

Optimising Communication Routines in Parallel X-Agents

CCLRC Software Engineering Group
L.S. Chin, C. Greenough, and D.J. Worth

January 16, 2007

Abstract:

This report describes the work done in analysing the communication routines within parallel x-agent models used for modelling biological systems. Based on the analysis, several changes are proposed to overcome possible deadlocks due to communication synchronisation, as well as to optimise performance and memory utilisation within the communication routines.

Keywords: X-Agents, Intelligent Agents, MPI, Message Passing

1 Background

Researchers at the University of Sheffield have developed a framework for building x-agent models of biological systems based on the X-Machine Agent Markup Language (XMML) [8, 9]. XMML files define the types of agents, their memory data and functions involved in a model. This file is parsed by the x-parser to create a C program for simulating the system.

The generated C program can be either a serial or a parallel program, with the parallel version using MPI [4] to pass messages between agents on different nodes.

A more detailed introduction to the topics of x-agents and computational systems biology is provided by [7], while notes on previous analysis and optimisation of the x-agent model is available in [1].

2 Overview of communication routines

MPI Communication routines appear in seven routines, and are used to perform three distinct operations:

1. Propagation of messages

- `propagate_messages()` – uses `MPI_Recv`, calls:
 - `sent_message_pack()` – uses `MPI_Send`
 - `receive_message_pack()` – uses `MPI_Recv`

2. Propagation of agents

- `propagate_agents()` – uses `MPI_Recv`, calls:
 - `sent_agent_pack()` – uses `MPI_Send`
 - `receive_agent_pack()` – uses `MPI_Recv`

3. Distribution of partitioning information

- `send_space_partitions()` – uses `MPI_Bcast()`

Propagation of agents and messages are both performed using similar methods, while the distribution of partitioning information involves only a simple broadcast call. As such, optimisation activities were focused on the propagation of messages, the results of which would be equally applicable to the propagation of agents.

3 Dissecting propagate_messages()

Figure 1 illustrates the steps involved in propagate_messages().

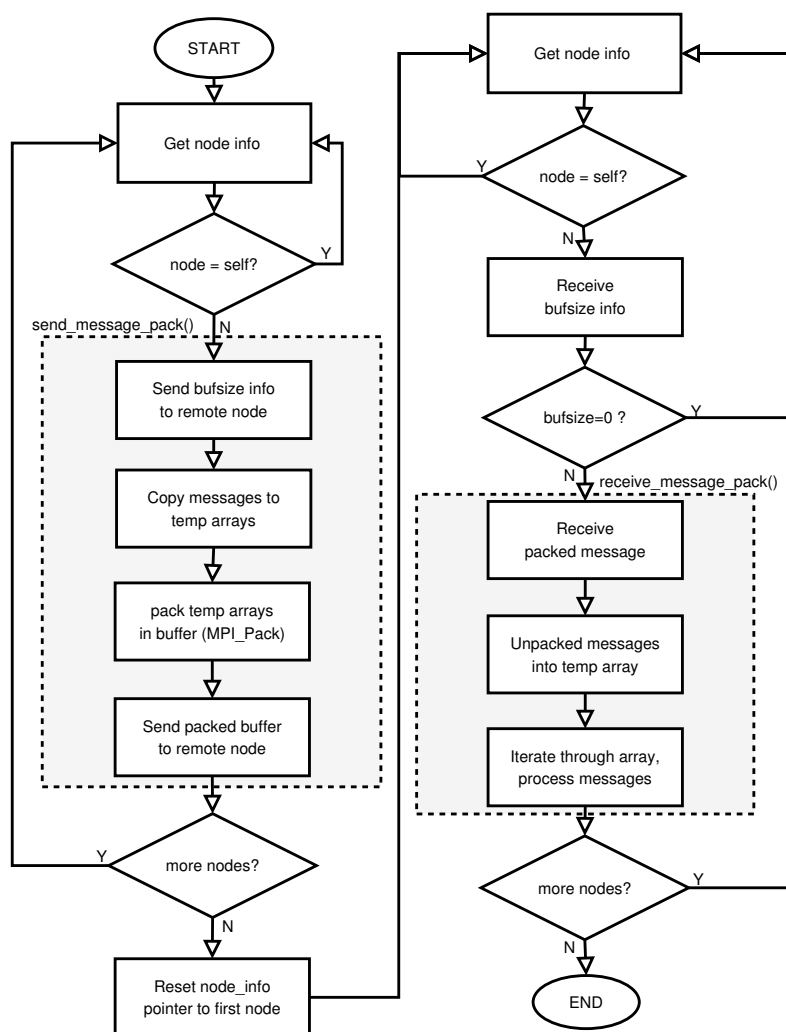


Figure 1: Flowchart for the propagate_messages() routine

The propagate_messages() routine, which calls send_message_pack() and receive_message_pack(), traverses the the list of nodes¹ and sends out the number of pending messages for each node. This is done by issuing an initial MPI_Send with information on the size of the pending messages. The actual messages are then packed into a buffer and sent using another MPI_Send call.

Once all messages have been sent out, the routine then traverses the list a second time to receive messages sent in by other nodes. For each remote node, it first receives information on the size of the incoming messages, allocates the appropriate buffers, and then issues matching MPI_Recv calls to complete the communication.

¹The term *nodes* is used in this text to represent instances of the model running in separate processes. In addition, we use *local node* to refer to the *node* performing the actions in question, and *remote node* to refer to the rest of the *nodes*

3.1 Possible deadlock

Note: Some of the terminology used in the following sections to describe MPI communication are defined in Appendix A (page 11)

In the current code, a series of `MPI_Send` calls are issued in advance followed by the `MPI_Recv` calls. Since `MPI_Send` is a blocking routine, the `MPI_Recv` calls will not be issued unless all the `MPI_Send` calls are completed. As such, the viability of this method essentially relies on `MPI_Send` being used in an asynchronous manner, whereby `MPI_Send` must complete even when a matching receive has not yet been issued on the remote node.

However, since asynchronous communication depends on message buffering which in turn is used at the discretion of the `MPI` implementation, it is unsafe to assume that communication will always be asynchronous. If the message being sent is larger than a predefined threshold, or when too many asynchronous calls deplete the system buffer, `MPI_Send` would resort to a synchronous communication mode. This results in a deadlock where every node would be waiting for the others to issue an `MPI_Recv`.

An example code which illustrates this behaviour is available in Appendix D (page 24). In this sample application, messages are repeatedly exchanged using `MPI_Send` followed by `MPI_Recv`. In each iteration, the message size gets bigger up to a point where `MPI_Send` becomes asynchronous and the application deadlocks.

