Technische Universität München

# Assignment 4: MPI Collectives and MPI-IO

Programming of Super Computers

Friedrich Menhorn, Benjamin Rüth, Erik Wannerberg
Team 12

February 2, 2016

Friedrich Menhorn, Benjamin Rüth, Erik Wannerberg: Assignment 4: MPI Collectives and MPI-IO
Programming of Super Computers, February 2, 2016

1

# Contents

ТИП

**Friedrich Menhorn, Benjamin Rüth, Erik Wannerberg: Assignment 4: MPI Collectives and MPI-IO**
**Programming of Super Computers, February 2, 2016**

**3**

# Optimizations

Sending the dimensions of the matrix to all processes:

**Original code:**

```
// send dimensions to all peers
if(rank == 0) {
    int i;
    for(i = 1; i < size; i++){
        MPI_Send(matrices_a_b_dimensions
            , 4, MPI_INT, i, 0,
            cartesian_grid_communicator
            );
    }
} else {
    MPI_Recv(matrices_a_b_dimensions, 4,
        MPI_INT, 0, 0,
        cartesian_grid_communicator, &
        status);
}
```

**With collective operation:**

```
// send dimensions to all peers
MPI_Bcast(matrices_a_b_dimensions, 4,
    MPI_INT, 0,
    cartesian_grid_communicator);
```

# Optimizations (contd.)

Distribute matrix blocks among all processes:

**Original code:**

```
// send a block to each process
if(rank == 0) {
    int i;
    for(i = 1; i < size; i++){
        MPI_Send(A_array, ...);
        MPI_Send(B_array, ...);
    }
    for(i = 0; i < A_local_block_size; i
        ++){
        A_local_block[i] = A_array[i];
    }
    for(i = 0; i < B_local_block_size; i
        ++){
        B_local_block[i] = B_array[i];
    }
} else {
    MPI_Recv(A_local_block, ...);
    MPI_Recv(B_local_block, ...);
}
```

**With collective operation:**

```
// send a block to each process
MPI_Scatter(A_array, A_local_block_size,
        MPI_DOUBLE, A_local_block,
        A_local_block_size, MPI_DOUBLE, 0,
        cartesian_grid_communicator);
MPI_Scatter(B_array, B_local_block_size,
        MPI_DOUBLE, B_local_block,
        B_local_block_size, MPI_DOUBLE, 0,
        cartesian_grid_communicator);
```

**Friedrich Menhorn, Benjamin Rüth, Erik Wannerberg: Assignment 4: MPI Collectives and MPI-IO**
**Programming of Super Computers, February 2, 2016**

5

# Optimizations (contd.)

Distribute matrix blocks among all processes and do initial rotation of blocks:

**Original code:**

```
// send a block to each process
...
// fix initial arrangements before the
    core algorithm starts
if(coordinates[0] != 0){
    MPI_Sendrecv_replace(A_local_block,
        ...);
}
if(coordinates[1] != 0){
    MPI_Sendrecv_replace(B_local_block,
        ...);
}
```

**With collective operation:**

```
// send a block to each process and fix
    initial arrangements before the
    core algorithm starts
int *send_count= (int *) malloc(...);
int *displs_A = (int *) malloc(...);
int *displs_B = (int *) malloc(...);
int displ_coord_A_col = ...
int displ_coord_B_row = ...
int displs_local_A = ...
int displs_local_B = ...

MPI_Gather(&displs_local_A, ...);
MPI_Gather(&displs_local_B, ...);

for(i = 0; i < size; i++){
send_count[i] = A_local_block_size;
}

MPI_Scatterv(A_array, ...);
MPI_Scatterv(B_array, ...);
```

# Optimizations (contd.)

Collect results after computation loop:

**Original code:**

```
// get C parts from other processes at
    rank 0
if(rank == 0) {
    for(i = 0; i < A_local_block_rows *
        B_local_block_columns; i++){
        C_array[i] = C_local_block[i];
    }
    int i;
    for(i = 1; i < size; i++){
        MPI_Recv(....);
    }
} else {
    MPI_Send(...);
}
```

**With collective operation:**

```
// get C parts from other processes at
    rank 0
int C_local_block_size =
    A_local_block_size;
MPI_Gather(C_local_block,
    C_local_block_size, MPI_DOUBLE,
    C_array, C_local_block_size,
    MPI_DOUBLE, 0,
    cartesian_grid_communicator);
```

**Friedrich Menhorn, Benjamin Rüth, Erik Wannerberg: Assignment 4: MPI Collectives and MPI-IO**
**Programming of Super Computers, February 2, 2016**

7

ПЛП

# Optimizations (contd.)

### Expectations

- **Performance:**
  Higher performance due to highly optimized MPI collectives (might also consider network topology etc.) vs. naive implementation via `MPI_Send`/`MPI_Recv`.

- **Scaling:**
  Efficient implementation of Broadcast, Gather, Scatter use treelike distribution of values, $\mathcal{O}(n) \rightarrow \mathcal{O}(\log(n))$. This means we also benefit with respect to scalability.

# Optimizations (contd.)

### Readability and Maintainability

- ⊕ Sending dimensions with `MPI_Bcast`
- ⊕ Distributing matrix blocks `MPI_Scatter`
- ⊖ Distributing matrix blocks and initial alignment with `MPI_Scatterv`
  - → introduces complicated offsets
  - → hard to see what actually happens (distribution & rotation)
- ⊕ Collecting results using `MPI_Gather`

**Friedrich Menhorn, Benjamin Rüth, Erik Wannerberg: Assignment 4: MPI Collectives and MPI-IO**
**Programming of Super Computers, February 2, 2016**

9

# Performance Measurements

Show some measurements for falsification/verification of our guess

**Friedrich Menhorn, Benjamin Rüth, Erik Wannerberg: Assignment 4: MPI Collectives and MPI-IO**
**Programming of Super Computers, February 2, 2016**

10

**Friedrich Menhorn, Benjamin Rüth, Erik Wannerberg: Assignment 4: MPI Collectives and MPI-IO**
**Programming of Super Computers, February 2, 2016**

11

# Data Sieving and 2-Phase IO

What is "Data Sieving" and "2-Phase IO"? How do they help improve IO performance?

# Scalability of original application

- Was the original implementation scalable in terms of IO performance?
- Was the original implementation scalable in terms of RAM storage?

# Impact of parallel IO

How much of the communication in the application was replaced with MPI-IO operations?