Creating MPI Variable Types

Why create your own MPI variable type?

- When using MPI it is best to do as few communications as possible
 - Rather send all required data in a single large communication than in a number of smaller communications
- If the data is continuous and all of the same type (e.g. an array) this is easy to achieve
- Often you want to send discontinuous information and/or variables of a mixture of types (e.g. only some of the member variables in an object)
- We can achieve this by making our own MPI variable types to encompass a number of simpler MPI types (or even other types that we have created our self)

Demonstrate with an example

 Creating an MPI variable to send some of the member variables within an object

```
#include <mpi.h>
#include <iostream>
#include <locale>

using namespace std;
int id, p;

MPI_Datatype MPI_my_class;

class my_class
{
  public:
        int I1, I2;
        int var_not_to_send;
        double D1;
        char S1[50];
        static void buildMPIType();
};
```

```
void my class::buildMPIType()
       int block lengths[4];
       MPI Aint offsets[4];
       MPI Aint addresses[4], add start;
       MPI_Datatype typelist[4];
       my class temp;
       typelist[0] = MPI INT;
       block lengths[0] = 1;
       MPI Get address(&temp.I1, &addresses[0]);
       typelist[1] = MPI INT;
       block lengths[1] = 1;
       MPI_Get_address(&temp.I2, &addresses[1]);
       typelist[2] = MPI DOUBLE;
       block lengths[2] = 1;
       MPI Get address(&temp.D1, &addresses[2]);
       typelist[3] = MPI CHAR;
       block lengths[3] = 50;
       MPI Get address(&temp.S1, &addresses[3]);
       MPI Get address(&temp, &add start);
       for (int i = 0; i < 4; i++) offsets[i] = addresses[i] - add start;
       MPI Type create struct(4, block lengths, offsets, typelist, &MPI my class);
       MPI Type commit(&MPI my class);
```

Continuing the example

```
int main(int argc, char *argv[])
       MPI Init(&argc, &argv);
       MPI Comm rank(MPI COMM WORLD, &id);
       MPI_Comm_size(MPI_COMM_WORLD, &p);
       my class::buildMPIType();
       my class data;
       if (id == 0)
              data.l1 = 6;
              data.I2 = 3;
              data.D1 = 10.0;
              data.var_not_to_send = 25;
              strncpy(data.S1, "My test string", 50);
       MPI_Bcast(&data, 1, MPI_my_class, 0, MPI_COMM_WORLD);
       cout << "On process " << id << " I1=" << data.I1 << " I2=" << data.I2 << " D1= " << data.D1 << " S1=" <<
       data.S1 << ". The unsent variable is " << data.var_not_to_send << endl;
       MPI_Type_free(&MPI_my_class);
       MPI Finalize();
```

What are we trying to do?

- When we have sent data previously we have sent a pointer to that data
- We still want to send a pointer to indicate which data is to be sent, but as the data can be discontinuous and/or of different types, we need to create a list of the data to be sent together with the offset of that data's memory location from the pointer
 - This must be done in a generic fashion (i.e. the set of offsets must be the same for every object of that type) We can't do it for, for instance, pointers stored within an object as the offset for the data pointed to will be different for each object
 - Later we will look at creating temporary MPI variable types to get around problems like this

How do we do this?

- We create a temporary object of the type we are trying to make the MPI datatype for:
 - The offsets of the member variables will be the same for all objects of the same class
 - my_class temp;
- We then store information about the individual variables that will make up the MPI datatype
 - Its MPI datatype typelist[2] = MPI_DOUBLE;
 - The number of continuous variable of that type block_lengths[2] = 1;
- We then get the pointer to that variable (stored in an MPI friendly format)
 - MPI_Get_address(&temp.D1, &addresses[2]);
 - Note that this is not yet the offset that we actually need

How do we do this? (continued)

- Once we have gathered this data for all the member variables that we wish to send using this MPI datatype we can calculate the offsets in the memory location
 - Obtain the memory location of the beginning of the object -MPI_Get_address(&temp, &add_start);
 - We then subtract this from all the addresses to get their offsets –
 offsets[i] = addresses[i] add_start;
- Once we have all this information we can use it to create the structure for the MPI_datatype (MPI_my_class in this example)
 - MPI_Type_create_struct(4, block_lengths, offsets, typelist, &MPI_my_class);
- Once the structure for the datatype has been created it must be committed before it can be used in any communications
 - MPI_Type_commit(&MPI_my_class);

