

**Message Passing Java
(MPJ)**

**Report
And
Performance Results for Java and GridTcp**

**Prepared By
Zhiji Huang
(6/15/2005)**

Table of Contents

1.0 Introduction	1
2.0 Design	1
3.0 Algorithms	2
3.1 MPJ.Init()	2
3.2 MPJ.Init() GridTcp	3
3.3 Send()	3
3.4 Recv()	3
3.5 Pack()	4
3.6 Unpack()	4
3.7 Barrier()	4
3.8 Bcast()	4
3.9 Reduce()	4
3.10 Other communications	5
3.11 SetBufferSize()	8
4.0 PingPong Performance	9
5.0 Performance Analysis	12
5.1 Object Serialization	12
5.2 Primitive Serialization	12
6.0 Initialization and Use	13
6.1 mpjrun	13
7.0 Continuation of MPJ	14

1.0 Introduction

Message Passing Java (MPJ) is being developed as a middleware between the user program and various communications protocols. Currently, MPJ supports GridTcp sockets as well as Java sockets. The objective is to provide a set of communication functionalities supporting distributed computing. This report describes the design, implementation, and performance results throughout this project.

2.0 Design

MPJ design is as follows:

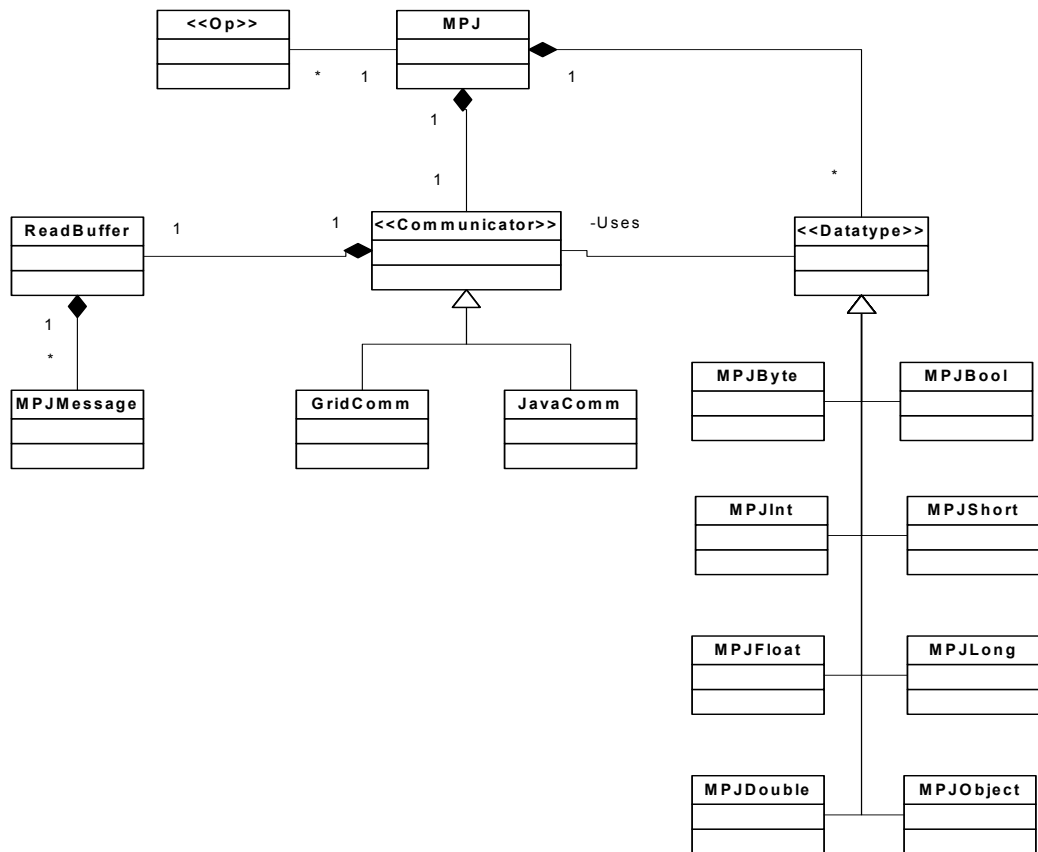


Diagram 1: Major MPJ classes currently in project.

The major classes shown in Diagram 1 (previous) are the implemented classes of MPJ.

MPJ is the main application. It contains a communicator of either JavaComm or GridComm, depending on the type of socket used. It also contains various data types supported by MPJ. In addition, MPJ provides initialization and finalization methods for the network connections.