3.2 Data duplication

An examination of the propagation routines reveal that data is duplicated up to four times during message staging.

Within each node, messages and agent memory are stored in *C structures* chained together as linked-lists. This means that data intended for propagation (in `propagate_messages()` or `propagate_agents()`) is fragmented and non-contiguous in memory. Since `MPI_Send` requires the outgoing data to be contiguous in memory, the packing of messages makes it possible to send a whole list of messages of different datatypes via the same call.

```
1 i = 0;
2 temp_message_location = node_info->location_messages;
3 while(temp_message_location)
4 {
5     location_message_list[i].id = temp_message_location->id;
6     location_message_list[i].x  = temp_message_location->x;
7     location_message_list[i].y  = temp_message_location->y;
8     location_message_list[i].z  = temp_message_location->z;
9     /* --- <snip> ---- */
10
11     i++;
12     temp_message_location = temp_message_location->next;
13 }
14
15 MPI_Pack(location_message_list, node_info->location_message_no, \
16         messagelocationType, buf, buffersize, &position, MPI_COMM_WORLD);
17
18 MPI_Send(buf, position, MPI_PACKED, node_info->node_id, 19, MPI_COMM_WORLD);
```

As seen above in the code snippet taken from the current implementation, data from the linked-list referenced to by `temp_message_location` is first copied to the `location_message_list` array. This ensures that data is contiguous in memory and ready for packing.

Data from `location_message_list` is then packed into `buf`, and eventually buffered by the `MPI_Send` call into the system buffer².

²If system buffer is unavailable, the `MPI_Send` call would block leading to a possible deadlock (as described in section 3.1)

Figure 2 depicts the data duplication that would occur when multiple message types are packed for sending.

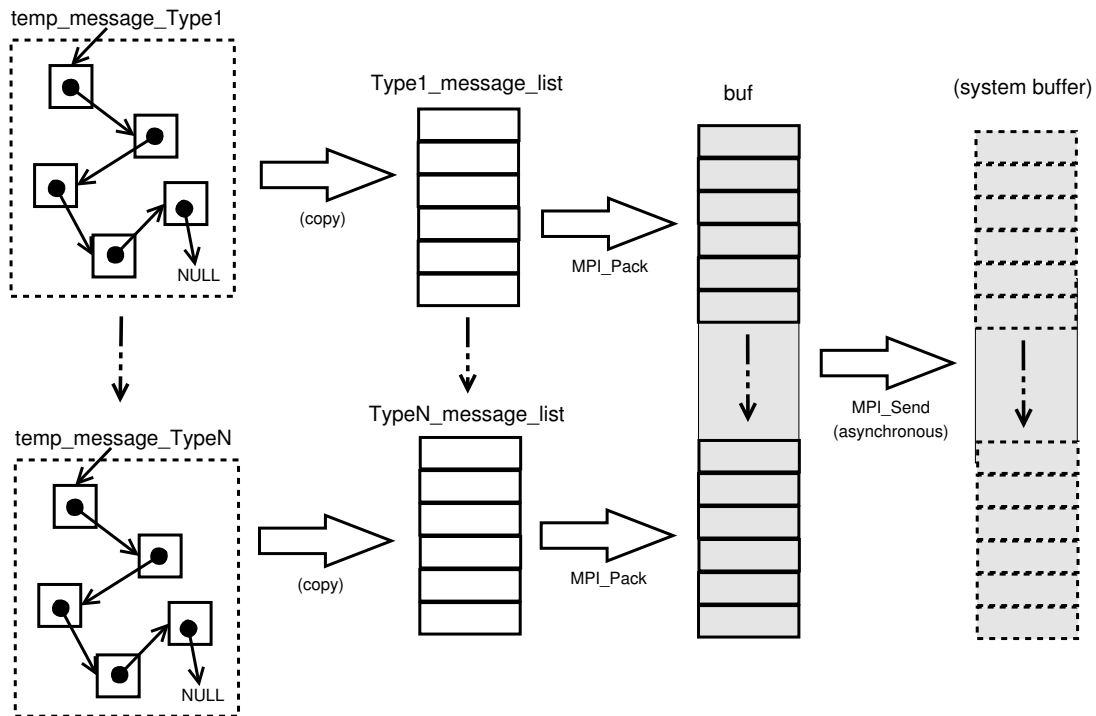


Figure 2: Memory utilisation of `send_message_pack()` for multiple messages types

Similar steps are performed (in reverse order) in `receive_message_pack()`. Data is unpacked from the receive buffer into a temporary array, which is then duplicated and attached to an appropriate linked-list.

When large amounts of messages are involved, this repeated duplication of memory is likely to have an impact on performance.

4 Proposed optimisation

This chapter presents several optimisation techniques than may overcome issues discussed in previous sections.

A sample implementation which incorporates these ideas are included in Appendix B and Appendix C. This code has not yet been thoroughly tested or benchmarked, and as such does not represent an end-product, but rather a reference point for further analysis and optimisation.

4.1 Using non-blocking communication

Non-blocking routines should be used for communication in order to obtain the following benefits:

- **Reduce synchronisation overheads:** By decoupling the completion of each send from the receipt, overhead imposed by the synchronisation between communicating nodes³ can be reduced.
- **Avoid deadlocks:** Non-blocking routine calls return immediately even if the message cannot yet be sent. This prevents the deadlock condition described in section 3.1 from occurring.

³Asynchronous communication can also be achieved using buffer-oriented sends, e.g. `MPI_Bsend`, and under certain circumstances, `MPI_Send`

- **Latency hiding:** After a non-blocking routine call returns, actual communication would continue in the background while the calling process can proceed with other tasks, returning at a later point to confirm that communication has completed successfully. This technique, often known as latency hiding, masks the overhead of waiting for communication completion by overlapping the idle time with useful computation.

In the sample implementation, `propagate_messages()` is split up into two routines – `propagate_messages_init()` which packs outgoing messages and issues the non-blocking communication calls, and `propagate_messages_complete()` which processes the input messages as well as confirms the completion of all communication.

Splitting the routine as so would simplify the overlapping of communication with other routines. Routines that can be scheduled during communication would include file I/O routines, run-time analysis routines for determining load balancing strategies, or even stages of the *X-machine* that are independent of the messages being communicated.

```

1 /* initiate propagation of messages */
2 propagate_messages_init();
3
4 /* write iteration data to file */
5 if (iteration_loop%output_frequency == 0)
6 {
7     saveiterationdata(iteration_loop);
8 }
9
10 /* other routines that do not depend on the messages */
11 analyse_load_balance(); /* fictional routine */
12
13 /* complete the propagation of messages */
14 propagate_messages_complete();

```

One drawback of using non-blocking communication is the increased memory requirements. Between the time when communication is initiated and when it is completed, application buffers that are allocated for every send and receive has to be preserved.

