

Week 6

Matrix-vector Multiplication

Research Computing @ CU Boulder

Lecture Objectives

- Review matrix-vector multiplication
- Propose replication of vectors
- Develop 2 parallel programs, each based on a different data decomposition

Research Computing @ CU Boulder

Outline

- Sequential algorithm and its complexity
- Design, analysis, and implementation of three parallel programs
 - Rowwise block striped
 - Checkerboard block

Research Computing @ CU Boulder

Sequential Algorithm

$$\begin{bmatrix} 2 & 1 & 0 & 4 \\ 3 & 2 & 1 & 1 \\ 4 & 3 & 1 & 2 \\ 3 & 0 & 2 & 0 \end{bmatrix} \times \begin{bmatrix} 1 \\ 3 \\ 4 \\ 1 \end{bmatrix} = \begin{bmatrix} 9 \\ 14 \\ 19 \\ 11 \end{bmatrix}$$

Research Computing @ CU Boulder

Storing Vectors

- Divide vector elements among processes
- Replicate vector elements
- Vector replication acceptable because vectors have only n elements, versus n^2 elements in matrices

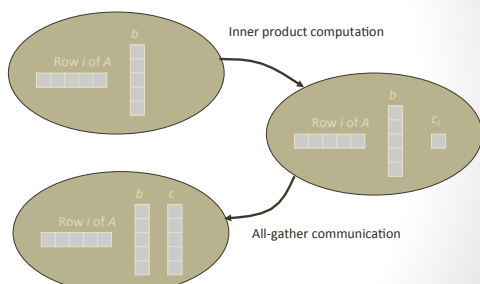
Research Computing @ CU Boulder

Rowwise Block Striped Matrix

- Partitioning through domain decomposition
- Primitive task associated with
 - Row of matrix
 - Entire vector

Research Computing @ CU Boulder

Phases of Parallel Algorithm



Research Computing @ CU Boulder

Agglomeration and Mapping

- Static number of tasks
- Regular communication pattern (all-gather)
- Computation time per task is constant
- Strategy:
 - Agglomerate groups of rows
 - Create one task per MPI process

Research Computing @ CU Boulder

Complexity Analysis

- Sequential algorithm complexity: $\Theta(n^2)$
- Parallel algorithm computational complexity: $\Theta(n^2/p)$
- Communication complexity of all-gather: $\Theta(\log p + n)$
- Overall complexity: $\Theta(n^2/p + \log p)$

Research Computing @ CU Boulder

Isoefficiency Analysis

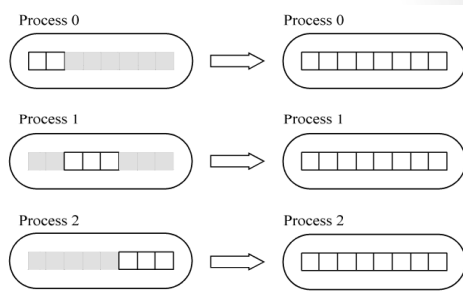
- Sequential time complexity: $\Theta(n^2)$
- Only parallel overhead is all-gather
 - When n is large, message transmission time dominates message latency
 - Parallel communication time: $\Theta(n)$
- $n^2 \geq Cpn \Rightarrow n \geq Cp$ and $M(n) = n^2$

$$M(Cp)/p = C^2 p^2 / p = C^2 p$$

- System is not highly scalable

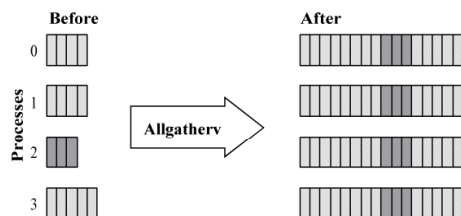
Research Computing @ CU Boulder

Block-to-replicated Transformation



Research Computing @ CU Boulder

MPI_Allgather



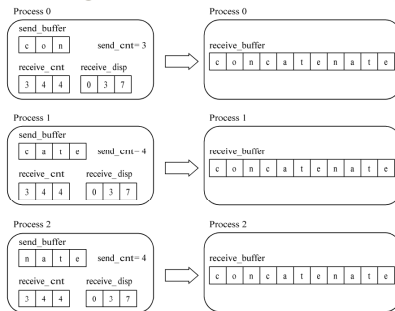
Research Computing @ CU Boulder

MPI_Allgatherv

```
int MPI_Allgatherv (
    void          *send_buffer,
    int           send_cnt,
    MPI_Datatype   send_type,
    void          *receive_buffer,
    int           *receive_cnt,
    int           *receive_disp,
    MPI_Datatype   receive_type,
    MPI_Comm       communicator)
```

