

## PVM vs. MPI

## Jack Dongarra



Distinguished Professor, University of Tennessee  
Distinguished Scientist, Oak Ridge National Laboratory  
Also known for netlib, lapack, etc. Tutorial presentation:  
<http://www.netlib.org/utk/people/jd-tutorial/Presentation.html>

## PVM vs. MPI

- MPI = Message-Passing Interface
  - SPMD (Single program, multiple data)
  - Each node runs the same program
  - The program “just exists”, it is not spawned explicitly
- PVM = Parallel Virtual Machine
  - “MPMD” (Multiple program, multiple data)
  - Processes are explicitly spawned
  - Processes are assigned to nodes in separate layer, possibly multiple per node

## PVM

- Can be used for heterogeneous network
- Arbitrary topology
- Messages can cross outside of host boundaries.
- Explicit packing and unpacking of messages required in code
- Fault tolerance features
- PVM came before MPI.
- Lower level, but more flexible

## PVM

- In PVM *daemon* processes must be resident on nodes prior to spawning PVM processes there.
- Upon command, the daemon launches the process.
- The PVM host file identifies participating nodes, or they can be added manually from the command line.
- Root process is started from pvm console command-line on one host.

## PVM

- Processes explicitly spawn child processes.
- Child can determine its parent.
- Processes have their own “task id”.
- Point-to-point send/receive similar to MPI.
- Tags, wildcards similar to MPI.

### Code Fragment for simple PVM process (1) (I have left out the error checking, etc.)

```
int main(int argc, char* argv[]) {
/* find out my task id number */
mytid = pvm_mytid();

/* find my parent's task id number */
myparent = pvm_parent();

/* if I don't have a parent then I am the parent */
if (myparent == PvmNoParent) {

/* spawn the child tasks */
info = pvm_spawn(argv[0], (char**)0, PvmTaskDefault, (char*)0,
```

### Code Fragment (2)

```
/* I'm still the parent */

for (i = 0; i < ntask; i++) {
/* recv a message from any child process */
buf = pvm_recv(-1, JOINTAG);

info = pvm_buinfo(buf, &len, &tag, &tid);

info = pvm_upkint(&mydata, 1, 1);
}
pvm_exit();
}
```

### Code Fragment (3)

```
/* I'm a child */

info = pvm_initsend(PvmDataDefault);
info = pvm_pkint(&mytid, 1, 1);
info = pvm_send(myparent, JOINTAG);
pvm_exit();
}
```

See <http://www.netlib.org/pvm/book> for more details.

### PVM groups

- Processes explicitly join and leave groups, named symbolically.
- Multicast, gather, barriers, etc. are done relative to group.
- Multicast can be into group from outside.
- Reduce operator for +, \*, max, min, or *user-defined*.
- A process can be in multiple groups.

### PVM

- Multicast to explicit receives, unlike MPI.

### PVM

- For further info, examples, and on-line manual, see:
  - [http://www.coe.uncc.edu/~abw/parallel/orig\\_pvm/using\\_pvm.html](http://www.coe.uncc.edu/~abw/parallel/orig_pvm/using_pvm.html)

## Timing Analysis for Parallel Applications

## Time Decomposition

- Parallel execution time can be divided into:

- Actual computation time +
- Communication time

$$t_{\text{parallel}} = t_{\text{comp}} + t_{\text{comm}}$$

- If there are  $m$  non-parallel message steps overall, then

$$t_{\text{comm}} = m * t_{\text{message}}$$

## Message Time Decomposition

- Message time can be divided into:
  - **Latency** (or start-up time) +
  - (number of data communicated)\*(delay per datum)

$$t_{\text{message}} = t_{\text{startup}} + n * t_{\text{datum}}$$

$1/t_{\text{datum}}$  is often called "**bandwidth**", the number of data per unit time.