JavaComm and GridComm hold Java sockets and GridTcp sockets, respectively. These two classes do not provide any major functionality other than maintaining the input and output streams of their respective sockets.

Communicator is the class that provides the primary communications capabilities of MPJ. As of the writing of this report, Communicator contains the major functions like Send(), Recv(), Barrier(), Bcast(), Pack() and Unpack().

MPJMessage is a wrapper around each message received by the Recv() functions. It holds the message's status and the actual message itself.

The various Datatype subclasses provide serialization and deserialization of their respective types for Communicator.

Various Operation classes are wrappers around operations that are implemented in the Datatype subclasses.

3.0 Algorithms

This section describes the important algorithms used in the major functions of MPJ.

3.1 MPJ.Init(String [] args) for Java

This Init() function establishes all to all connections using Java sockets. First, MPJ receives initialization commands such as rank, am slave, master node, number of processes, etc. Such arguments are mainly used to identify each process.

On the master process, Init() creates a ServerSocket and accepts connections from slaves until the number of connections equal the number of processes. Following each connection, the master process reads the connecting process's rank and identifies each connection with that rank, storing the rank-connection information in JavaComm. Then, the master broadcasts the ranks and their corresponding hostname to all slave processes. At this point, the master process's Init() is complete.

For the slave processes, they first connect to the master process, send their rank, and then receive a table with other slaves' ranks and their hostnames. Once such information is exchanged with the master, the slaves will connect with each other. First, all of the slaves except for the highest ranking slave will create ServerSockets. Afterwards, the lowest ranking slave will accept connections from higher ranking slaves. The lowest

ranking slave will then receive a rank from the connection it received and update its rank-connection table (much like the master). When the lowest ranking slave has received all connections, its Init() is complete. The second lowest ranking slave then accepts connections from higher ranking slaves, and the process repeats until the second highest ranking slave has accepted a connection from the highest ranking slave, indicating all slaves are connected to all other slaves. At this point, the Init() process is complete for Java sockets.

3.2 MPJ.Init(String [] args, IpTableEntry[] ipTable, GridTcp)

This Init() function establishes all to all connections for GridTcp sockets. The algorithm for this case is very similar to the Init() for Java sockets (Section 3.1). However, the difference is that since the UserProgWrapper of GridTcp pre-initializes an ipTable that corresponds to each of the hosts, all processes already know about each other. Thus, the process simply utilizes the connection algorithm described in Section 3.1 without distinguishing between the master and slaves (and also no broadcast of rank-connection information by master). Each process simply connects to all lower ranking processes while accepting connections from higher ranking processes.

3.3 Communicator.Send(Object[] buf, int offset, int count, Datatype type, int dest, int tag)

The Send() function takes in various parameters describing the datatype, the send count, the send buffer, the offset, the destination rank, and the message tag. The send buffer must be an array. The datatype is actually a Datatype object from MPJ, such as MPJ.SHORT.

Before sending a message, the Send() function first creates a header for the message, including the message's type, size in bytes, count, and the tag. Send() then serializes the message using the Datatype specified in the function parameters, and writes the header along with the serialized message to the output stream corresponding to the destination rank.

Note that this a blocking operations.

3.4 Communicator.Recv(Object[] buf, int offset, int count, Datatype type, int src, int tag)

When a user calls the Recv() function, Recv will perform a blocking read operation on the input stream corresponding to the source rank. First, Recv() reads 16 bytes of the message header and then reads the rest of the body with respect to the size defined in the header. Then, Recv() will deserialize the message using the correct Datatype. The completed message is stored in an MPJMessage, which is then checked against the parameters of Recv(). If the tag or datatype does not match, the message is stored in a message queue, and Recv() will read again for a new message.

If MPJ.ANY_SOURCE is specified, Recv() will poll each socket and read from the first socket with available data, and then check the tag.

If MPJ.ANY_TAG is specified, then Recv() will return the message if the datatype matches the parameter.

Recv() will crash if the count parameter is smaller than the actual message's count.

3.5 Pack(Object[] inbuf, int offset, int incount, Datatype type, byte[] outbuf, int position)

Pack is similar to Send(). It uses the Datatype to serialize the message into the provided output buffer, and then returns the updated position. If the Datatype is an MPJ.OBJECT, each object is serialized individually (rather than as a buffer), so they can be deserialized or extracted individually.

