

Design of scalable Java message-passing communications over InfiniBand

Roberto R. Expósito · Guillermo L. Taboada · Juan Touriño · Ramón Doallo

Submitted: December 15, 2010

Abstract This paper presents `ibvdev` a scalable and efficient low-level Java message-passing communication device over InfiniBand. The continuous increase in the number of cores per processor underscores the need for efficient communication support for parallel solutions. Moreover, current system deployments are aggregating a significant number of cores through advanced network technologies, such as InfiniBand, increasing the complexity of communication protocols, especially when dealing with hybrid shared/distributed memory architectures such as clusters. Here, Java represents an attractive choice for the development of communication middleware for these systems, as it provides built-in networking and multithreading support. As the gap between Java and compiled languages performance has been narrowing for the last years, Java is an emerging option for High Performance Computing (HPC).

The developed communication middleware `ibvdev` increases Java applications performance on clusters of multi-core processors interconnected via InfiniBand through: (1) providing Java with direct access to InfiniBand using InfiniBand Verbs API, somewhat restricted so far to MPI libraries; (2) implementing an efficient and scalable communication protocol which obtains start-up latencies and bandwidths similar to MPI performance results; and (3) allowing its integration in any Java parallel and distributed application. In fact, it has been successfully integrated in the Java messaging library MPJ Express.

The experimental evaluation of this middleware on an InfiniBand cluster of multi-core processors has shown significant point-to-point performance benefits, up to 85% start-up latency reduction and twice the bandwidth compared to previous Java middleware on InfiniBand. Additionally, the impact of

Roberto R. Expósito · Guillermo L. Taboada · Juan Touriño · Ramón Doallo
Computer Architecture Group, Dept. of Electronics and Systems,
University of A Coruña (Spain)
E-mail: {rreye,taboada,juan,doallo}@udc.es

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

ibvdev on message-passing collective operations is significant, achieving up to one order of magnitude performance increases compared to previous Java solutions, especially when combined with multithreading. Finally, the efficiency of this middleware, which is even competitive with MPI in terms of performance, increments the scalability of communications intensive Java HPC applications.

Keywords Message-Passing in Java (MPJ) · InfiniBand · Multi-core Architectures · High Performance Computing · Remote Direct Memory Access (RDMA) · Performance Evaluation

1 Introduction

Java is the leading programming language both in academia and industry environments, and it is an emerging alternative for High Performance Computing (HPC) [1] due to its appealing characteristics: built-in networking and multithreading support, object orientation, automatic memory management, platform independence, portability, security, an extensive API and a wide community of developers. Furthermore, in the era of multi-core processors, the use of Java threads is considered a feasible option to harness the performance of these processors.

Java initially was severely criticized for its poor computational performance [2], but the performance gap between Java and native (compiled) languages likes C and Fortran has been narrowing for the last years. The main reason is that the Java Virtual Machine (JVM), which executes Java applications, is now equipped with Just-in-Time (JIT) compilers that obtain native performance from Java bytecode. Nevertheless, the tremendous improvement in its computational performance is not enough for Java to be a successful language in the area of parallel computing, as the performance of the communications is also essential to achieve high scalability in Java for HPC.

Message-passing is the most widely used parallel programming paradigm as it is highly portable, scalable and usually provides good performance. It is the preferred choice for parallel programming distributed memory systems such as multi-core clusters, currently the most popular system deployments due to their scalability, flexibility and interesting cost/performance ratio. Here, Java represents an attractive alternative to languages traditionally used in HPC, such as C or Fortran, for the development of applications for these systems as it provides built-in networking and multithreading support, key features for taking full advantage of hybrid shared/distributed memory architectures. Thus, Java can use threads in shared memory (intra-node) and its networking support for distributed memory (inter-node) communications.