4.2 Message buffering strategies

To reduce the overhead incurred during data packing (covered in section 3.2), the following strategies were used:

- Packing the data *by hand* rather than using `MPI_Pack` – this involves allocating a single buffer that acts as both the temporary array as well as the send buffer. Addresses of strategic points in the buffer are casted to the different message types, which can then be used as the temporary array for copying data from the linked-list.
- Use of non-blocking *synchronous MPI* sends to prevent system buffering.
- Freeing up memory used by the linked-list as soon as possible.

Figure 3 illustrates the memory usage of this buffering strategy when applied to `propagate_messages()`. Similar techniques can be applied to the packing of different types of agents in `propagate_agents()`.

The following code snippet is an example of how the different message arrays can be assigned to the appropriate portions in `sendbuf`.

```

1 /* declaration of variables used in the following code snippet */
2 int bufsize;
3 void * sendbuf;
4 struct_message_type1 * msg_type1_array;

```

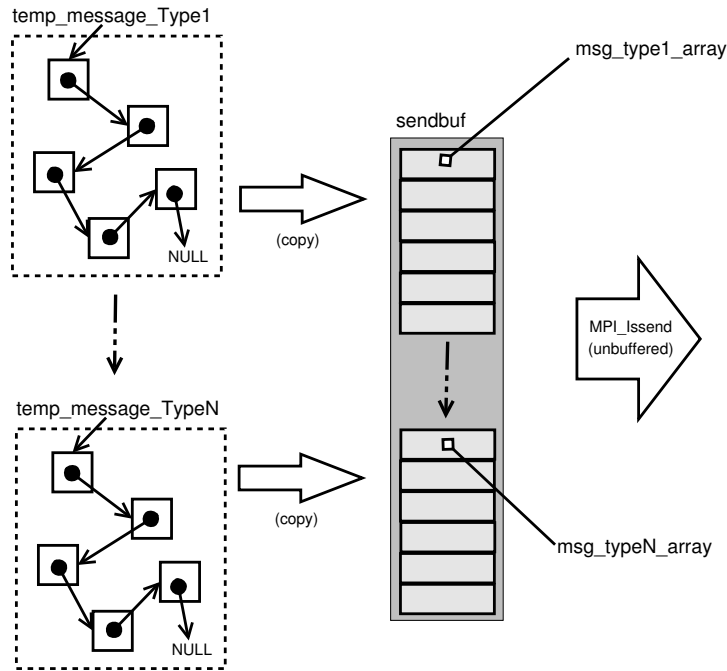


Figure 3: Memory usage of proposed optimisation for multiple message types

```

5 struct_message_type2 * msg_type2_array;
6 /* ... for subsequent message types ....
7 struct_message_typeN * msg_typeN_array;
8 */
9 int msg_type1_count, msg_type1_count /* ..., msg_typeN_count*/;
10
11 /* ---- <snip> ----*/
12
13 /* Determine required buffer size */
14 bufsize = 0;
15 bufsize += sizeof(struct_message_type1) * msg_type1_count;
16 bufsize += sizeof(struct_message_type2) * msg_type2_count;
17 /* .... for subsequent message types ....
18 bufsize += sizeof(struct message_typeN) * msg_typeN_count;
19 */
20
21 /* allocate memory for send buffer */
22 sendbuf = (void *) malloc(bufsize);
23
24 /* assign message array pointers */
25 msg_type1_array = (struct_message_type1 *) sendbuf;
26 msg_type2_array = (struct_message_type2 *) \
27     &msg_type1_array[msg_type1_count];
28 /* .... for subsequent message types ....
29 msg_typeN_array = (struct_message_typeN *) \
30     &msg_type(N-1)_array[msg_type(N-1)_count];
31 */
32
33 /* ---- < copy data from linked list into arrays > ---- */
34 /* ---- < free memory of linked list > ---- */
35 /* ---- < send data in 'sendbuf' to recipient > ---- */
36
37 free(sendbuf);

```

4.3 Reducing communication cost

In the current implementation, every node would issue an `MPI_Send` to all other nodes with information on how much data to expect in the actual propagation of messages. This means that at each propagation step, $N(N - 1)$ sets of communication has to be established followed by $\sum M_{ij}$ more for sending the actual messages (where N is the number of nodes, and $\sum M_{ij}$ the sum of all outgoing messages for every node). This would be a huge performance bottleneck, especially when the simulation is scaled up to large numbers of nodes.

Nearest Neighbour communication

This overhead can be greatly reduced if the communication pattern can be predetermined. For example, if the simulation domain is spatially decomposed and communication is limited to only the nearest neighbours, communication cost can be reduced from that of $O(n^2)$ to $O(n)$.

Each node would have to keep track of nodes that are adjacent to it, and limit its communication to only these nodes. Neighbouring nodes can be determined manually during the partitioning stage, or obtained automatically if the *Virtual Topology* mechanism in *MPI* were used (see Chapter 7 of [6]).

While we believe that the biological model being analysed can greatly benefit from this optimisation, the required changes has not been incorporated into the sample code provided (Appendix B and Appendix C) as it would involve modification of other routines which are not directly related to communication.

For other simulations where the communication pattern is unknown and the above-mentioned optimisation not applicable, we have considered several other strategies for making the propagation more efficient.

Strategy A: Using `MPI_Alltoall`

The first strategy is to propagate information on message sizes using `MPI_Alltoall`. This would allow us to replace the $(N - 1)$ `MPI_Sends` and $(N - 1)$ `MPI_Recvs` on each node with a single collective operation. Collective routines are generally more efficient compared to functionally equivalent sets of point-to-point routines as the communication algorithms used are often optimised based on message sizes and characteristics of the hardware.