3.6 Unpack (byte[] inbuf, int position, java.lang.Object outbuf, int offset, int outcount, Datatype datatype)

Unpack performs the opposite of the Pack() operation. It will deserialize the input buffer into the output buffer using the corresponding Datatype's deserialize operation.

3.7 Barrier()

Barrier() is simple in that rank 0 will receive a message from all other ranks, and broadcast another message once it has received from all other ranks. Thus, each process is blocked until all processes have called Barrier().

3.8 Bcast(Object[] buf, int offset, int count, Datatype type, int root)

Bcast() will broadcast the message specified in buf from rank root to all other processes. Bcast currently follows a tree-structured algorithm, which reduces the number of send stages to $\log_2(n)$. Such an algorithm should be much better performance when broadcasting large messages. Bcast() on the root process will send the buffer to all other processes in the communicator, while the non-root processes will receive.

3.9 Reduce(Object sendbuf, int sendoffset, Object recvbuf, int recvoffset, int count, Datatype datatype, Op op, int root)

Reduce is a collective communication call, meaning all processes within the communicator are involved. Reduce will perform an operation defined by the op parameter. Operations are done element-wise on every send buffer and the result is stored in the root's buffer.

For example with MPJ.SUM as the operation:

	send[0]	send[1]	send[2]	send[3]
p0	1	2	3	4
p1	1	2	3	4
p2	1	2	3	4
recvbuf	[0] = 3	[1] = 6	[2] = 9	[3] = 12

Valid arguments for the op parameter are the following:

MPJ.BAND
MPJ.BOR
MPJ.BXOR
MPJ.LAND
MPJ.LOR
MPJ.LXOR
MPJ.MAXLOC
MPJ.MAX
MPJ.MINLOC
MPJ.MIN
MPJ.PROD
MPJ.SUM

Note that the MAX/MIN LOC operations requires an DataLoc (provided) array as the root's buffer.

3.10 Other Communications Algorithms

Other communications algorithms are compatible with mpiJava. Note that all of these functions require an array as the input and output buffers (even though the function prototype may just require an Object, this is type-casted into an array of appropriate type later). Their documentation can be found at www.hpjava.org. Please see online examples or the provided example code on how to use the following functions:

Isend() – Spawns a thread to send a message, like send().

```
public Request Isend(Object buf,  
                    int offset,  
                    int count,  
                    Datatype type,  
                    int dest,  
                    int tag)
```

Irecv() – Spawns a thread to receive a message, like recv().

```
public Request Irecv(Object buf,  
                   int offset,  
                   int count,  
                   Datatype type,  
                   int src,  
                   int tag)
```

Gather() – Gather all data in the inbuffers at all the processes to the outbuf of the root in rank order.

```
public void Gather(Object sendbuf,  
                  int sendoffset,  
                  int sendcount,  
                  Datatype sendtype,  
                  Object recvbuf,  
                  int recvoffset,  
                  int recvcount,  
                  Datatype recvtype,  
                  int root)
```

Gatherv() – Variation of Gather(), where different displaced input elements are allowed.

```
public void Gatherv(Object sendbuf,  
                   int sendoffset,  
                   int sendcount,  
                   Datatype sendtype,  
                   Object recvbuf,  
                   int recvoffset,  
                   int[] recvcounts,  
                   int[] displs,  
                   Datatype recvtype,  
                   int root)
```

AllGather() – Like Gather(), but result of outbuf is sent to all processes.

```
public void Allgather(Object sendbuf,  
                     int sendoffset,  
                     int sendcount,  
                     Datatype sendtype,  
                     Object recvbuf,  
                     int recvoffset,  
                     int recvcount,  
                     Datatype recvtype)
```


AllGatherv() – Like **Gatherv()**, but result of **outbuf** is sent to all processes.

```
public void Allgatherv(Object sendbuf,  
    int sendoffset,  
    int sendcount,  
    Datatype sendtype,  
    Object recvbuf,  
    int recvoffset,  
    int[] recvcounts,  
    int[] displs,  
    Datatype recvtype)
```

Scatter() – Opposite of **Gather()**. Scatter the buffer on the root to all processes in rank order.

```
public void Scatter(Object sendbuf,  
    int sendoffset,  
    int sendcount,  
    Datatype sendtype,  
    Object recvbuf,  
    int recvoffset,  
    int recvcount,  
    Datatype recvtype,  
    int root)
```