The increasing number of cores per system demands efficient and scalable message-passing communication middleware. However, up to now Message-Passing in Java (MPJ) implementations have been focused on providing portable communication devices, rather than concentrate on developing efficient low-level communication devices on high-speed networks. The lack of efficient high-speed networks in Java, due to its inability to control the underlying

specialized hardware, results in lower performance than MPI, especially for short messages. This paper presents a scalable and efficient Java low-level message-passing communication device, `ibvdev`, aiming to its integration in MPJ implementations in order to provide higher performance on InfiniBand multi-core clusters. In fact, it has been already integrated successfully in the MPJ library MPJ Express [3] (<http://mpj-express.org>).

The structure of this paper is as follows: Section 2 presents InfiniBand background information. Section 3 introduces the related work. Section 4 describes the design and implementation of the efficient `ibvdev` middleware, covering in detail the operation of the communication algorithms that provide the highest performance over InfiniBand. Section 5 shows the performance results of the implemented library on an InfiniBand multi-core cluster. The evaluation consists of a micro-benchmarking of point-to-point and collectives primitives, as well as a kernel/application benchmarking in order to analyze the impact of the use of the library on their overall performance. Section 6 summarizes our concluding remarks.

2 Java Communications over InfiniBand

2.1 InfiniBand Architecture

The InfiniBand Architecture (IBA) [4] defines a System Area Network (SAN) for interconnecting processing nodes and I/O nodes. In an InfiniBand network, processing nodes and I/O nodes are connected to the fabric by Channel Adapters (CA). Channel Adapters usually have programmable DMA engines with protection features. There are two kinds of channel adapters: Host Channel Adapter (HCA) and Target Channel Adapter (TCA). HCAs sit on processing nodes and TCAs connect I/O nodes to the fabric.

The InfiniBand communication stack consists of different layers. The interface presented by Channel adapters to consumers belongs to the transport layer. A queue-based model is used in this interface. A Queue Pair (QP) in InfiniBand Architecture consists of two queues: a send queue and a receive queue. The send queue holds instructions to transmit data and the receive queue holds instructions that describe where received data has to be placed. Communication operations are described in Work Queue Requests (WQR), or descriptors, and submitted to the work queue. Once submitted, a Work Queue Request becomes a Work Queue Element (WQE). WQEs are executed by Channel Adapters. The completion of work queue elements is reported through Completion Queues (CQs). Once a work queue element is finished, a completion entry is placed in the associated completion queue. Applications can check the completion queue to see if any work queue request has been finished.

2.1.1 Channel and Memory Semantics

InfiniBand Architecture supports both channel and memory semantics. In channel semantics, send/receive operations are used for communication. To receive a message, the programmer posts a receive descriptor which describes where the message should be put at the receiver side. At the sender side, the programmer initiates the send operation by posting a send descriptor. The send descriptor describes where the source data is but does not specify the destination address at the receiver side. When the message arrives at the receiver side, the hardware uses the information in the receive descriptor to put data in the destination buffer. Multiple send and receive descriptors can be posted and they are consumed in FIFO order. The completion of descriptors are reported through CQs.

In memory semantics, Remote Direct Memory Access (RDMA) write and RDMA read operations are used instead of send and receive operations. These operations are one-sided and do not incur software overhead at the other side. The sender initiates RDMA operations by posting RDMA descriptors. A RDMA descriptor contains both the local data source address and the remote data destination address. At the sender side, the completion of a RDMA operation can be reported through CQs. The operation is transparent to the software layer at the receiver side.

Both communication semantics require communication memory to be registered with InfiniBand hardware and pinned in memory. The registration operation involves informing the network-interface of the virtual to physical address translation of the communication memory. The pinning operation requires the operating system to mark the pages corresponding to the communication memory as non-swappable. Thus, communication memory stays locked in physical memory, and the network-interface can access it as desired.