Research Computing @ CU Boulder

MPI_Allgatherv in Action



Research Computing @ CU Boulder

Function create_mixed_xfer_arrays

- First array
 - How many elements contributed by each process
 - Uses utility macro BLOCK_SIZE
- Second array
 - Starting position of each process' block
 - Assume blocks in process rank order

Research Computing @ CU Boulder

Function replicate_block_vector

- Create space for entire vector
- Create “mixed transfer” arrays
- Call `MPI_Allgatherv`

Research Computing @ CU Boulder

Function read_replicated_vector

- Process $p-1$
 - Opens file
 - Reads vector length
- Broadcast vector length (root process = $p-1$)
- Allocate space for vector
- Process $p-1$ reads vector, closes file
- Broadcast vector

Research Computing @ CU Boulder

Function print_replicated_vector

- Process 0 prints vector
- Exact call to `printf` depends on value of parameter `datatype`

Research Computing @ CU Boulder

Run-time Expression

- χ : inner product loop iteration time
- Computational time: $\chi n[p]$
- All-gather requires $\lceil \log p \rceil$ messages with latency λ .
- Total vector elements transmitted:
 $(2^{\lceil \log p \rceil} - 1) / 2^{\lceil \log p \rceil}$
- Total execution time: $\chi n[p] + \lambda \lceil \log p \rceil + (2^{\lceil \log p \rceil} - 1) / (2^{\lceil \log p \rceil} \beta)$

Research Computing @ CU Boulder

Benchmarking Results

p	Execution Time (msec)		Speedup	Mflops
	Predicted	Actual		
1	63.4	63.4	1.00	31.6
2	32.4	32.7	1.94	61.2
3	22.3	22.7	2.79	88.1
4	17.0	17.8	3.56	112.4
5	14.1	15.2	4.16	131.6
6	12.0	13.3	4.76	150.4
7	10.5	12.2	5.19	163.9
8	9.4	11.1	5.70	180.2
16	5.7	7.2	8.79	277.8

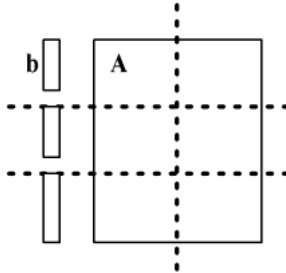
Research Computing @ CU Boulder

Checkerboard Block Decomposition

- Associate primitive task with each element of the matrix **A**
- Each primitive task performs one multiply
- Agglomerate primitive tasks into rectangular blocks
- Processes form a 2-D grid
- Vector **b** distributed by blocks among processes in first column of grid

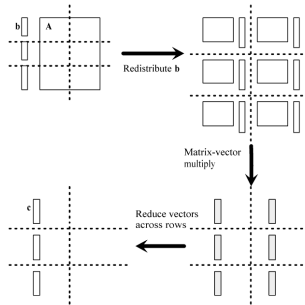
Research Computing @ CU Boulder

Tasks after Agglomeration



Research Computing @ CU Boulder

Algorithm's Phases

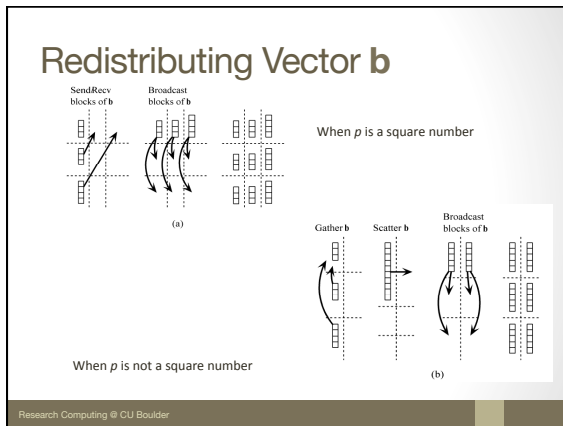


Research Computing @ CU Boulder

Redistributing Vector \mathbf{b}

- Step 1: Move \mathbf{b} from processes in first row to processes in first column
 - If p square
 - First column/first row processes send/receive portions of \mathbf{b}
 - If p not square
 - Gather \mathbf{b} on process 0, 0
 - Process 0, 0 broadcasts to first row procs
- Step 2: First row processes scatter \mathbf{b} within columns