Scatterv() – Variation of **Scatter()** where differently displaced input elements are allowed.

```
public void Scatterv(Object sendbuf,  
    int sendoffset,  
    int[] sendcounts,  
    int[] displs,  
    Datatype sendtype,  
    Object recvbuf,  
    int recvoffset,  
    int recvcount,  
    Datatype recvtype,  
    int root)
```

AllToAll() – Sends each element of inbuf to the corresponding ranking process. Do on all processes, store in rank order.

```
public void Alltoall(Object sendbuf,  
                    int sendoffset,  
                    int sendcount,  
                    Datatype sendtype,  
                    Object recvbuf,  
                    int recvoffset,  
                    int recvcount,  
                    Datatype recvtype)
```

AllToAllv() – Variation of AllToAll() where differently displaced input elements are allowed.

```
public void Alltoallv(Object sendbuf,  
                     int sendoffset,  
                     int[] sendcount,  
                     int[] sdispls,  
                     Datatype sendtype,  
                     Object recvbuf,  
                     int recvoffset,  
                     int[] recvcount,  
                     int[] rdispls,  
                     Datatype recvtype)
```

3.11 SetBufferSize(int numbytes)

MPJ uses permanent buffers to help with serialization and deserialization performance.

As such, communicator has a SetBufferSize(int numbytes) function.

Example:

```
MPJ.COMM_WORLD.SetBufferSize(655365);
```

4.0 PingPong Performance

Java Sockets

Optimal:

loop = 25, msgSize = 512bytes, elapsedTime=8msec
transfer rate = 25.6Mbps
loop = 25, msgSize = 1024bytes, elapsedTime=8msec
transfer rate = 51.2Mbps
loop = 25, msgSize = 2048bytes, elapsedTime=9msec
transfer rate = 91.02222222222223Mbps
loop = 25, msgSize = 4096bytes, elapsedTime=9msec
transfer rate = 182.04444444444445Mbps
loop = 25, msgSize = 8192bytes, elapsedTime=11msec
transfer rate = 297.89090909090913Mbps
loop = 25, msgSize = 16384bytes, elapsedTime=10msec
transfer rate = 655.36Mbps
loop = 25, msgSize = 32768bytes, elapsedTime=18msec
transfer rate = 728.1777777777778Mbps
loop = 25, msgSize = 65536bytes, elapsedTime=38msec
transfer rate = 689.8526315789475Mbps
loop = 25, msgSize = 131072bytes, elapsedTime=58msec
transfer rate = 903.944827586207Mbps
loop = 25, msgSize = 262144bytes, elapsedTime=98msec
transfer rate = 1069.9755102040817Mbps
loop = 25, msgSize = 524288bytes, elapsedTime=174msec
transfer rate = 1205.2597701149425Mbps
loop = 25, msgSize = 1048576bytes, elapsedTime=309msec
transfer rate = 1357.3799352750812Mbps

MPJ:

loop = 25, msgSize = 512bytes, elapsedTime=26msec
transfer rate = 7.8769230769230765Mbps
loop = 25, msgSize = 1024bytes, elapsedTime=9msec
transfer rate = 45.51111111111111Mbps
loop = 25, msgSize = 2048bytes, elapsedTime=9msec
transfer rate = 91.02222222222223Mbps
loop = 25, msgSize = 4096bytes, elapsedTime=9msec
transfer rate = 182.04444444444445Mbps
loop = 25, msgSize = 8192bytes, elapsedTime=10msec
transfer rate = 327.68Mbps
loop = 25, msgSize = 16384bytes, elapsedTime=15msec
transfer rate = 436.9066666666667Mbps
loop = 25, msgSize = 32768bytes, elapsedTime=20msec
transfer rate = 655.36Mbps
loop = 25, msgSize = 65536bytes, elapsedTime=36msec
transfer rate = 728.1777777777778Mbps
loop = 25, msgSize = 131072bytes, elapsedTime=56msec
transfer rate = 936.2285714285714Mbps
loop = 25, msgSize = 262144bytes, elapsedTime=109msec
transfer rate = 961.9963302752293Mbps
loop = 25, msgSize = 524288bytes, elapsedTime=188msec
transfer rate = 1115.5063829787234Mbps
loop = 25, msgSize = 1048576bytes, elapsedTime=337msec
transfer rate = 1244.60059347181Mbps