Sending data with our new type

- We can send them using this MPI datatype like we would any other MPI variable
 - E.g. MPI_Bcast(&data, 1, MPI_my_class, 0, MPI_COMM_WORLD);
 - ...or MPI_Send(&data, 1, MPI_my_class, i, tag_num, MPI_COMM_WORLD

- Types should be freed once they are no longer needed
 - MPI Type free(&MPI my class);
 - Not very important in this example, but not doing so when repeatedly using a temporary MPI datatype to send data will result in memory leaks and potential crashes

Creating temporary MPI datatypes

- Creating an object for a class is not dissimilar to having MPI_datatypes for more primitive variables and can then be used in the same way
- Sometimes, though, you want to send a set of unrelated information at it is still best to send it at the same time
 - To do this we can create an that MPI_datatype that includes all the separate variables
 - These variables may not have a fixed relationship to one another in memory (e.g. memory that is reallocated as it changes size – includes most stl containers)
 - Can't simple create the type once and forget about it

An example – Sending disparate data

• In this example assume that the class and associated functions from the previous example exist:

```
my class my data[10];
int value top, value bottom;
void Send Data()
       int block lengths[3];
       MPI Aint addresses[3];
       MPI Datatype typelist[3];
       MPI Datatype MPI Temp;
       typelist[0] = MPI my class;
       block lengths[0] = 10;
       MPI Get address(my data, &addresses[0]);
       typelist[1] = MPI INT;
       block lengths[1] = 1;
       MPI Get address(&value top, &addresses[1]);
       typelist[2] = MPI INT;
       block lengths[2] = 1;
       MPI Get address(&value bottom, &addresses[2]);
       MPI Type create struct(3, block lengths, addresses, typelist, &MPI Temp);
       MPI Type commit(&MPI Temp);
       MPI Bcast(MPI BOTTOM, 1, MPI Temp, 0, MPI COMM WORLD);
       MPI Type free(&MPI Temp);
```

```
int main(int argc, char *argv[])
       MPI Init(&argc, &argv);
       MPI Comm rank(MPI COMM WORLD, &id);
       MPI Comm size(MPI COMM WORLD, &p);
       my class::buildMPIType();
       if (id == 0)
              for (int i = 0: i < 10: i++)
                      my data[i].11 = 6+i*25;
                      my data[i].12 = 3-i*4;
                      mv data[i].D1 = 10.0+31.*i;
                      my data[i].var not to send = 25;
                      strncpy(my data[i].S1, "My test string", 50);
              value top = 16;
              value bottom = 5;
       Send Data();
       cout << "On process " << id << endl;
       for (int i = 0; i < 10; i++)
              cout << "\t" << i << ": I1=" << my_data[i].I1 << " I2=" << my_data[i].I2 << " D1=" << my_data[i].D1 << " S1=" <<
                      my data[i].S1 << ". The unsent variable is " << my data[i].var not to send << endl;
       cout << "\ttop value: " << value top << "\tbottom value: " << value bottom << endl;</pre>
       MPI Type free(&MPI my class);
       MPI Finalize();
```

Sending disparate data

- Similar to creating an MPI type for a class, but some differences
 - Don't subtract the location of the beginning of the object The data being sent is not related in terms of memory location
 - When sending the data the pointer to the data to be sent (or received) is now the bottom of memory - MPI_BOTTOM
- Our temporary data structure includes an MPI type that we have created ourself
- Note that in this example the variables for sending and receiving are the same – This need not be the case – All that is required is that the order of the data and the size of the data on the sending and receiving side be the same

Writing to File when using MPI and Post-Processing

Writing data from multiple nodes

- Data can be transferred back to a single node and written as a single file
 - Fine for small amounts of output data
 - Can become very expensive in communication terms if there is a lot of output data and/or data needs to be frequently outputted (e.g. outputting data at multiple timesteps)
- Alternative is to have each process write its own data to file
 - Still need to still be transferred to a single destination, but often to a file server rather than one of the other nodes (which would still need to write it somewhere)
 - Need to carry out post-processing to combine the data from the multiple files

Numbered File Names

- File names need to be numbered according to the processor that created it and (potentially) the output/iteration/timestep number
 - You don't want to overwrite data from other processes and trying to have different processes write to the same file can result in either errors or serious blocking issues
- Using a string stream to create the file name and open the file:

```
stringstream fname;
fstream f1;
fname << "output" << num << "_" << id << ".dat";
f1.open(fname.str().c_str(), ios_base::out);</pre>
```