Research Computing @ CU Boulder



Complexity Analysis

- Assume p is a square number
 - If grid is $1 \times p$, devolves into columnwise block striped
 - If grid is $p \times 1$, devolves into rowwise block striped

Research Computing @ CU Boulder

Complexity Analysis (continued)

- Each process does its share of computation: $\Theta(n^2/p)$
- Redistribute **b**: $\Theta(n / \sqrt{p} + \log p(n / \sqrt{p})) = \Theta(n \log p / \sqrt{p})$
- Reduction of partial results vectors: $\Theta(n \log p / \sqrt{p})$
- Overall parallel complexity: $\Theta(n^2/p + n \log p / \sqrt{p})$

Research Computing @ CU Boulder

Isoefficiency Analysis

- Sequential complexity: $\Theta(n^2)$
- Parallel communication complexity: $\Theta(n \log p / \sqrt{p})$
- Isoefficiency function:
 $n^2 \geq Cn \sqrt{p} \log p \Rightarrow n \geq C \sqrt{p} \log p$

$$M(C\sqrt{p} \log p) / p = C^2 p \log^2 p / p = C^2 \log^2 p$$
- This system is much more scalable than the previous implementation

Research Computing @ CU Boulder

Creating Communicators

- Want processes in a virtual 2-D grid
- Create a custom communicator to do this
- Collective communications involve all processes in a communicator
- We need to do broadcasts, reductions among subsets of processes
- We will create communicators for processes in same row or same column

Research Computing @ CU Boulder

What's in a Communicator?

- Process group
- Context
- Attributes
 - Topology (lets us address processes another way)
 - Others we won't consider

Research Computing @ CU Boulder

Creating 2-D Virtual Grid of Processes

- `MPI_Dims_create`
 - Input parameters
 - Total number of processes in desired grid
 - Number of grid dimensions
 - Returns number of processes in each dim
- `MPI_Cart_create`
 - Creates communicator with cartesian topology

Research Computing @ CU Boulder

`MPI_Dims_create`

```
int MPI_Dims_create (
    int nodes,
        /* Input - Procs in grid */

    int dims,
        /* Input - Number of dims */

    int *size)
    /* Input/Output - Size of
       each grid dimension */
```

Research Computing @ CU Boulder

`MPI_Cart_create`

```
int MPI_Cart_create (
    MPI_Comm old_comm, /* Input - old communicator */

    int dims, /* Input - grid dimensions */

    int *size, /* Input - # procs in each dim */

    int *periodic,
        /* Input - periodic[j] is 1 if dimension j
           wraps around; 0 otherwise */

    int reorder,
        /* 1 if process ranks can be reordered */

    MPI_Comm *cart_comm)
    /* Output - new communicator */
```

Research Computing @ CU Boulder

Using MPI_Dims_create and MPI_Cart_create

```
MPI_Comm cart_comm;
int p;
int periodic[2];
int size[2];
...
size[0] = size[1] = 0;
MPI_Dims_create (p, 2, size);
periodic[0] = periodic[1] = 0;
MPI_Cart_create (MPI_COMM_WORLD, 2, size,
                 1, &cart_comm);
```

Research Computing @ CU Boulder

Useful Grid-related Functions

- MPI_Cart_rank
 - Given coordinates of process in Cartesian communicator, returns process rank
- MPI_Cart_coords
 - Given rank of process in Cartesian communicator, returns process' coordinates

Research Computing @ CU Boulder

Header for MPI_Cart_rank

```
int MPI_Cart_rank (
    MPI_Comm comm,
    /* In - Communicator */
    int *coords,
    /* In - Array containing process'
       grid location */
    int *rank)
    /* Out - Rank of process at
       specified coords */
```