A possible implementation of this strategy might involve the following code:

```
1 /* Constants generated by parser depending on number of message types */
2 #define MESSAGE_TYPE1 0
3 #define MESSAGE_TYPE2 1
4 #define MESSAGE_TYPE_COUNT 2
5 /* --- <snip> --- */
6
7 /* for in/out message count table, we need 2d arrays that are
8  contiguous in memory (for MPI_Alltoall). First dim needs to be
9  dynamically sized based on node count, while second dim is
10  pre-defined by the model and can be static.
11 */
12 int (*in_mcount)[MESSAGE_TYPE_COUNT];
13 int (*out_mcount)[MESSAGE_TYPE_COUNT];
14
15 /* Allocate required memory for in/out message count table */
16 in_mcount = (int (*)[MESSAGE_TYPE_COUNT]) malloc\
17     (sizeof(int) * totalnodes * MESSAGE_TYPE_COUNT);
18 out_mcount = (int (*)[MESSAGE_TYPE_COUNT]) malloc\
19     (sizeof(int) * totalnodes * MESSAGE_TYPE_COUNT);
20
21 /* How many messages do we have to send for each node? */
22 node_info = *p_node_info;
23 while(node_info)
24 {
```

```

25     out_mcount[node_info->node_id][MESSAGE_TYPE1] = \
26         node_info->location_type1_no;
27     out_mcount[node_info->node_id][MESSAGE_TYPE2] = \
28         node_info->location_type2_no;
29     node_info = node_info->next;
30 }
31
32 /* Propagate list */
33 MPI_Alltoall(out_mcount, MESSAGE_TYPE_COUNT, MPI_INT, in_mcount, \
34             MESSAGE_TYPE_COUNT, MPI_INT, MPI_COMM_WORLD);
35
36 /* we would now be able to obtain expected message count from
37    each node via in_mcount[ node_id ][ MESSAGE_TYPE? ]
38 */
39
40 /* ---- <snip> ---- */
41
42 free(in_count);
43 free(out_count);

```

Strategy B: Using MPI.Probe

Another possible strategy is to pack message count information and actual message data within the same communication. This would remove the need for establishing a second connection, and thus reduce the total communication count from $[N(N-1) + \sum M_{ij}]$ to only $[N(N-1)]$ (where N is the number of nodes, and $\sum M_{ij}$ the sum of all outgoing messages for every node).

Since $\sum M_{ij} \leq N(N-1)$, this strategy could potentially reduce the number of required communication by up to 50%.

On each node, messages for each remote node are packed into the send buffers along with a counter of messages for each message type. If there are no messages for a particular remote node, an empty (zero bytes) data is sent.

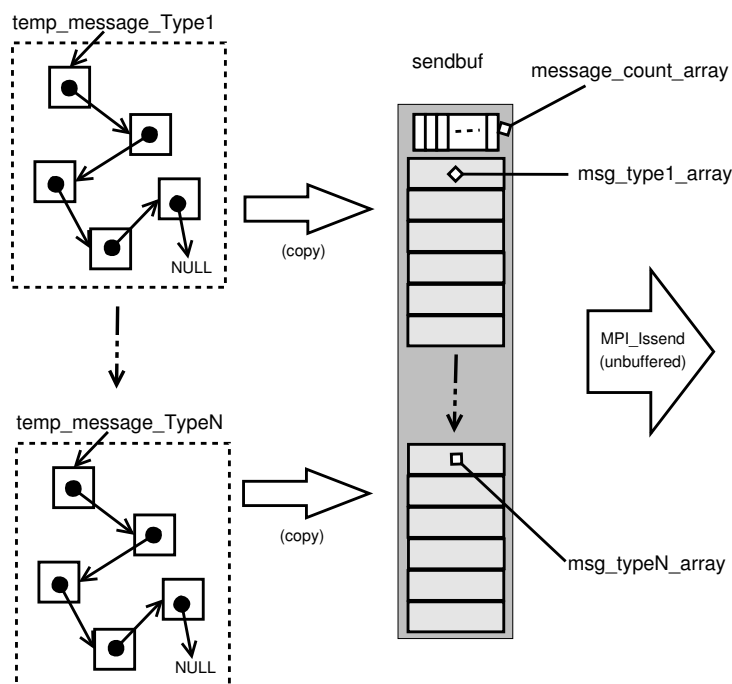


Figure 4: Counter of each message type is included at the beginning of the buffer. This assists the recipient in unpacking the rest of the buffer.

After issuing all non-blocking sends, each node then proceeds to receive data from every remote node. `MPI_Probe` is used to determine the size of data to expect from each communication. If the input data size is non-zero, an appropriate amount of memory is allocated for the receive buffer. All communication can then be completed using non-blocking receives.

Based on our preliminary benchmarks⁴, this strategy led to slightly better performance compared to *Strategy A*, and has therefore been implemented in our sample implementation (Appendix B and Appendix C).

⁴performed on SCARF (<http://hpcsg.esc.rl.ac.uk/scarf/index.html>), a 128 node dual Opteron cluster with 8GB memory per node, using *MPICH-GM* (MPI over Myrinet).

References

- [1] D.J. Worth and C. Greenough, “*Optimising X-Agent Models in Computational Biology*”, Software Engineering Group Note SEG-N-001 (2006).
http://www.cse.clrc.ac.uk/seg/html/wrap_pubs.shtml
- [2] M. Snir, S. Otto, S. Huss-Lederman, D. Walker, and J. Dongarra, “*MPI – The Complete Reference*”
<http://www.netlib.org/utk/papers/mpi-book/mpi-book.html>
- [3] “*MPI Performance Topics*”
http://www.llnl.gov/computing/tutorials/mpi_performance/
- [4] “*The Message Passing Interface (MPI) standard*”
<http://www-unix.mcs.anl.gov/mpi/>
- [5] Supercomputing Institute, University of Minnesota, “*Point to Point Communications with MPI*”
<http://www.msi.umn.edu/tutorial/scicomp/general/MPI/content1.html>
- [6] N. MacDonald, E. Minty, T. Harding, and S. Brown, “*Writing Message-Passing Parallel Programs with MPI*”
<http://www2.epcc.ed.ac.uk/epic/mpi/notes/>
- [7] “*Agent-based modelling – a tool for replicating biological systems*”
http://www.dcs.shef.ac.uk/~rod/Integrative_Systems_Biology.html
- [8] “*X-Machine Agents*”
<http://www.dcs.shef.ac.uk/~stc/x-agents/index.shtml>
- [9] “*FLAME – The Flexible Agent Modelling Environment*”
<http://www.flame.ac.uk>

Appendix A MPI Terminology

Application Buffer :

Address space allocated by user program and used to store data that is to be sent or received.

System Buffer :

System address space used for buffering messages. This buffer is not visible to the programmer, and accessible only by the *MPI* library. The primary purpose of system buffer space is to enable asynchronous communications.

Asynchronous Communication :

In an asynchronous communication, send operations may complete even when the receiving process has not received the message. Messages are copied to the system buffer, and scheduled for transmission when a matching receive is issued.

After the completion of an asynchronous send, the application send buffer can be safely overwritten, but there is no guarantee that the receiving process has received the message.

`MPI_Send` can be used to attempt an asynchronous send. However, depending on message sizes and the availability of system buffer, `MPI_Send` may fall back to a synchronous mode.

Synchronous Communication :

In a synchronous communication, send operations will only be completed when the message has been completely transferred to the receiving process.

After the completion of a synchronous send, it is guaranteed that the application send buffer can be safely overwritten, and that the receiving process has a copy of the message.

