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.







### 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? White) Civic 2001 Green White OR Green Civic AND 2001 Civic AND 2001 AND (White OR Green)

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

### **Another Solution**



# 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.





### 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.



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.







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.





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.





- 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.



- 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?





**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
  - Iocal sums
  - combine results in a tree structure
  - follows recipe

### **Recursive Decomposition**







## 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,
  - Iocal 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
- ocmbine(left,right): left+right sent to
  parent
- reduce-gen(root): return

## **Reduce Basic Structure**



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

```
struct tally
       1
       2
          {
                                                                    Smallest element
            float smallest1;
       3
                                                                    Second smallest
            float smallest2;
       4
       5
          };
       6
                                                                    Initialize tally
          tally init()
       7
       8
          {
2<sup>nd</sup> Min in Peril-L
      9
            tally t;
     10
           t.smallest1=MAX FLOAT;
     11
            t.smallest2=MAX FLOAT;
     12
            return t;
     13
          }
     14
                                                                    Local accumulation
     15
          tally accum(tally t, float elem)
     16
          {
                                                                    Is this a new smallest?
     17
            if(t.small1>elem)
     18
             {
     19
               t.smallest2=t.smallest1;
     20
               t.smallest1=elem;
     21
            }
     22
            else
     23
             {
                                                                    Is it a new second smallest?
     24
               if(t.smallest2>elem)
     25
               {
     26
                 t.smallest2=elem;
     27
               }
     28
               return t;
11+1429
```
```
24
             if(t.smallest2>elem)
                                                               Is it a new second smallest?
    25
             {
    26
               t.smallest2=elem;
    27
             }
    28
             return t;
    29
           }
    30
         }
    31
                                                               Combine into "left" by
    32
         tally combine(tally left, tally right)
Peril-I
                                                               accumulating right values
    33
         {
    34
           tally t;
    35
         t=accum(left, right.smallest1);
    36
           t=accum(t, right.smallest2);
Min in
    37
           return t;
    38
        }
    39
    40
         float reduceGen(tally t)
    41
         {
2nd
    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













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

|           | 16 | {  |                                |
|-----------|----|--|--------------------------------|
|           | 17 | if(index%(2*stride)==0)  |                                |
|           | 18 | {  | Combine                        |
|           | 19 | <pre>ltally[index+stride]=nodeval'[index];</pre>                   |                                |
|           | 20 | <u>nodeval'</u> [index]= <b>combine</b> ( <u>ltally</u> [index+str | ride],                         |
|           | 21 | <u>nodeval'</u> [index+s   | stride]);                      |
|           | 22 | <pre>stride=2*stride;</pre>  |                                |
|           | 23 | }  |                                |
|           | 24 | else   |                                |
|           | 25 | {  |                                |
|           | 26 | break;   |                                |
|           | 27 | }  |                                |
|           | 28 | }  |                                |
|           | 29 | <pre>stride=P/2;</pre>   |                                |
|           | 30 | if(index==0)   |                                |
|           | 31 | {  |                                |
|           | 32 | ptally= <u>nodeval'</u> [0];                                       | Clear existing up sweep value  |
|           | 33 | <pre>nodeval'[0]=init ();</pre>                                    | Set init() as parent input     |
|           | 34 | }  |                                |
|           | 35 | while(stride>1)  | Begin logic for tree descent   |
|           | 36 | {  |                                |
|           | 37 | ptally= <u>nodeval'</u> [index];                                   | Grab parent value              |
|           | 38 | <pre>nodeval'[index]=ptally;</pre>                                 | Send it down to left           |
|           | 39 | <pre>nodeval'[index+stride]=</pre>                                 | Send parent + left child right |
|           |    | <pre>combine (ptally, <u>ltally[index+stride]);</u></pre>          |                                |
|           | 40 | <pre>stride=stride/2;</pre>  | Go down to next level          |
|           | 41 | }  |                                |
|           | 42 | <pre>for(i=0; i<size; i++)<="" pre=""></size;></pre>               |                                |
|           | 43 | {  |                                |
|           | 44 | <pre>myResult[i]=scanGen (ptally, myData[i]);</pre>                | Generate Scan                  |
|           | 45 | }  |                                |
| 11+14-03- | 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



## **Overlap Regions**

- Obtain data from neighbors.
- Compute locally.
- Avoid false sharing.
- Use local matrix
  - no special edge cases
  - uniform indices
  - batch communication cheaper





## 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.





#### In practice we work on A.







#### Matrix inversion - similar.





# Block Cyclic Distribution



#### Julia Sets Assignment: Mandelbrot



# Irregular Allocations



## Randomized Distributions



| $P_0$    | $P_1$                  | $P_2$                  | <i>P</i> <sub>3</sub>  | $P_0$    | $P_1$                  | $P_2$    | <i>P</i> <sub>3</sub>  |
|----------|------------------------|------------------------|------------------------|----------|------------------------|----------|------------------------|
| $P_4$    | $P_5$                  | $P_6$                  | $P_7$                  | $P_4$    | $P_5$                  | $P_6$    | $P_7$                  |
| $P_8$    | $P_9$                  | $P_{10}$               | <i>P</i> <sub>11</sub> | $P_8$    | $P_9$                  | $P_{10}$ | <i>P</i> <sub>11</sub> |
| $P_{12}$ | <i>P</i> <sub>13</sub> | <i>P</i> <sub>14</sub> | $P_{15}$               | $P_{12}$ | <i>P</i> <sub>13</sub> | $P_{14}$ | <i>P</i> <sub>15</sub> |
| $P_0$    | $P_1$                  | $P_2$                  | $P_3$                  | $P_0$    | $P_1$                  | $P_2$    | $P_3$                  |
| $P_4$    | $P_5$                  | $P_6$                  | $P_7$                  | $P_4$    | $P_5$                  | $P_6$    | $P_7$                  |
| $P_8$    | <i>P</i> 9             | $P_{10}$               | $P_{11}$               | $P_8$    | <i>P</i> 9             | $P_{10}$ | $P_{11}$               |
| $P_{12}$ | <i>P</i> <sub>13</sub> | <i>P</i> <sub>14</sub> | $P_{15}$               | $P_{12}$ | <i>P</i> <sub>13</sub> | $P_{14}$ | <i>P</i> <sub>15</sub> |

(a)

(b)

#### Irregular distribution with regular mapping! Not good.



$$V = [0, 1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11]$$
random(V) = [8, 2, 6, 0, 3, 7, 11, 1, 9, 5, 4, 10]
$$mapping = 8 \ 2 \ 6 \ 0 \ 3 \ 7 \ 111195410$$

$$P_0 \ P_1 \ P_2 \ P_3$$

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

# 2-D Randomized Distribution



2-D block random distribution.



## Irregular Allocations

- Same idea as overlap regions:
  - get data local inspector
  - Iocal 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.





Figure 3.34 A mesh used to model Lake Superior.





Random partitioning. Partitioning with minimum edge cut.

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



# Exploration of states.





- Useful data-structures.
- Usually constructed with pointers.
- Challenges for
  - communication
  - Ioad balance on irregular trees
- If little communication among sub-trees:
  - Allocate sub-trees to processes, copy the "cap".
  - All processes know the structure.





## **Dynamic Allocations**

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









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

11+14-03-2011
## 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.