## Some Comparative Times (Pacheco 1997)

Machine	Arithmetic Op	Latency	Delay per Double	Ops/Latency	Period
Cray T3D	0.011 ms	21 ms	0.3 ms		1909
IBM SP-2	0.0042 ms	35 ms	0.23 ms		8333
Ethernet	N/A	500 ms	8.9 ms		N/A

## Latency Hiding

- In order to prevent  $t_{\text{message}}$  from destroying any speedup due to parallelism, we can try the following:
  - While a processing element is awaiting a message, perform some computation that doesn't require the message.
- Note that we are really trying to hide the entire communication cost, not just the "latency" component of it.

## Latency Hiding (2)

- One technique for hiding latency is "multiprogramming":
  - On a single processor, run more than one process.
  - While one process is awaiting a message, another could be doing useful computational work.
  - This requires that process-switching be relatively efficient (e.g. using threads rather than processes).
- The ratio of processes to processors is sometimes called the "parallel slackness".

## Parallel Time Complexity

- We assume familiarity with  $O$ ,  $\Omega$ , and  $\Theta$  notation.
- Their use is to bound the time complexity as a function of the problem size “ $n$ ”.

## Complexity Example

- Matrix-vector multiplication:
  - $n \times n$  matrix
  - $n$  element vector
- Assume  $n$  processors
  - Every processor has a row of the matrix
  - Each row is multiplied by the vector simultaneously
  - It takes  $O(n)$  to multiply one row, so  $t_{\text{comp}}(n) \in O(n)$

## Matrix-Vector Multiplication

- If the matrix first had to be distributed in order for the multiplication to take place, then the cost of distributing the rows from one processing element is  $O(n^2)$ , while the cost of collecting the result is  $O(n)$ .
- Therefore, the parallel cost is the same as the obvious sequential cost.

## Matrix-Vector Multiplication

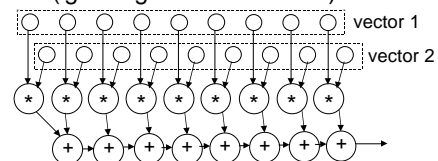
- If the same matrix is to be used many times, then the overhead of distribution gets less and less significant as the number of multiplications increases.
- In this case, the parallel cost approaches  $O(n)$ , which is an improvement over the  $O(n^2)$  serial method.

## Cost Optimality

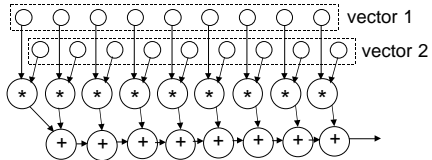
- A **cost-optimal algorithm** is defined to be one in which the **effort**, as a function of problem size, is bounded by a **constant** times the sequential effort.
- One-shot matrix-vector multiplication is not cost-optimal for distributed memory using the technique described, whereas multiplication repeated at least  $n$  times is.

## Using Graphs to Illustrate Algorithms

- The vector inner product can be shown thus (ignoring distribution cost):



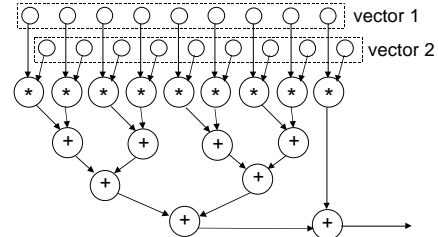
## Using Graphs to Illustrate Algorithms



- Assume unit time for each operation.
- The time is proportional to path length.
- The longest path length for an  $n$ -element vector is  $O(n)$ , sim. to serial.

## Using Graphs to Illustrate Algorithms

- Restructuring the  $+$  nodes as a **tree** gives us faster performance on  $n$  processors.



## Algorithm Analysis

- The previous tree gives us  $O(\log n)$  on  $n$  processors.
- Is it cost optimal?

## Scaling Down Processors

- As the size of the vector grows very large, we can divide the additions up among  $p$  processors,  $p \ll n$ , adding the elements within a processor sequentially and only using the tree at the end.
- The time is dominated by the sequential adds, which is  $O(n)$  for a given  $p$ .
- Is this cost optimal?