`MPI_Ssend` (note the extra 's' which stands for *synchronous*) can be used to issue a synchronous send.

Blocking :

A communication routine is blocking if its call only returns when the communication is completed. For example, `MPI_Send` will only return when the message has been moved to the system buffer or when the receiving process has received the message, while `MPI_Recv` will only return when incoming data is safely stored in the application buffer so that it is ready for use.

Non-Blocking :

A communication routine is non-blocking if its call returns without waiting for the communication to complete. It is not safe to modify or use the application buffer after a non-blocking call returns, and it is the programmer's responsibility to ensure that the communication has completed before the application buffer is free for reuse.

Non-blocking routines are primarily used to overlap computation with communication. Every *MPI* send/receive routine has a non-blocking counterpart, e.g. `MPI_Isend` and `MPI_Irecv`.

Appendix B Source Code for modified propagate_messages()

```
1 /* ===== *
2 * File: propagate_messages.c
3 * Desc: Modified version of propagate_messages(). Uses non-blocking
4 *       communication to avoid deadlocks, and customised buffering
5 *       strategies to reduce memory requirements.
6 *       Routine is split into propagate_messages_init() and
7 *       propagate_messages_complete() so that non-dependent computation
8 *       may be overlapped with the communication.
9 *
10 * Auth: L.S. Chin (l.s.chin@rl.ac.uk)
11 *
12 * NOTES: Annotations on 'dynamic' (model dependent) sections of this code
13 *        are included in annotated_propagate_messages.c
14 * ===== */
15
16 #include "header.h"
17
18 /* ===== Declaration of required constants ===== */
19
20 /* These values would be set in header.h by x-parser, depending on number of
21    message types */
22 #define MESSAGE_LOCATION 0
23 #define MESSAGE_TYPE_COUNT 1
24
25
26 /* message tag used for propagating messages */
27 #define TAG 12
28
29
30 /* ===== Declaration of file level global vars ===== */
31
32 /* buffer space for packed message from/to each node */
33 void ** inbuf;
34 void ** outbuf;
35
36 /* vars associated to MPI communication */
37 MPI_Request *in_req;
38 MPI_Request *out_req;
39
40
41 /* ===== Begin implementation code ===== */
42
43 /** \fn void propagate_messages_init()
44  * \brief Initiate propagation of messages between nodes using
45  *       non-blocking communication routines.
46  */
47 void propagate_messages_init() {
48     /* <NOTE> This routine uses the following external global variables
49      * + xmachine_message_location ** p_location_message (Pointer to
50      *       message struct for looping through location message list)
51      * + node_information ** p_node_info (Pointer to head of node list)
52      * + int totalnodes (Result of MPI_Comm_size)
53      * + int node_number (Result of MPI_Comm_rank)
54      */
55
56     /* ----- Begin Variable declarations ----- */
57     int i, j, outcount, bufsize;
58     int probed;
59     int *message_count_list;
60     MPI_Status status;
61
62     node_information * node_info;
```

```

63  /* pointers to temporary messages. One set for each message type */
64  xmachine_message_location *message_location_temp;
65  xmachine_message_location_data *message_location_list;
66
67
68  /* ----- Allocate Required Memory ----- */
69
70  /* Memory for MPI send/receive requests */
71  in_req = (MPI_Request *)malloc(sizeof(MPI_Request) * totalnodes);
72  out_req = (MPI_Request *)malloc(sizeof(MPI_Request) * totalnodes);
73
74  /* Memory for buffer arrays */
75  inbuf = (void **)malloc(sizeof(void *) * totalnodes);
76  outbuf = (void **)malloc(sizeof(void *) * totalnodes);
77
78
79  /* ----- Fill send buffers and post non-blocking sends ---- */
80
81  /* iterate through node list */
82  node_info = *p_node_info;
83  while(node_info)
84  {
85      i = node_info->node_id;
86      /* don't send messages to self */
87      if (i == node_number)
88      {
89          outbuf[i] = NULL;
90          out_req[i] = MPI_REQUEST_NULL;
91          node_info = node_info->next;
92          continue;
93      }
94
95      bufsize = 0;
96      outcount = 0;
97
98      /* for each message type */
99      outcount += node_info->location_message_no;
100
101
102      /* build output buffer */
103      if (outcount > 0)
104      {
105          /* mem requirements for message_count_list array */
106          bufsize += sizeof(int) * MESSAGE_TYPE_COUNT;
107
108          /* mem requirements for each message type messages */
109          bufsize += sizeof(xmachine_message_location_data) * \
110                  node_info->location_message_no;
111
112          /* allocate required memory */
113          outbuf[i] = (void *) malloc(bufsize);
114
115          /* assign array pointers to relevant points in buffer */
116          message_count_list = (int *) outbuf[i];
117
118          message_count_list[MESSAGE_LOCATION] = node_info->location_message_no;
119          message_location_list = (xmachine_message_location_data *)\
120                                  &message_count_list[MESSAGE_TYPE_COUNT];
121
122
123          /* traverse message list and populate buffer */
124          message_location_temp = node_info->location_messages;
125          j = 0;
126          while (message_location_temp)

```

```

127     {
128         message_location_list[j].id = message_location_temp->id;
129         message_location_list[j].x = message_location_temp->x;
130         message_location_list[j].y = message_location_temp->y;
131
132         message_location_temp = message_location_temp->next;
133         j++;
134     }
135     /* clear location messages from internal list */
136     p_location_message = &node_info->location_messages;
137     freelocationmessages();
138     node_info->location_message_no = 0;
139     p_location_message = &current_node->location_messages;
140
141
142 }
143 else /* nothing to send */
144 {
145     outbuf[i] = NULL;
146 }
147
148 /* post non-blocking send */
149 MPI_Issend(outbuf[i], bufsize, MPI_BYTE, i, TAG, \
150           MPI_COMM_WORLD, &out_req[i]);
151
152 /* select next node */
153 node_info = node_info->next;
154 }
155
156
157
158 /* ----- Prepare and post non-blocking receives ----- */
159
160 /* we expect (totalnodes - 1) incoming messages */
161 probed = 0;
162 inbuf[node_number] = NULL;
163 in_req[node_number] = MPI_REQUEST_NULL;
164 while (probed < (totalnodes - 1))
165 {
166     /* probe incoming messages for bufsize */
167     MPI_Probe(MPI_ANY_SOURCE, TAG, MPI_COMM_WORLD, &status);
168
169     /* get sender's id */
170     i = status.MPI_SOURCE;
171
172     /* get bufsize */
173     MPI_Get_count(&status, MPI_BYTE, &bufsize);
174
175     /* allocate memory for recv buffer */
176     if (bufsize > 0)
177     {
178         inbuf[i] = (void *)malloc(bufsize);
179     }
180     else
181     {
182         inbuf[i] = NULL;
183     }
184
185     /* post non-blocking receive */
186     MPI_Irecv(inbuf[i], bufsize, MPI_BYTE, i, TAG, \
187             MPI_COMM_WORLD, &in_req[i]);
188
189     probed ++;
190 }

