# CME213/ME339
# Lecture 20

Eric Darve
Erich Elsen
Austin Gibbons

Department of Mechanical Engineering
Institute for Computational and Mathematical Engineering
Stanford University

Spring 2012

# Example: row-wise matrix-vector product

```
1   RowMatrixVectorMultiply(int n, double *a, double *b, double *x,
2                           MPI_Comm comm) {
3   /* a stores the rows of the matrix that this process
4      is going to compute with.
5      b stores only the part of b owned by this process */
6   int i, j;
7   int nlocal; /* Number of locally stored rows of A */
8   double *fb; /* Will store the entire vector b */
9   int nprocs, myrank;
10  MPI_Status status;
11
12  /* Get information about the communicator */
13  MPI_Comm_size(comm, &nprocs);
14  MPI_Comm_rank(comm, &myrank);
15
16  /* Allocate the memory used to store the entire b */
17  fb = (double *)malloc(n*sizeof(double));
18  nlocal = n/nprocs; // nprocs must divide n
```

```
1   /* Gather entire vector b on each processor using Allgather */
2   MPI_Allgather(b, nlocal, MPI_DOUBLE, fb, nlocal, MPI_DOUBLE,
3                 comm);
4
5   /* Perform the matrix-vector multiplication involving the
6      locally stored submatrix. */
7   for (i=0; i<nlocal; i++) {
8     x[i] = 0.0;
9     for (j=0; j<n; j++)
10       x[i] += a[i*n+j]*fb[j];
11  }
12  /* Done! */
13  free(fb);
14  }
```

# Example: column-wise matrix-vector product

```
1   ColMatrixVectorMultiply(int n, double *a, double *b, double *x,
2                          MPI_Comm comm) {
3   int i, j;
4   int nlocal;
5   double *px;
6   double *fx;
7   int nprocs, myrank;
8   MPI_Status status;
9   /* Get identity and size information from the communicator */
10  MPI_Comm_size(comm, &nprocs);
11  MPI_Comm_rank(comm, &myrank);
12
13  nlocal = n/nprocs;
14
15  /* Allocate memory for arrays storing intermediate results. */
16  px = (double *)malloc(n*sizeof(double));
17  fx = (double *)malloc(n*sizeof(double));
```

```
1    /* Compute the partial-dot products that correspond to the
2       local columns of A.*/
3    for (i=0; i<n; i++) {
4      px[i] = 0.0;
5      for (j=0; j<nlocal; j++)
6        px[i] += a[i*nlocal+j]*b[j];
7    }
8
9    /* Sum-up the results by performing an element-wise
10      reduction operation. Result is stored at 0. */
11   MPI_Reduce(px, fx, n, MPI_DOUBLE, MPI_SUM, 0, comm);
12
13   /* Redistribute fx in a fashion similar to that of vector b.
14      Data is sent from 0 to all other processes. */
15   MPI_Scatter(fx, nlocal, MPI_DOUBLE, x, nlocal, MPI_DOUBLE, 0,
16               comm);
17   free(px); free(fx);
18   }
```

# Virtual topologies

Many problems are naturally mapped to certain topologies such as grids. This is the case for example for matrices, or for 2D or 3D structured grids.

The two main types of topologies supported by MPI are Cartesian grids and graphs.

MPI topologies are virtual — there may be no relation between the physical structure of the parallel machine and the process topology.

# Why?

Convenience: virtual topologies may be useful for applications with specific communication patterns — patterns that match an MPI topology structure.

Communication efficiency: a particular implementation may optimize process mapping based upon the physical characteristics of a given parallel machine. For example nodes that are nearby on the grid (East/West/North/South neighbors) will be as close as possible in the network (lowest communication time).

The mapping of processes into an MPI virtual topology is dependent upon the MPI implementation.

# MPI functions for virtual topologies

```
1  int MPI_Cart_create(MPI_Comm comm_old, int ndims, int *dims,
2                      int *periods, int reorder,
3                      MPI_Comm *comm_cart)
```

The topology information is attached to a new communicator
comm_cart that is created by MPI_Cart_create.

The array periods is used to specify whether or not the topology
has wraparound connections. In particular, if periods[i] is
non-zero, then the topology has wraparound connections along
dimension i, otherwise it does not.