Research Computing @ CU Boulder

Header for MPI_Cart_coords

```
int MPI_Cart_coords (
    MPI_Comm comm,
    /* In - Communicator */
    int rank,
    /* In - Rank of process */
    int dims,
    /* In - Dimensions in virtual grid */
    int *coords)
    /* Out - Coordinates of specified
    process in virtual grid */
```

Research Computing @ CU Boulder

MPI_Comm_split

- Partitions the processes of a communicator into one or more subgroups
- Constructs a communicator for each subgroup
- Allows processes in each subgroup to perform their own collective communications
- Needed for columnwise scatter and rowwise reduce

Research Computing @ CU Boulder

Header for MPI_Comm_split

```
int MPI_Comm_split (
    MPI_Comm old_comm,
    /* In - Existing communicator */
    int partition, /* In - Partition number */
    int new_rank,
    /* In - Ranking order of processes
    in new communicator */
    MPI_Comm *new_comm)
    /* Out - New communicator shared by
    processes in same partition */
```

Research Computing @ CU Boulder

Example: Create Communicators for Process Rows

```
MPI_Comm grid_comm; /* 2-D process grid */
MPI_Comm grid_coors[2];
    /* Location of process in grid */

MPI_Comm row_comm;
    /* Processes in same row */

MPI_Comm_split (grid_comm, grid_coors[0],
    grid_coors[1], &row_comm);
```

Research Computing @ CU Boulder

Run-time Expression

- Computational time: $\chi \lceil n/\sqrt{p} \rceil \lceil n/\sqrt{p} \rceil$
- Suppose p a square number
- Redistribute **b**
 - Send/recv: $\lambda + 8 \lceil n/\sqrt{p} \rceil / \beta$
 - Broadcast: $\log \sqrt{p} (\lambda + 8 \lceil n/\sqrt{p} \rceil / \beta)$
- Reduce partial results:
 $\log \sqrt{p} (\lambda + 8 \lceil n/\sqrt{p} \rceil / \beta)$

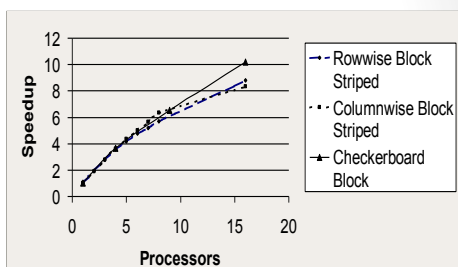
Research Computing @ CU Boulder

Benchmarking

Procs	Predicted (msec)	Actual (msec)	Speedup	Megaflops
1	63.4	63.4	1.00	31.6
4	17.8	17.4	3.64	114.9
9	9.7	9.7	6.53	206.2
16	6.2	6.2	10.21	322.6

Research Computing @ CU Boulder

Comparison of Three Algorithms



Research Computing @ CU Boulder

Summary (1/3)

- Matrix decomposition \Rightarrow communications needed
 - Rowwise block striped: all-gather
 - Columnwise block striped: all-to-all exchange
 - Checkerboard block: gather, scatter, broadcast, reduce
- All three algorithms: roughly same number of messages
- Elements transmitted per process varies
 - First two algorithms: $\Theta(n)$ elements per process
 - Checkerboard algorithm: $\Theta(n/\sqrt{p})$ elements
- Checkerboard block algorithm has better scalability

Research Computing @ CU Boulder

Summary (2/3)

- Communicators with Cartesian topology
 - Creation
 - Identifying processes by rank or coords
- Subdividing communicators
 - Allows collective operations among subsets of processes

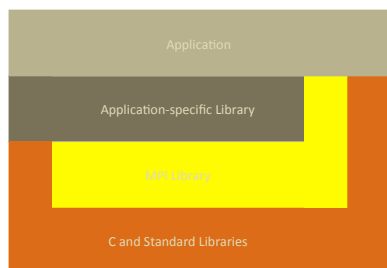
Research Computing @ CU Boulder

Summary (3/3)

- Parallel programs and supporting functions much longer than C counterparts
- Extra code devoted to reading, distributing, printing matrices and vectors
- Developing and debugging these functions is tedious and difficult
- Makes sense to generalize functions and put them in libraries for reuse

Research Computing @ CU Boulder

MPI Application Development



Research Computing @ CU Boulder