2.1.2 Transport Services

There are five transport modes defined by the InfiniBand specification: Reliable Connection (RC), eXtended Reliable Connection (XRC), Reliable Datagram (RD), Unreliable Connection (UC), and Unreliable Datagram (UD). All transports provide a checksum verification.

Reliable Connection (RC) is the most popular transport service for implementing MPI over InfiniBand. As a connection-oriented service, a QP with RC transport must be dedicated to communicating with only one other QP. A process that communicates with N other peers must have at least N QPs created. The RC transport provides almost all the features available in InfiniBand, most notably reliable send/receive, RDMA and atomic operations.

RC transport made no distinction between connecting a process (generally one per core for MPI) and connecting a node. Thus, the associated resource consumption increased directly in relation to the number of cores in the system. To address this problem eXtended Reliable Connection (XRC) was introduced. Instead of having a per-process cost, XRC was designed to allow a single

connection from one process to an entire node. XRC provides the services of the RC transport, but defines a very different connection model and method for determining data placement on the receiver in channel semantics. When using the RC transport, the connection model is purely based on processes. By contrast, XRC allows connection optimization based on the location of a process. The node of the peer to connect to is now taken into account, so instead of requiring a new QP for each process, now each process needs only have one QP per node to be fully connected. This reduces the number of QPs required by a factor of number of cores per node.

Unreliable Connection (UC) provides connection-oriented service with no guarantees of ordering or reliability. It supports RDMA write capabilities and sending messages larger than the Maximum Transmission Unit (MTU) size. Being connection-oriented in nature, every communicating peer requires a separate QP. In regard to resources required, it is identical to RC, while no providing reliable service. Thus, it appears unattractive for implementing MPI over this transport.

Unreliable Datagram (UD) is a connection-less and unreliable transport, the most basic transport specified for InfiniBand. As a connection-less transport, a single UD QP can communicate with any number of other UD QPs. However, the UD transport has a number of limitations. The UD transport does not provide any reliability: lost packets are not reported and the arrival order is not guaranteed. However, this can be solved relying on Reliable Datagram (RD). Moreover, UD transport does not enable RDMA. All communication must be performed using channel semantics, i.e. send/receive.

Table 1 shows the available operations for each transport service, since not all transport service support all operations, which has to be taken into account for a message-passing middleware implementation.

Table 1 Operations available for each transport service

Operation	RC	XRC	UC	RD	UD
SEND (WITH IMMEDIATE)	X	X	X	X	X
RECEIVE	X	X	X	X	X
RDMA WRITE (WITH IMMEDIATE)	X	X	X	X	
RDMA READ	X	X		X	
ATOMIC	X	X		X	

2.1.3 Shared Receive Queues

Shared Receive Queues (SRQs) were introduced in the InfiniBand 1.2 specification to address scalability issues with InfiniBand memory usage. In order to receive a message on a QP, a receive buffer must be posted in the Receive Queue (RQ) of that QP. To achieve high-performance, MPI implementations pre-post buffers to the RQ to accommodate unexpected messages. When using the RC transport of InfiniBand, one QP is required per communicating

peer. However, this task of pre-posting receives on each QP can have very high memory requirements for communication buffers. Recognizing that such buffers could be pooled, SRQ support was added so instead of connecting a QP to a dedicated RQ, buffers could be shared across QPs. In this method, a smaller pool can be allocated and then refilled on demand instead of pre-posting on each connection.

2.2 Message-passing Communication Devices

Message-passing libraries usually support new transport protocols through the use of pluggable low-level communication devices, such as Abstract Device Interface (ADI) in MPICH, Byte Transfer Layer (BTL) in OpenMPI, and `xdev` [5] in MPJ Express. These communication devices abstract the particular operation of a communication protocol, such Myrinet eXpress (MX), uDAPL (user Direct Access Programming Library), InfiniBand Verbs (IBV), Shared Memory or SCTP (Stream Control Transmission Protocol), conforming to an API on top of which the message-passing library implements its communications.