The argument reorder is used to determine if the processes in the new group (i.e., communicator) are to be reordered or not.

If reorder is false, then the rank of each process in the new group is identical to its rank in the old group.

Otherwise, MPI_Cart_create may reorder the processes if that leads to a better embedding of the virtual topology onto the parallel computer.

It will result in an error if the total number of processes specified by dims is greater than the number of processes in the comm_old communicator.

# Other useful functions

```
1  int MPI_Cart_rank(MPI_Comm comm_cart, int *coords, int *rank)
2  int MPI_Cart_coords(MPI_Comm comm_cart, int rank, int maxdims,
3                      int *coords)
```

coords are the Cartesian coordinates of a process. Its size is the number of dimensions.

```
1  int MPI_Cart_shift(MPI_Comm comm_cart, int dir, int s_step,
2                     int *rank_source, int *rank_dest)
```

dir: direction; s_step: length shift.

rank_dest contains the group rank of the neighboring process in
the specified dimension and distance. The rank of the process for
which the calling process is the neighboring process in the specified
dimension and distance is returned rank_source.

Thus, the group ranks returned in rank dest and rank source can
be used as parameters for MPI_Sendrecv().

# Groups and communicators

This is a more advanced topic. I will only cover it superficially. More information can be found online.

The main point is that collective communications assume that we have identified a certain group of processes that want to communicate.

By default this group contains all the processes.

This can be changed using groups. For example we may want to communicate along the rows of a matrix or the columns.
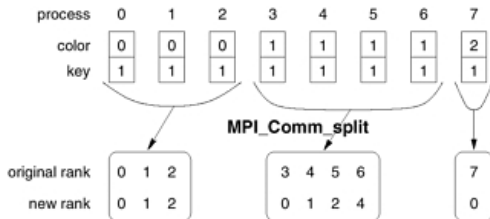
# Split

The simplest method to create groups:

```
int MPI_Comm_split(MPI_Comm comm, int color, int key,
                   MPI_Comm *newcomm)
```

Each subgroup contains all processes that have supplied the same value for the color parameter. Within each subgroup, the processes are ranked in the order defined by the value of the key parameter, with ties broken according to their rank in the old communicator (i.e., comm). A new communicator for each subgroup is returned in the newcomm parameter.

# Splitting a Cartesian topology

It is very common that one wants to split a Cartesian topology along certain dimensions.
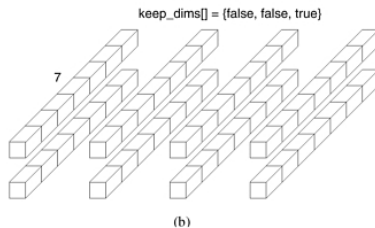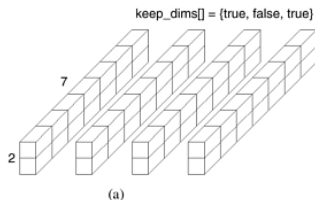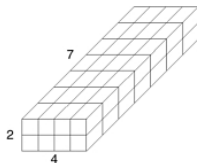
```
1  int MPI_Cart_sub(MPI_Comm comm_cart, int *keep_dims,
2                   MPI_Comm *comm_subcart)
```

`keep_dims`: boolean flag that determines whether that dimension is retained in the new communicators or split, e.g., if false then a split occurs.

# Example

The grid has size $2 \times 4 \times 7$.



keep_dims[] = {true, false, true}

(a)

keep_dims[] = {false, false, true}

(b)