```

```

191
192 }
193
194
195 /** \fn void propagate_messages_complete()
196  * \brief Complete propagation of messages between nodes
197  */
198 void propagate_messages_complete() {
199
200     /* <NOTE> This routine uses the following external global variables
201      * + node_information ** p_node_info (Pointer to head of node list)
202      * + int totalnodes (Result of MPI_Comm_size)
203      */
204
205     /* ----- Begin Variable declarations ----- */
206     int i, j, bufsize, sender;
207     int *message_count_list;
208     MPI_Status status;
209     node_information * node_info;
210
211     /* pointers to temporary messages. One set for each message type */
212     xmachine_message_location *message_location_temp;
213     xmachine_message_location_data *message_location_list;
214
215
216
217     /* ----- Complete and process pending receives ----- */
218
219     /* process received messages */
220     /* instead of a Waitall, use Waitany so that we may overlap data
221      handling time with other pending receives */
222     for (j = 0; j < totalnodes - 1; j++)
223     {
224         /* wait for any receive to complete */
225         MPI_Waitany(totalnodes, in_req, &sender, &status);
226
227         /* determine size of message */
228         MPI_Get_count(&status, MPI_BYTE, &bufsize);
229
230         /* if size = 0, there's nothing for us to do */
231         if (bufsize == 0) continue;
232
233         /* Assign array pointers to appropriate locations in buffer */
234         message_count_list = (int *)inbuf[sender];
235
236         /* for each message type */
237         message_location_list = (xmachine_message_location_data *)\
238                                 &message_count_list[MESSAGE_TYPE_COUNT];
239
240
241         /* process received messages */
242         for (i = 0; i < message_count_list[MESSAGE_LOCATION]; i++)
243         {
244             message_location_temp = (xmachine_message_location*) \
245                                     add_location_message_internal();
246
247             message_location_temp->id = message_location_list[i].id;
248             message_location_temp->x = message_location_list[i].x;
249             message_location_temp->y = message_location_list[i].y;
250         }
251
252         free(inbuf[sender]);
253     }
254

```

```
255     free(inbuf);
256     free(in_req);
257
258
259
260     /* ----- Complete non-blocking sends ----- */
261
262     /* wait for all non-blocking sends to complete */
263     MPI_Waitall(totalnodes, out_req, MPI_STATUSES_IGNORE);
264
265     /* free allocated send buffers */
266     for(i = 0; i < totalnodes; i++)
267     {
268         if (outbuf[i]) free(outbuf[i]);
269     }
270     free(outbuf);
271     free(out_req);
272
273 }
```


Appendix C Source Code for modified propagate_agents()

```
1 /* ===== *
2 * File: propagate_agents.c
3 * Desc: Modified version of propagate_agents(). Uses non-blocking
4 *       communication to avoid deadlocks, and customised buffering
5 *       strategies to reduce memory requirements.
6 *
7 *       While not implemented here, routine can be easily split up into
8 *       propagate_agents_init() and propagate_agents_complete() to allow
9 *       for overlapping of computation with agent propagation.
10 *
11 *       Algorithm for moving agents from main list to outbound list has
12 *       also been modified (simplified?)
13 *
14 * Auth: L.S. Chin (l.s.chin@rl.ac.uk)
15 *
16 * NOTES: Techniques for generating this code dynamically with the x-parser
17 *        can be implied from the annotations in annotated_propagate_messages.c
18 * ===== */
19
20 #include "header.h"
21
22 /* ===== Declaration of required constants ===== */
23
24 /* These values would be set in header.h by x-parser, depending on number of
25    agent types */
26 #define AGENT_CIRCLE 0
27 #define AGENT_CIRCLE2 1
28 #define AGENT_TYPE_COUNT 2
29
30 /* message tag used for propagating agents */
31 #define TAG 11
32
33
34
35 /* ===== Begin implementation code ===== */
36
37 /** \fn void propagate_agents()
38  * \brief Check agent positions to see if any need to be moved to a another node.
39  */
40 void propagate_agents() {
41
42     /* ----- Begin Variable declarations ----- */
43
44     int agent_type;
45     int remove_agent;
46     int sender, outcount;
47     int *agent_count_list;
48     int probed, i, j, bufsize;
49     node_information *node_info;
50     xmachine *prev_xmachine, *next_xmachine;
51     double x_xmachine, y_xmachine, z_xmachine;
52
53
54     /* variables used for MPI calls */
55     MPI_Status status;
56     MPI_Request * in_req;
57     MPI_Request * out_req;
58
59     /* buffer space for packed agents from/to each PE */
60     void ** inbuf;
61     void ** outbuf;
62
```

```

63  /* pointers to temporary agents */
64  xmachine_memory_Circle *agent_Circle_list;
65  xmachine_memory_Circle2 *agent_Circle2_list;
66
67
68
69  /* ---- determine which agents should be moved ----- */
70
71  /* loop through agent list */
72  prev_xmachine = NULL;
73  remove_agent = 0;
74  current_xmachine = current_node->agents;
75  while(current_xmachine)
76  {
77      next_xmachine = current_xmachine->next;
78
79      if(current_xmachine->xmachine_Circle != NULL)
80      {
81          x_xmachine = current_xmachine->xmachine_Circle->x;
82          y_xmachine = current_xmachine->xmachine_Circle->y;
83          z_xmachine = 0.0;
84          agent_type = 0;
85      }
86      else if(current_xmachine->xmachine_Circle2 != NULL)
87      {
88          x_xmachine = current_xmachine->xmachine_Circle2->x;
89          y_xmachine = current_xmachine->xmachine_Circle2->y;
90          z_xmachine = 0.0;
91          agent_type = 1;
92      }
93
94  /* if any agent is located beyond our partition data, move it */
95  if (
96      x_xmachine < current_node->partition_data[0] ||
97      x_xmachine > current_node->partition_data[1] ||
98      y_xmachine < current_node->partition_data[2] ||
99      y_xmachine > current_node->partition_data[3] ||
100     z_xmachine < current_node->partition_data[4] ||
101     z_xmachine > current_node->partition_data[5] )
102  {
103      /* determine which node we should move agent to */
104
105      /* loop through node list */
106      node_info = *p_node_info;
107      while(node_info)
108      {
109          if (node_info->node_id == node_number)
110          { /* skip my own node */
111              node_info = node_info->next;
112              continue;
113          }
114
115          /* should agent be in this node? */
116          if (
117              x_xmachine > node_info->partition_data[0] &&
118              x_xmachine < node_info->partition_data[1] &&
119              y_xmachine > node_info->partition_data[2] &&
120              y_xmachine < node_info->partition_data[3] &&
121              z_xmachine > node_info->partition_data[4] &&
122              z_xmachine < node_info->partition_data[5] )
123          {
124              /* move agent to target node */
125              if (agent_type == 0)
126              {

