

Introduction to Parallel and Distributed Processing

Custom Datatypes and I/O

Jaroslav 'Jaric' Zola

<http://www.jzola.org/>

Partitioning Communication

- Often, it is advantageous to organize processors into logical groups, e.g. processors in the same row/column
- Each group can act “independently” of any other group (collectives)
- MPI provides mechanism to create new communicator from the existing one

Creating Communicator

```
1 // ex01.cpp
2 #include <iostream>
3 #include <mpi.h>
4
5 int main(int argc, char* argv[]) {
6     MPI_Init(&argc, &argv);
7
8     int rank, size;
9
10    MPI_Comm_size(MPI_COMM_WORLD, &size);
11    MPI_Comm_rank(MPI_COMM_WORLD, &rank);
12
13    int col = rank % 4;
14    int row = rank >> 2;
15
16    MPI_Comm col_comm;
17    MPI_Comm_split(MPI_COMM_WORLD, col, rank, &col_comm);
18
19    int nrank;
20    MPI_Comm_rank(col_comm, &nrank);
21
22    std::cout << rank << " " << row << " " << col << " "
23              << nrank << std::endl;
24
25    MPI_Comm_free(&col_comm);
26
27    return MPI_Finalize();
28 } // main
```

Custom Datatypes

- Same but different: `sizeof(A)=12`, `sizeof(B)=8`

```
1 struct A {  
2     char c0;  
3     short int si0;  
4     int i0;  
5     char c1;  
6 };
```

```
1 struct B {  
2     char c0;  
3     char c1;  
4     short int si0;  
5     int i0;  
6 };
```

MPI Datatypes

- Abstractly, type represented via type map:
$$Type = \{(type_0, disp_0), \dots, (type_{n-1}, disp_{n-1})\}$$
- Examples:
 - `MPI_DOUBLE`
 $\{(double, 0)\}$
 - `struct X { double d; char c; };`
 $\{(double, 0), (char, 8)\}$

Easiest Type Creator

- Replicate *count* times old type
- Example: $old = \{(int, 0), (double, 8)\}$, with $count = 2$
 $new = \{(int, 0), (double, 8), (int, 16), (double, 24)\}$

```
1 | MPI_Datatype MPI_INT2;  
2 | MPI_Type_contiguous(2, MPI_INT, &MPI_INT2);  
3 | MPI_Type_commit(&MPI_INT2);  
4 |  
5 | if (rank == 0) {  
6 |     int tab[256];  
7 |     MPI_Send(tab, 1, MPI_INT2, 1, 11, MPI_COMM_WORLD);  
8 | } else if (rank == 1) {  
9 |     int buf[2];  
10 |     MPI_Status stat;  
11 |     MPI_Recv(buf, 1, MPI_INT2, 0, 11, MPI_COMM_WORLD, &stat);  
12 | }  
13 |  
14 | MPI_Type_free(&MPI_INT2);
```

Brute-force Approach

- Consider each type as a stream of bytes:

```
1  struct X {  
2      char c;  
3      double d;  
4  };  
5  
6  MPI_Datatype MPI_X;  
7  MPI_Type_contiguous(sizeof(X), MPI_BYTE, &MPI_X);  
8  MPI_Type_commit(&MPI_X);  
9  // ...  
10 MPI_Type_free(&MPI_X);
```

Data Blocks

- Stack together blocks of data with some stride
- Example, type describing column of a matrix:

```
1 | int n = 2;
2 | std::vector<int> A{0, 1, 2, 3};
3 |
4 | MPI_Datatype MPI_COL;
5 | MPI_Type_vector(n, 1, n, MPI_INT, &MPI_COL);
6 | MPI_Type_commit(&MPI_COL);
7 |
8 | if (rank == 0) {
9 |     MPI_Send(A.data(), 1, MPI_COL, 1, 11, MPI_COMM_WORLD);
10 | } else if (rank == 1) {
11 |     std::vector<int> b(n);
12 |     MPI_Status stat;
13 |     MPI_Recv(b.data(), n, MPI_INT, 0, 11, MPI_COMM_WORLD, &stat);
14 | }
15 |
16 | MPI_Type_free(&MPI_COL);
```