# Matrix-vector multiplication: 2D

```
1  MatrixVectorMultiply_2D(int n, double *a, double *b, double *x,
2      MPI_Comm comm) {
3  int ROW=0, COL=1; /* Improve readability */
4  int i, j, nlocal;
5  double *px; /* Will store partial dot products */
6  int nprocs, dims[2], periods[2], keep_dims[2];
7  int myrank, my2drank, mycoords[2];
8  int drank, col0rank, coords[2];
9  MPI_Status status;
10 MPI_Comm comm_2d, comm_row, comm_col;
11
12 /* Get information about the communicator */
13 MPI_Comm_size(comm, &nprocs); MPI_Comm_rank(comm, &myrank);
14
15 /* Compute the size of the square grid.
16    We assume that nprocs is a square and that the matrix size
17    is a multiple of sqrt(nprocs). */
18 dims[ROW] = dims[COL] = sqrt(nprocs);
19 nlocal    = n/dims[ROW];
20 ...
```

```
1   /* Allocate memory for the array that will hold
2      the partial dot-products */
3   px = malloc(nlocal*sizeof(double));
4
5   /* Set up the Cartesian topology and get the rank &
6      coordinates of the process in this topology */
7   periods[ROW] = periods[COL] = 1;
8   /* Periods for wrap-around connections */
9
10  MPI_Cart_create(MPI_COMM_WORLD, 2, dims, periods, 1, &comm_2d);
11
12  /* Get my rank in the new topology */
13  MPI_Comm_rank(comm_2d, &my2drank);
14
15  /* Get my coordinates */
16  MPI_Cart_coords(comm_2d, my2drank, 2, mycoords);
17  ...
```

```
1   /* Create the row-based sub-topology */
2   keep_dims[ROW] = 0;
3   keep_dims[COL] = 1;
4   MPI_Cart_sub(comm_2d, keep_dims, &comm_row);
5
6   /* Create the column-based sub-topology */
7   keep_dims[ROW] = 1;
8   keep_dims[COL] = 0;
9   MPI_Cart_sub(comm_2d, keep_dims, &comm_col);
10   ...
```

```
1   /* Redistribute the b vector. */
2   /* Step 1. The processors along the 0th column
3      send their data to the diagonal processor. */
4   if (mycoords[COL] == 0 && mycoords[ROW] != 0) {
5     /* I'm in the first column */
6     coords[ROW] = mycoords[ROW];
7     coords[COL] = mycoords[ROW]; // coords of diagonal block
8     MPI_Cart_rank(comm_2d, coords, &drank);
9     /* Send data to the diagonal block */
10    MPI_Send(b, nlocal, MPI_DOUBLE, drank, 1, comm_2d);
11  }
12
13  if (mycoords[ROW] == mycoords[COL] && mycoords[ROW] != 0) {
14    /* I am a diagonal block */
15    coords[ROW] = mycoords[ROW];
16    coords[COL] = 0; // Receiving from column 0
17    MPI_Cart_rank(comm_2d, coords, &col0rank);
18    MPI_Recv(b, nlocal, MPI_DOUBLE, col0rank, 1, comm_2d,
19        &status);
20  }
21  ...
```

```
1   /* Step 2. The diagonal processors perform a
2               column-wise broadcast */
3   coords[0] = mycoords[COL];
4   /* Note the 1D rank in the column sub-topology */
5   MPI_Cart_rank(comm_col, coords, &drank);
6   MPI_Bcast(b, nlocal, MPI_DOUBLE, drank, comm_col);
7
8   /* Get into the main computational loop */
9   for (i=0; i<nlocal; i++) {
10    px[i] = 0.0;
11    for (j=0; j<nlocal; j++)
12      px[i] += a[i*nlocal+j]*b[j];
13  }
14  ...
```

```
1   /* Perform the sum-reduction along the rows to add up
2      the partial dot-products */
3   coords[0] = 0;
4   /* Note the 1D rank in the row sub-topology */
5   MPI_Cart_rank(comm_row, coords, &col0rank);
6   MPI_Reduce(px, x, nlocal, MPI_DOUBLE, MPI_SUM, col0rank,
7              comm_row);
8
9   /* Free up communicators */
10  MPI_Comm_free(&comm_2d);
11  MPI_Comm_free(&comm_row);
12  MPI_Comm_free(&comm_col);
13  free(px);
14  }
```

## Derived data types

We do not have time to cover that topic in detail.

This can be used to communicate data that is not simply contiguous in memory.

