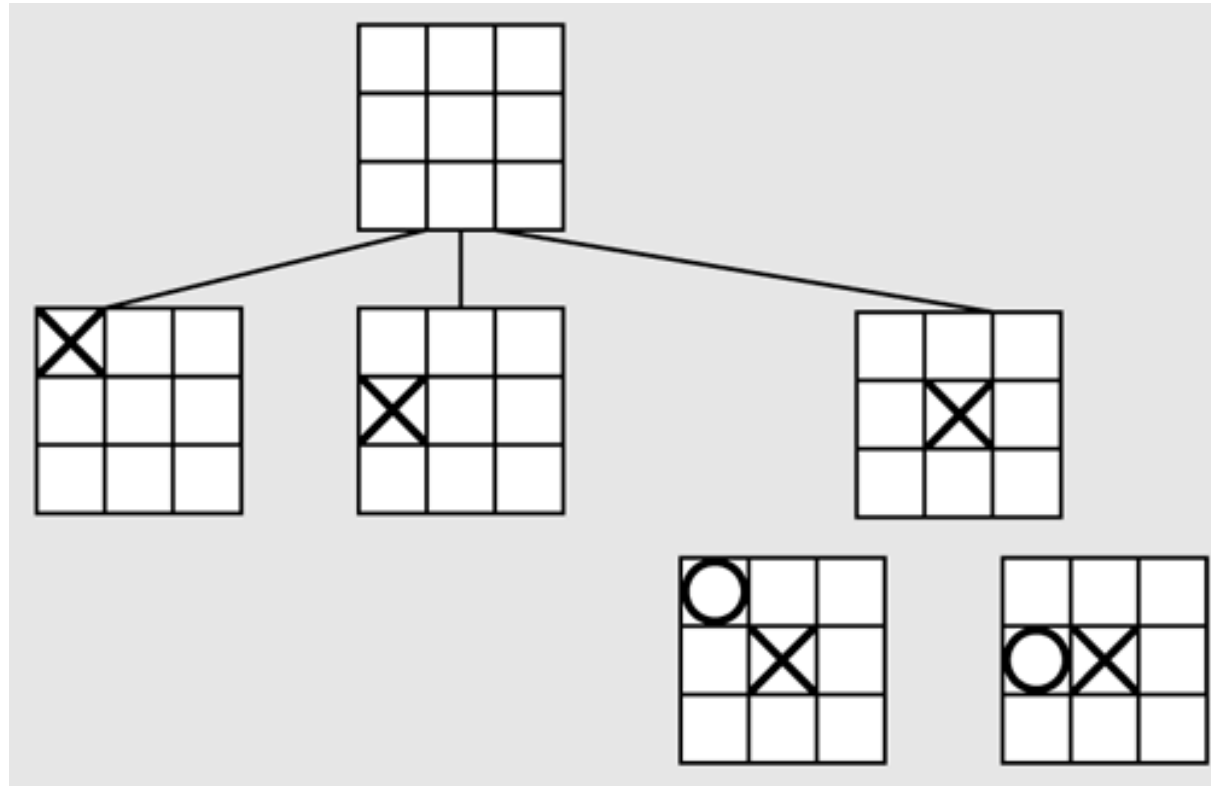
# Cap Copy & Subtree Allocation

# Dynamic Allocations

- Dynamic & unpredictable trees.
  - Search with different algorithms.
  - Work queues useful.
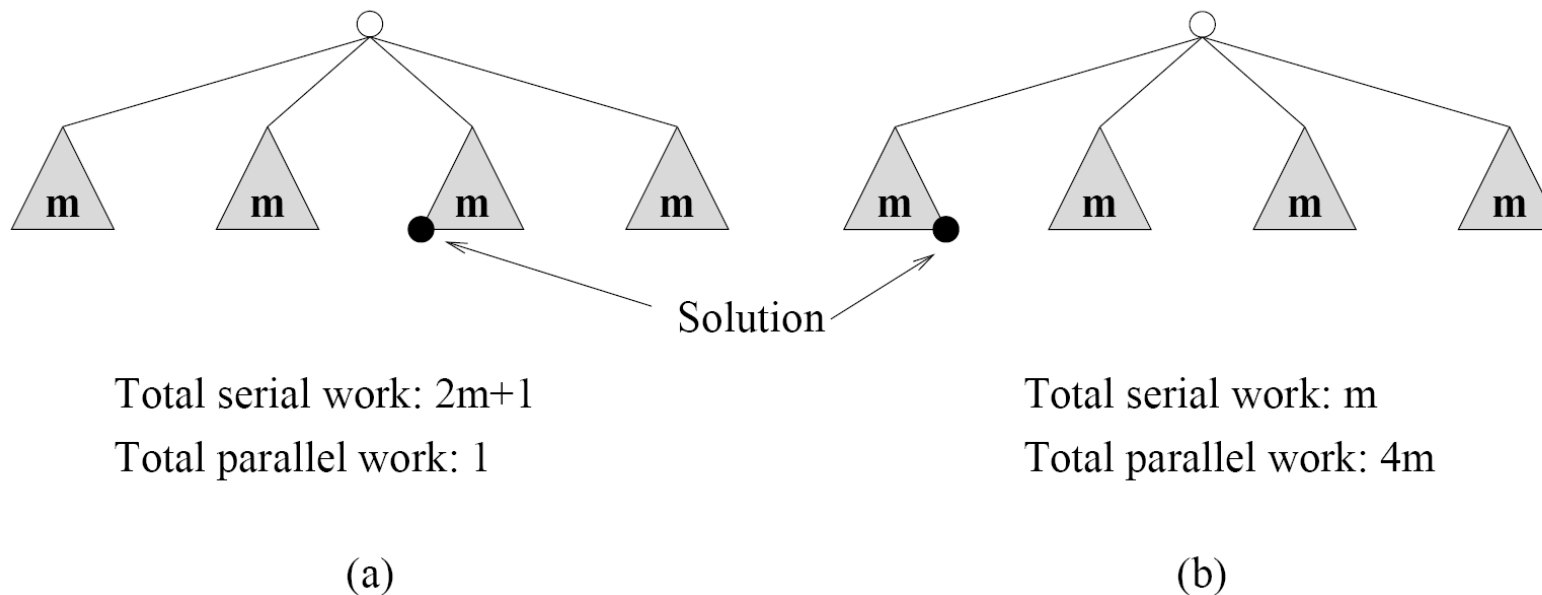  - Pruning involves communication.
  - Termination may be an issue!

# Application to Game Search

# Performance Anomalies
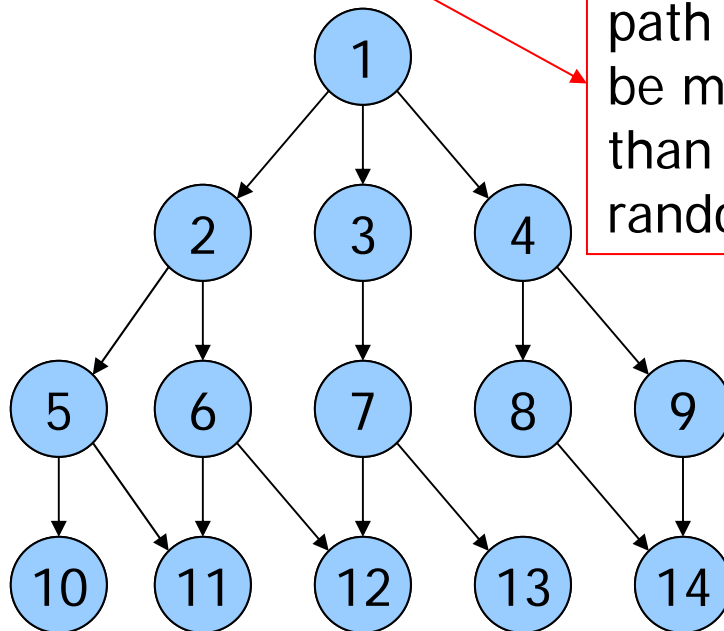
Work depends on the order of the search!



Total serial work: 2m+1
Total parallel work: 1

Total serial work: m
Total parallel work: 4m

(a)

(b)

**Figure 3.19** An illustration of anomalous speedups resulting from exploratory decomposition.

# Search Orderings Issues

# Speculative Decomposition

- Dependencies between tasks are not known a-priori.
  - How to identify independent tasks?
  - Conservative approach: identify tasks that are *guaranteed* to be independent.
  - Optimistic approach: schedule tasks even if we are not sure – may roll-back later.