Figure 1 presents an overview of the communications support of MPJ applications on the high-speed Myrinet network, on Gigabit Ethernet and on shared memory. From top to bottom, MPJ applications rely on MPJ libraries, whose communication support is implemented in the device layer. Current Java communication devices are implemented either on JVM threads (`smpdev`, a multithreading device), on sockets over the TCP/IP stack (`niodev` on Java NIO sockets and `iodev` on Java IO sockets), or on native communication layers such as Myrinet eXpress (`mxdev`, a device on MX).

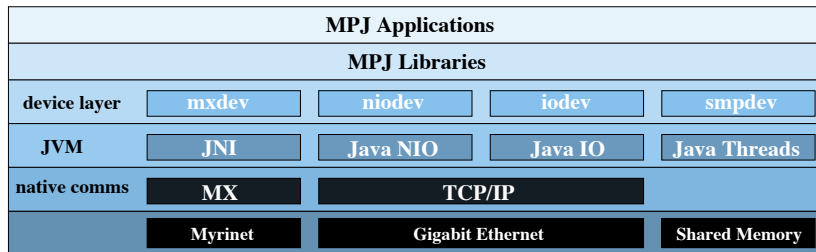


Fig. 1 Communications Support of MPJ Applications

Regarding InfiniBand, up to now no direct support was made available for MPJ applications to fully exploit the communication capability of InfiniBand networks. This lack of direct InfiniBand support in Java requires the use of upper layer protocols such as IPoIB [6] (IP over InfiniBand) TCP emulation, as shown in Figure 2, or SDP (Sockets Direct Protocol), the high performance native sockets library on InfiniBand. However, the use of IPoIB, the

only communication library that fully supports Java over InfiniBand, shows quite poor performance [7]. Moreover, when relying on SDP the performance generally improves, but this is not always possible. Regarding MPI libraries, its direct InfiniBand support has been implemented some years ago on top of InfiniBand Verbs (IBV) API (see Figure 2), achieving very high performance results. Therefore, our objective is the implementation of the direct InfiniBand support in Java on IBV through the development of a low-level Java communication device that can take advantage of InfiniBand RDMA transfers, thus outperforming significantly previous Java support on InfiniBand.

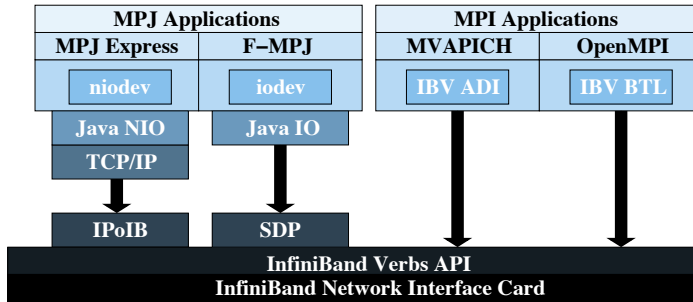


Fig. 2 MPI/MPJ Applications support on InfiniBand

3 Related Work

Current research on efficient Java communication libraries over InfiniBand is, to our knowledge, restricted to Jackal, Aldeia, Java Fast Sockets (JFS), Jdib and uStream projects, next presented. Jackal [8] is a Java DSM (Distributed Shared Memory) middleware for clusters with InfiniBand Verbs support, embracing also RDMA transfers, but it does not provide any API to Java developers as it only implements data transfers specifically for Jackal. Aldeia [9] is a proposal of an asynchronous sockets communication layer over InfiniBand whose preliminary results were encouraging, but requires an extra-copy, which incurs an important overhead to provide asynchronous write operations, whereas the read method is synchronous.

JFS [10] is our high performance Java socket implementation for efficient shared memory and high-speed networks support. JFS relies on SDP (see Figure 2) to support Java communication over InfiniBand. Moreover, JFS avoids the need for primitive data type array serialization and reduces buffering and unnecessary copies. Nevertheless, the use of the sockets API is a significant drawback to support efficient message-passing communications.