Parallel I/O

- After computing and communication, forgotten element
- What if your input data consists of 10-100 TB?
- If sequential I/O becomes serious performance bottleneck (Amdahl's law)

I/O – Nasty Business

- Lack of really good hardware/software solutions
- Hard to scale
- Can go wrong in many places: client code, OS level, metadata server, storage node, network

MPI I/O

- Provides middle-level abstraction for I/O on top of the FS
- Hides low-level details
- But still requires seeks and data views

Basic MPI I/O Write

```
1 // ex05.cpp
2 #include <vector>
3 #include <mpi.h>
4
5 int main(int argc, char* argv[]) {
6     MPI_Init(&argc, &argv);
7
8     int rank;
9     MPI_Comm_rank(MPI_COMM_WORLD, &rank);
10
11     MPI_File fh;
12     MPI_File_open(MPI_COMM_WORLD, argv[1], MPI_MODE_CREATE|MPI_MODE_WRONLY,
13                   MPI_INFO_NULL, &fh);
14
15     MPI_Status stat;
16     std::vector<int> buf(8, rank);
17
18     MPI_File_seek(fh, rank * buf.size() * sizeof(int), MPI_SEEK_SET);
19     MPI_File_write_all(fh, buf.data(), buf.size(), MPI_INT, &stat);
20
21     MPI_File_close(&fh);
22
23     return MPI_Finalize();
24 } // main
```

Basic MPI I/O Read

```
1 // ex06.cpp
2 #include <vector>
3 #include <mpi.h>
4
5 int main(int argc, char* argv[]) {
6     MPI_Init(&argc, &argv);
7
8     int rank, size;
9     MPI_Comm_rank(MPI_COMM_WORLD, &rank);
10    MPI_Comm_size(MPI_COMM_WORLD, &size);
11
12    MPI_File fh;
13    MPI_File_open(MPI_COMM_WORLD, argv[1], MPI_MODE_RDONLY, MPI_INFO_NULL, &fh);
14
15    MPI_Offset fs;
16    MPI_File_get_size(fh, &fs);
17
18    MPI_Status stat;
19
20    int bsz = fs / size;
21    std::vector<int> buf(bsz / sizeof(int));
22
23    MPI_File_seek(fh, rank * bsz, MPI_SEEK_SET);
24    MPI_File_read_all(fh, buf.data(), buf.size(), MPI_INT, &stat);
25
26    MPI_File_close(&fh);
27
28    return MPI_Finalize();
29 } // main
```

More Complex Reading

```
1 // ex07.cpp
2 #include <vector>
3 #include <mpi.h>
4
5 int main(int argc, char* argv[]) {
6     MPI_Init(&argc, &argv);
7
8     int rank, size;
9     MPI_Comm_rank(MPI_COMM_WORLD, &rank);
10    MPI_Comm_size(MPI_COMM_WORLD, &size);
11
12    MPI_Datatype MPI_COL;
13    MPI_Type_vector(8, 1, 8, MPI_INT, &MPI_COL);
14    MPI_Type_commit(&MPI_COL);
15
16    MPI_File fh;
17    MPI_File_open(MPI_COMM_WORLD, argv[1], MPI_MODE_RDONLY, MPI_INFO_NULL, &fh);
18
19    MPI_Status stat;
20    std::vector<int> buf(8);
21
22    MPI_File_set_view(fh, rank * sizeof(int), MPI_INT, MPI_COL, "native", MPI_INFO_NULL);
23    MPI_File_read_all(fh, buf.data(), 8, MPI_INT, &stat);
24
25    MPI_Type_free(&MPI_COL);
26    MPI_File_close(&fh);
27
28    return MPI_Finalize();
29 } // main
```

Reading “Plain” Text

- Plain text is horrible for parallel computing
- No fixed-size separation between records, etc.
- Possible approaches:
 - Preprocess to binary format (e.g. index + data)
 - Preprocessor can be in OpenMP
 - Or:
 - Split file evenly at bytes level, read in parallel
 - Parse on each processor, communicate truncated lines

For Fun

- Implement 2D matrix-vector product with parallel I/O and custom communicators.