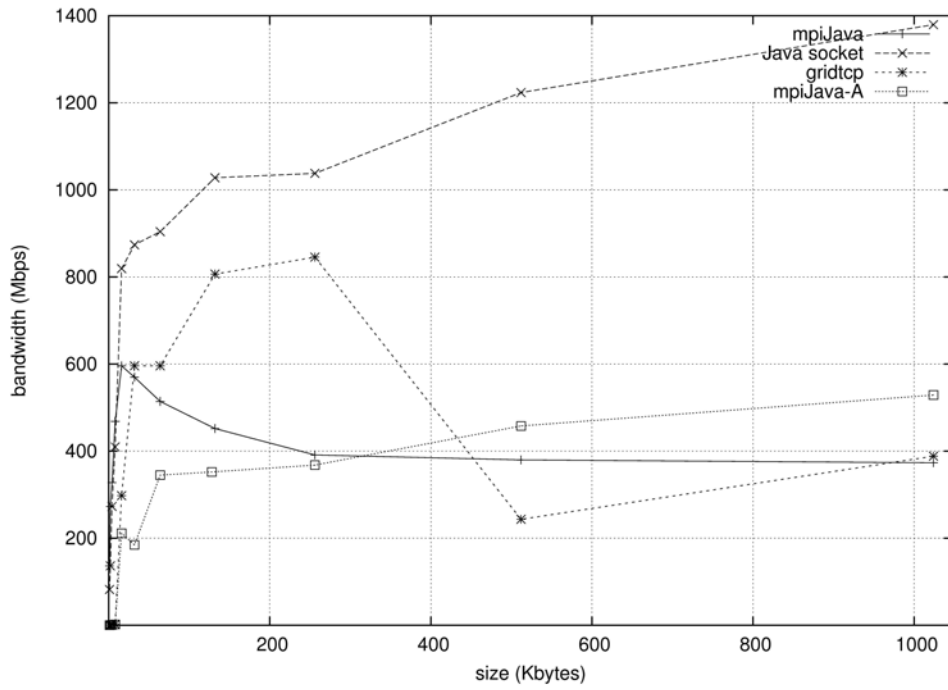
GridTcp

Optimal:

loop = 25, msgSize = 512bytes, elapsedTime=1789msec
transfer rate = 0.1144773616545556Mbps
loop = 25, msgSize = 1024bytes, elapsedTime=1749msec
transfer rate = 0.23419096626643798Mbps
loop = 25, msgSize = 2048bytes, elapsedTime=1789msec
transfer rate = 0.4579094466182224Mbps
loop = 25, msgSize = 4096bytes, elapsedTime=1743msec
transfer rate = 0.9399885255306942Mbps
loop = 25, msgSize = 8192bytes, elapsedTime=1758msec
transfer rate = 1.8639362912400457Mbps
loop = 25, msgSize = 16384bytes, elapsedTime=22msec
transfer rate = 297.89090909090913Mbps
loop = 25, msgSize = 32768bytes, elapsedTime=22msec
transfer rate = 595.7818181818183Mbps
loop = 25, msgSize = 65536bytes, elapsedTime=41msec
transfer rate = 639.3756097560976Mbps
loop = 25, msgSize = 131072bytes, elapsedTime=61msec
transfer rate = 859.4885245901639Mbps
loop = 25, msgSize = 262144bytes, elapsedTime=128msec
transfer rate = 819.2Mbps
loop = 25, msgSize = 524288bytes, elapsedTime=864msec
transfer rate = 242.7259259259259Mbps
loop = 25, msgSize = 1048576bytes, elapsedTime=1045msec
transfer rate = 401.3688038277512Mbps

MPJ:

loop = 25, msgSize = 512bytes, elapsedTime=1815msec
transfer rate = 0.1128374655647383Mbps
loop = 25, msgSize = 1024bytes, elapsedTime=1755msec
transfer rate = 0.2333903133903134Mbps
loop = 25, msgSize = 2048bytes, elapsedTime=1759msec
transfer rate = 0.46571915861284824Mbps
loop = 25, msgSize = 4096bytes, elapsedTime=1798msec
transfer rate = 0.9112347052280311Mbps
loop = 25, msgSize = 8192bytes, elapsedTime=1747msec
transfer rate = 1.875672581568403Mbps
loop = 25, msgSize = 16384bytes, elapsedTime=31msec
transfer rate = 211.40645161290325Mbps
loop = 25, msgSize = 32768bytes, elapsedTime=71msec
transfer rate = 184.60845070422536Mbps
loop = 25, msgSize = 65536bytes, elapsedTime=76msec
transfer rate = 344.92631578947373Mbps
loop = 25, msgSize = 131072bytes, elapsedTime=148msec
transfer rate = 354.2486486486486Mbps
loop = 25, msgSize = 262144bytes, elapsedTime=285msec
transfer rate = 367.9214035087719Mbps
loop = 25, msgSize = 524288bytes, elapsedTime=458msec
transfer rate = 457.8934497816594Mbps
loop = 25, msgSize = 1048576bytes, elapsedTime=793msec
transfer rate = 528.9160151324086Mbps



For Java sockets, PingPong performance is about 1-5% slower than optimal performance.