```

```

127         node_info->Circle_agent_no ++;
128         current_xmachine->next = node_info->Circle_agents;
129         node_info->Circle_agents = current_xmachine;
130     }
131     if (agent_type == 1)
132     {
133         node_info->Circle2_agent_no ++;
134         current_xmachine->next = node_info->Circle2_agents;
135         node_info->Circle2_agents = current_xmachine;
136     }
137
138     /* flag agent for removal from current node */
139     remove_agent = 1;
140
141     /* no need to traverse the node list further */
142     node_info = NULL;
143 }
144 else
145 { /* continue searching */
146     node_info = node_info->next;
147 }
148 }
149
150 }
151
152 /* remove agent from current node if necessary */
153 if (remove_agent)
154 {
155     remove_agent = 0; /* reset value */
156
157     /* decrement agent count */
158     current_node->agent_total --;
159
160     /* if first agent is removed */
161     if (prev_xmachine == NULL)
162     {
163         current_node->agents = next_xmachine;
164     }
165     else
166     {
167         prev_xmachine->next = next_xmachine;
168     }
169 }
170 else
171 {
172     prev_xmachine = current_xmachine;
173 }
174
175 /* move on to next agent in list */
176 current_xmachine = next_xmachine;
177 }
178
179
180
181 /* ----- Allocate Required Memory ----- */
182
183 /* MPI related dynamic memory requirements */
184 in_req = (MPI_Request *)malloc(sizeof(MPI_Request) * totalnodes);
185 out_req = (MPI_Request *)malloc(sizeof(MPI_Request) * totalnodes);
186
187 /* Memory for buffer arrays */
188 inbuf = (void **)malloc(sizeof(void *) * totalnodes);
189 outbuf = (void **)malloc(sizeof(void *) * totalnodes);
190

```

```

191
192 /* ----- fill send buffers and post non-blocking sends ----- */
193 node_info = *p_node_info;
194 while(node_info)
195 {
196     i = node_info->node_id;
197     if (i == node_number)
198     {
199         outbuf[i] = NULL;
200         out_req[i] = MPI_REQUEST_NULL;
201         node_info = node_info->next;
202         continue;
203     }
204
205     bufsize = 0;
206     outcount = 0;
207
208     /* for each agent type */
209     outcount += node_info->Circle_agent_no;
210     outcount += node_info->Circle2_agent_no;
211
212     /* build output buffer */
213     if (outcount > 0)
214     {
215         /* mem requirements for agent_count_list array */
216         bufsize += sizeof(int) * AGENT_TYPE_COUNT;
217
218         /* for each message type */
219         /* mem requirements for Circle agents */
220         bufsize += sizeof(xmachine_memory_Circle) * \
221                 node_info->Circle_agent_no;
222         /* mem requirements for Circle2 agents */
223         bufsize += sizeof(xmachine_memory_Circle2) * \
224                 node_info->Circle2_agent_no;
225
226         /* allocate required memory */
227         outbuf[i] = (void *) malloc(bufsize);
228
229         /* assign array pointers to relevant points in buffer */
230         agent_count_list = (int *) outbuf[i];
231
232         agent_count_list[AGENT_CIRCLE] = node_info->Circle_agent_no;
233         agent_Circle_list = (xmachine_memory_Circle *) \
234                 &agent_count_list[AGENT_TYPE_COUNT];
235
236         agent_count_list[AGENT_CIRCLE2] = node_info->Circle2_agent_no;
237         agent_Circle2_list = (xmachine_memory_Circle2 *) \
238                 &agent_Circle_list[agent_count_list[AGENT_CIRCLE]];
239
240
241         /* traverse agent list and populate buffer */
242         j = 0;
243         temp_xmachine = node_info->Circle_agents;
244         while (temp_xmachine)
245         {
246             agent_Circle_list[j].id = temp_xmachine->xmachine_Circle->id;
247             agent_Circle_list[j].x = temp_xmachine->xmachine_Circle->x;
248             agent_Circle_list[j].y = temp_xmachine->xmachine_Circle->y;
249             agent_Circle_list[j].fx = temp_xmachine->xmachine_Circle->fx;
250             agent_Circle_list[j].fy = temp_xmachine->xmachine_Circle->fy;
251             agent_Circle_list[j].radius = temp_xmachine->xmachine_Circle->radius;
252             agent_Circle_list[j].iradius = temp_xmachine->xmachine_Circle->iradius;
253
254             temp_xmachine = temp_xmachine->next;

```

```

255         j++;
256     }
257     /* free list early to conserve memory */
258     p_xmachine = &node_info->Circle_agents;
259     freexmachines();
260     node_info->Circle_agent_no = 0;
261
262
263     /* repeat for next agent type */
264     j = 0;
265     temp_xmachine = node_info->Circle2_agents;
266     while(temp_xmachine)
267     {
268         agent_Circle2_list[j].id = temp_xmachine->xmachine_Circle2->id;
269         agent_Circle2_list[j].x = temp_xmachine->xmachine_Circle2->x;
270         agent_Circle2_list[j].y = temp_xmachine->xmachine_Circle2->y;
271         agent_Circle2_list[j].fx = temp_xmachine->xmachine_Circle2->fx;
272         agent_Circle2_list[j].fy = temp_xmachine->xmachine_Circle2->fy;
273         agent_Circle2_list[j].radius = temp_xmachine->xmachine_Circle2->radius;
274         agent_Circle2_list[j].iradius = temp_xmachine->xmachine_Circle2->iradius;
275
276         temp_xmachine = temp_xmachine->next;
277         j++;
278     }
279     /* free list early to conserve memory */
280     p_xmachine = &node_info->Circle2_agents;
281     freexmachines();
282     node_info->Circle2_agent_no = 0;
283
284
285     /* Reset xmachine pointer */
286     p_xmachine = &current_node->agents;
287
288 }
289 else
290 {
291     outbuf[i] = NULL;
292 }
293
294 /* post non-blocking send */
295 MPI_Issend(outbuf[i], bufsize, MPI_BYTE, i, TAG, \
296           MPI_COMM_WORLD, &out_req[i]);
297
298 /* select next node */
299 node_info = node_info->next;
300 }
301
302
303
304 /* ----- Prepare and post non-blocking receives ----- */
305
306 /* we expect (totalnodes - 1) incoming messages */
307 probed = 0;
308 inbuf[node_number] = NULL;
309 in_req[node_number] = MPI_REQUEST_NULL;
310
311 while (probed < (totalnodes - 1))
312 {
313     /* probe incoming messages for bufsize */
314     MPI_Probe(MPI_ANY_SOURCE, TAG, MPI_COMM_WORLD, &status);
315
316     /* get sender's id */
317     i = status.MPI_SOURCE;
318