The most useful are:

```
1   int MPI_Type_vector(
2         int count, int blocklen, int stride,
3         MPI_Datatype old_type, MPI_Datatype *newtype)
4
5   int MPI_Type_indexed(
6         int count, int blocklens[], int indices[],
7         MPI_Datatype old_type, MPI_Datatype *newtype)
```

# Short example

```
1   ...
2   float a[SIZE][SIZE] =
3     {1.0, 2.0, 3.0, 4.0,
4      5.0, 6.0, 7.0, 8.0,
5      9.0, 10.0, 11.0, 12.0,
6      13.0, 14.0, 15.0, 16.0};
7   ...
8   MPI_Datatype columntype;
9   MPI_Type_vector(SIZE, 1, SIZE, MPI_FLOAT, &columntype);
10  MPI_Type_commit(&columntype);
11  ...
12  MPI_Send(&a[0][i], 1, columntype, i, tag, MPI_COMM_WORLD);
13  ...
```

# Performance metrics

There are a few key concepts that help understand the running time of parallel applications.

There are several components in the total run time $T_p(n)$:

- Local computations
- Data exchange and communication
- Waiting time because of load imbalance

## Cost

The cost $C_p(n)$ of a parallel program is:

$$C_p(n) = p\ T_p(n)$$

This is a measure of the total amount of work done for all processors.

This cost is optimal if

$$C_p(n) = T^\star(n),$$

the optimal runtime of the sequential algorithm.

# Speedup

Defined as:
$$S_p(n) = \frac{T^\star(n)}{T_p(n)}$$

This represents the savings obtained from running in parallel.

In some rare cases superlinear effect can be observed, that is $S_p(n) > p$. This is in most cases due to cache effect. With $p$ processes we have effectively access to more cache space than on a single processor, resulting in some cases in better performance.

# Efficiency

It makes sense to compare the speed-up against $p$. Ideally the speed-up scales like $p$. To better illustrate this dependence, it is common to define the efficiency of a parallel code:

$$E_p(n) = \frac{T^\star(n)}{C_p(n)} = \frac{S_p(n)}{p} = \frac{T^\star(n)}{p T_p(n)}$$

The ideal speed-up corresponds to $E_p = 1$.

# Amdahl's law

Although this is a crude model, it can provide a general sense of what speed-up can be achieved. This is a good measure to understand how much of the code one should try to parallelize.

Assume that a fraction $f$ of the code is executed sequentially. Then the speed-up is given by:

$$S_p(n) = \frac{T^\star(n)}{f\ T^\star(n) + ((1-f)/p)\ T^\star(n)} = \frac{1}{f + (1-f)/p} \leq \frac{1}{f}$$

For example, if $f = 20\%$, the maximum speed-up that can be achieved is 5, no matter how good our MPI implementation is, or how many processors we can employ.

This sounds a bit depressing. How small can $f$ be?

The good news is that $f$ typically depends on $n$ and goes to 0 as $n$ increases.

This means that to observe a certain speed-up on $p$ processes we may need to increase the problem size $n$ sufficiently.

Increasing $n$ often also allows reducing the overhead in the calculation from communication and idle times.

# Gustafson's law

This model assumes that the sequential part has a constant execution time $\tau_f$ irrespective of $n$. Then if the parallel part can be perfectly parallelized:

$$S_p(n) = \frac{T^\star(n)}{\tau_f + (T^\star(n) - \tau_f)/p}$$

Therefore in this model we always have

$$\lim_{n \to \infty} S_p(n) = p$$

# Short example: dot product

Goal: calculate a dot product and store the result on one node.

$\alpha$: time to perform a multiplication or addition (we assume these times are equal for simplicity).
$\beta$: time to communicate a float to a neighboring process.

Two steps algorithm:

1. Calculate local dot product: $\sum_j a_j b_j$
2. Use a spanning tree: $\ln_2 p$ passes, each with one addition and one communication

Total serial run time: $T^\star = 2\alpha n$
Total parallel run time: $T_p = 2\alpha n/p + (\alpha + \beta)\ln_2 p$

# Efficiency and iso-efficiency

Efficiency:

$$E_p(n) = \frac{2\alpha n}{2\alpha n + (\alpha + \beta)\, p \ln_2 p} = \frac{1}{1 + ((\alpha + \beta)/2\alpha)\, p \ln_2 p / n}$$

Efficiency is maintained (iso-efficiency) as long as:

$$p \ln_2 p = \Theta(n) \quad \text{or} \quad p = \Theta(n / \ln_2 n)$$