For GridTcp sockets, PingPong performance is significantly slower, at around 30-50% less than optimal.

5.0 Performance Analysis

At this point, the performance for Java sockets seem very reasonable. However, it should be noted that such performance is the optimal MPJ performance, in that no data conversion is needed.

When conversion of data is required, MPJ performance with Java sockets drop to slower levels than the original MPI-Java. Unfortunately this means that when dealing with data that needs to be serialized, MPJ-GridTCP runs even slower.

It should also be noted that to achieve good performance with java sockets, permanent send and receive buffers are used. It has been found that such buffers are very expensive to create and discard, and will cause a performance degradation if they are created and discarded for every Send() and Recv() pair, especially for buffers greater than 512kb. How such permanent buffers will affect GridTcP's checkpointing performance should be examined more closely

5.1 Object Serialization

Object serialization and deserialization is done using Object IO streams. It should be noted that any user objects should overload the default serialization methods if object serialization performance is critical, as the default methods provided by java are very slow. However, object serialization in MPJ-Java is still 20-30% faster than mpiJava object serialization.

5.2 Primitive Serialization

Even after some tuning of serialization of primitives, the performance cost is still very heavy. This is due to the high number of instructions needed to serialize a piece of data. For example, look at the following Java code:

```
int x; //for just 1 integer
byte[] arr[4]; //extra memory cost
arr[3] = (byte) ( x );
arr[2] = (byte) ( x >>> 8); //shift, cast, copy
arr[1] = (byte) (x >>> 16); //repeat
arr[0] = (byte) (x >>> 24);
```

To serialize a single integer, there needs to be extra storage and quite a few instructions, including assignment, cast, and shift used multiple times. In C/C++, this operation would be simply a pointer to the int, followed by the number of bytes of the int.

This creates a major problem with collective communications, where multiple send/recv pairs may be needed. There is a trade off between creating byte buffers and repeatedly serializing/deserializing into the user's buffers, both of which are extremely slow.

6.0 Initialization and Use

Currently, MPJ files are not in any sort of package, so you will either need to put your program files into a folder with MPJ files or use the `-classpath java` command to point to a folder containing MPJ files.

The following is the command line usage of MPJ.

Master process (start this first):

```
> java -classpath <paths to needed folders> <master node name> port -np <number of processes>
```

Slave processes:

```
> java -classpath <paths> <master node name> port -yourrank <rank> -amslave -yourname <hostname> -np <number of processes>
```

For example, when trying to run the pingpong program between `mnode0` and `mnode0`:

```
Mnode0> java -classpath ~uwagent/uwagent/GridTcp/GridTcp.jar:. PingPong mnode0 10000 -np 2
```

```
Mnode1> java -classpath ~uwagent/uwagent/GridTcp/GridTcp.jar:. PingPong mnode0 10000 -yourrank 1 -yourname mnode1 -amslave -np 2
```

Note that when initializing a GridTcp session, `-yourname <name>` is *not* needed. Also, you will need to pass the appropriate arguments to `UserProgWrapper` first.

After that, any program written for the original `mpiJava` should run correctly, provided that the needed communications functions are implemented. You should only need to remove the packages (unless this is changed later) and change `MPI.something` to `MPJ.something`.

6.1 mpjrun

Since its inception, `mpjrun` has become slightly outdated with the continuous changes made to MPJ. However, it should still perform automatic initialization of MPJ programs with some slight tweaking. The drawback to using `mpjrun` is that it requires multiple threads per process to monitor `stderr` and `stdout`. This will result in significant performance hits.

7.0 Continuation of MPJ

Future improvements include implementing an MPJException class to handle various exceptions and also additional communications algorithms. Also, it may be useful to improve the various collective communication algorithms for better performance, such as using permanent buffers for serialization and storage if memory allows.