```

```

319     /* get bufsize */
320     MPI_Get_count(&status, MPI_BYTE, &bufsize);
321
322     if (bufsize > 0)
323     {
324         inbuf[i] = (void *)malloc(bufsize);
325     }
326     else
327     {
328         inbuf[i] = NULL;
329     }
330
331     /* post non-blocking receive */
332     MPI_Irecv(inbuf[i], bufsize, MPI_BYTE, i, TAG, \
333             MPI_COMM_WORLD, &in_req[i]);
334
335     probed ++;
336 }
337
338
339 /* NOTE: routine can be split up here if we want to separate
340 *       into _init() and _complete(). Some of the variables
341 *       would have to be declared as global within this file.
342 *       (as implemented in propagate_messages.c)
343 */
344
345
346 /* ----- Complete and process pending receives ----- */
347
348 /* process received messages */
349 /* instead of a Waitall, use Waitany so that we may overlap data
350    handling time with other pending receives */
351 for (j = 0; j < totalnodes - 1; j++)
352 {
353     /* wait for any receive to complete */
354     MPI_Waitany(totalnodes, in_req, &sender, &status);
355
356     /* determine size of message */
357     MPI_Get_count(&status, MPI_BYTE, &bufsize);
358
359     /* if size = 0, there's nothing for us to do */
360     if (bufsize == 0) continue;
361
362     /* Assign array pointers to appropriate locations in buffer */
363     agent_count_list = (int *) inbuf[sender];
364     /* for each message type */
365     agent_Circle_list = (xmachine_memory_Circle *) \
366         &agent_count_list[AGENT_TYPE_COUNT];
367     agent_Circle2_list = (xmachine_memory_Circle2 *) \
368         &agent_Circle_list[agent_count_list[AGENT_CIRCLE]];
369
370
371     /* process received messages */
372     for (i = 0; i < agent_count_list[AGENT_CIRCLE]; i++)
373     {
374         add_Circle_agent( \
375             agent_Circle_list[i].id, \
376             agent_Circle_list[i].x, \
377             agent_Circle_list[i].y, \
378             agent_Circle_list[i].fx, \
379             agent_Circle_list[i].fy, \
380             agent_Circle_list[i].radius, \
381             agent_Circle_list[i].iradius \
382         );

```

```

383
384     }
385
386     /* repeat for next agent type */
387     for (i = 0; i < agent_count_list[AGENT_CIRCLE2]; i++)
388     {
389         add_Circle2_agent( \
390             agent_Circle2_list[i].id, \
391             agent_Circle2_list[i].x, \
392             agent_Circle2_list[i].y, \
393             agent_Circle2_list[i].fx, \
394             agent_Circle2_list[i].fy, \
395             agent_Circle2_list[i].radius, \
396             agent_Circle2_list[i].iradius \
397         );
398     }
399
400     free(inbuf[j]);
401 }
402
403 free(inbuf);
404 free(in_req);
405
406
407
408
409 /* ----- Complete non-blocking sends ----- */
410
411 /* wait for all non-blocking sends to complete */
412 MPI_Waitall(totalnodes, out_req, MPI_STATUSES_IGNORE);
413
414 /* free allocated send buffers */
415 for(i = 0; i < totalnodes; i++)
416 {
417     if (outbuf[i]) free(outbuf[i]);
418 }
419 free(outbuf);
420 free(out_req);
421
422 }

```

Appendix D Source Code for test application – mpisend.c

```
1 #include <mpi.h>
2
3 #define BUFMIN 10
4 #define BUFMAX 3000
5 #define STEPSIZE 10
6
7 /* simple test to determine if and when MPI_Send becomes blocking */
8 /* Since send buffer implementation is not fixed by the MPI standards.
9    a program which runs successfully under one set of conditions may
10   fail under another set */
11
12 int main(int argc, char ** argv) {
13
14     /* dynamically sized buffer for send/receive */
15     int *outbuf, *inbuf;
16     int target;
17     int dsize, i, j;
18
19     /* MPI related vars */
20     MPI_Comm w_comm = MPI_COMM_WORLD;
21     MPI_Status status;
22     int w_size, w_rank;
23
24     MPI_Init(&argc, &argv);
25     MPI_Comm_size(w_comm, &w_size);
26     MPI_Comm_rank(w_comm, &w_rank);
27
28     if (w_size < 2)
29     {
30         printf("This test is meant to be run with at least 2 procs\n");
31         MPI_Finalize();
32         return 1;
33     }
34
35     /* prep */
36     dsize = sizeof(double);
37     target = (w_rank == 0) ? 1 : 0; /* select other proc as target */
38
39     if (w_rank == 0) printf("Testing message exchange using MPI_Send and \
40         MPI_Recv for data size %d to %d bytes (in steps of %d bytes) :\n" \
41         , BUFMIN * dsize, BUFMAX * dsize, STEPSIZE * dsize);
42
43     /* Have proc 0 and 1 repeatedly exchange messages of increasing size */
44     for (i = BUFMIN; i <= BUFMAX; i += STEPSIZE)
45     {
46         /* other procs do nothing except participate in Barriers. */
47         if (w_rank > 1)
48         {
49             MPI_Barrier(w_comm);
50             continue;
51         }
52
53         /* prepare dummy data */
54         outbuf = malloc(dsize * i);
55         inbuf = malloc(dsize * i);
56         for (j = 0; j < i; j++) inbuf[j] = (double)(w_rank + j);
57
58         if (w_rank == 0) printf("Exchanging %d bytes\n", i * dsize);
59
60         /* send message to target */
61         MPI_Send(outbuf, i, MPI_DOUBLE, target, 0, w_comm);
62
```



```
63
64     /* receive message from target */
65     MPI_Recv(outbuf, i, MPI_DOUBLE, target, 0, w_comm, &status);
66
67     if (w_rank == 0) printf("----- DONE -----\n");
68
69     /* Make sure exchange completed before proceeding to next round */
70     MPI_Barrier(w_comm);
71
72     /* release dummy data */
73     free(outbuf);
74     free(inbuf);
75
76 }
77
78 MPI_Finalize();
79
80 return 0;
81 }
```