## Generalization

- Whenever the number of operations (including communication as an operation) in the parallel case is proportional to the serial complexity, we can achieve cost optimality by scaling down.
- The general concept is captured by Brent's Lemma.

## Brent's Lemma

- If an algorithm  $A$  entails  $m$  operations and can be done in parallel time  $t$  with *some* number of processors,
- then  $p$  processors can execute the algorithm in time
- $$t + (m-t)/p$$

## Brent's Lemma Summarized

- $t$  = time on some number of processors
- $m$  = number of operations (unit time each)
- time on  $p$  processors is  $\leq t + (m-t)/p$

## Application of Brent's Lemma

- To achieve cost optimality for vector inner product, use  $n/(\log n)$  processors.
- Observed that product can be done in  $\log n$  with arbitrarily-many processors.
- Brent's lemma says it can be done with  $p$  processors in

$$\underbrace{\log n}_t + \underbrace{(2n-1-\log n)}_{(m-t)} / \underbrace{(n/\log n)}_p$$

## Application of Brent's Lemma

$$\begin{aligned} & \log n + (2n-1-\log n) / (n/\log n) \\ &= \log n + 2 \log n - (\log n)/n - (\log n)^2/n \\ & \text{which is } O(\log n). \end{aligned}$$

## Proof of Brent's Lemma (1)

- Consider the graph of the algorithm done with some number of processors in time  $t$ .
- Let  $s_i$  be the number of operations done at the  $i^{\text{th}}$  level, i.e. at "time"  $i$ .
- On  $p$  processors, we can reschedule the  $s_i$  operations in time  $\text{ceiling}(s_i/p)$ .

## Proof of Brent's Lemma (2)

- On  $p$  processors, we can reschedule the  $s_i$  operations in time  $\text{ceiling}(s_i/p)$ .
- The total computation can therefore be done on  $p$  processors in time

$$\text{sum}(i = 1 \text{ to } t, \text{ceiling}(s_i/p))$$

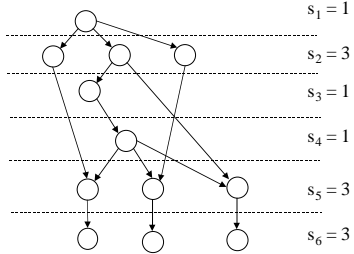
## Proof of Brent's Lemma (3)

$$\begin{aligned} & \text{sum}(i = 1 \text{ to } t, \text{ceiling}(s_i/p)) \\ & \text{is bounded by} \\ & \text{sum}(i = 1 \text{ to } t, (s_i+p-1)/p) \end{aligned}$$

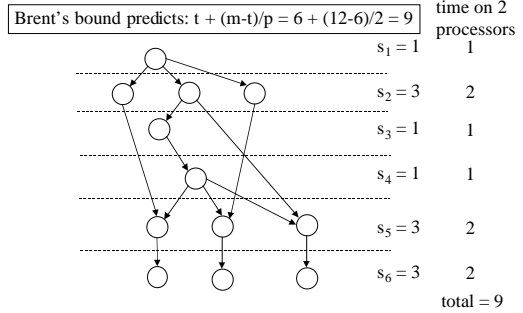
$$\begin{aligned} &= \text{sum}(i = 1 \text{ to } t, s_i/p) \\ &+ \text{sum}(i = 1 \text{ to } t, p/p) \\ &- \text{sum}(i = 1 \text{ to } t, 1/p) \end{aligned}$$

$$\begin{aligned} &= m/p + t - t/p \\ &= t + (m-t)/p, \text{ as advertised.} \end{aligned}$$

### Illustration of Brent



### Illustration of Brent for 2 processors



### Graph Exercise

- By the prefix sum problem, we mean that of computing from an array

$x_0, x_1, x_2, \dots, x_{n-1}$

the array

$(x_0), (x_0+x_1), (x_0+x_1+x_2), \dots, (x_0+x_1+x_2+\dots+x_{n-1})$

- Can this problem be sped up using parallelism?
- Is there a cost optimal version?