Jdib [11, 12] (Java Direct InfiniBand) is a Java encapsulation of IBV API which maximizes Java communication performance using directly, through

1 Java Native Interface (JNI), the InfiniBand RDMA mechanism. The main
2 contribution of Jdib is its direct access to RDMA, providing to performance-
3 concerned developers, for the first time, a Java RDMA API. Thus, Jdib signifi-
4 cantly outperforms its alternatives, currently limited to IPoIB- and SDP-based
5 solutions. The main drawbacks of Jdib are its low-level API and the JNI over-
6 head incurred for each Jdib operation.
7

8 uStream [13] is a user-level stream protocol implemented on top of IBV that
9 provides a higher level API than Jdib. In fact, uStream abstracts developers
10 from the most tedious operations in Jdib, such as the buffer management,
11 synchronization and the use of the IBV API, while fully exploiting InfiniBand
12 RDMA performance. Therefore, uStream is much more effective and easier to
13 use than Jdib for building parallel and distributed applications.
14

16 4 ibvdev: Efficient Java Communications over InfiniBand

18 This section presents the design and implementation of the `ibvdev` communi-
19 cation device, the Java message-passing middleware over InfiniBand developed
20 in this paper. Unlike VIA [14, 15], InfiniBand architecture does not specify
21 an API. Instead, it defines the functionality provided by HCAs to operating
22 systems in terms of Verbs (a “verb” is a semantic description of a function
23 that must be provided). The Verbs interface specifies such functionality as
24 transport resource management, multicast, work request processing and event
25 handling. The most important implementation used today of Verbs interface
26 is the IBV API provided by the OFED (OpenFabrics Enterprise Distribution)
27 driver distributed by the OpenFabrics Alliance [16]. IBV is also the lowest
28 level InfiniBand networking API for applications, available only in C language.
29 Therefore, any Java communication support on IBV must resort to JNI in order
30 to access IBV API and obtain the best possible performance, the target
31 of the communication middleware developed, `ibvdev`.
32
33

35 4.1 Message-Passing in Java Libraries

37 There have been several efforts [1] over the last decade to develop a Java
38 message-passing system since its introduction [17]. Most of these projects were
39 prototype implementations, without maintaining. Currently, the most relevant
40 ones in terms of uptake by the HPC community are `mpiJava` [18], `MPJ Ex-`
41 `press` [3], `MPJ/Ibis` [19] and `F-MPJ` [20].
42

43 `mpiJava` [18] is a Java messaging system that uses JNI to interact with
44 the underlying native MPI library. This project has been perhaps the most
45 successful Java HPC messaging system, in terms of uptake by the community.
46 However, although its performance is usually high, `mpiJava` currently only
47 supports some native MPI implementations, as wrapping a wide number of
48 function and heterogeneous runtime environments entails an important main-
49 taining effort. Additionally, this implementation presents instability problems,
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

1 derived from the native code wrapping (all MPJ methods are wrapped), and
2 has thread safety issues in the wrapper layer, being unable to take advantage of
3 multi-core systems through multithreading, even if the underlying MPI library
4 is thread safe.
5

6 MPJ Express is an MPJ implementation of the mpiJava 1.2 API [17] spec-
7 ification. MPJ Express is thread-safe and presents a modular design which
8 includes a pluggable architecture of communication devices that allows to com-
9 bine the portability of the “pure” Java New I/O package (Java NIO) communi-
10 cations (`niodev` device) with the high performance Myrinet support (through
11 the native Myrinet eXpress communication library in the `mxdev` device).
12

13 MPJ/Ibis [19] is an implementation of the JGF MPJ API [21] specifi-
14 cation on top of Ibis [22]. The design philosophy of Ibis is similar to MPJ
15 Express; it is possible to use 100% pure Java communication or use special
16 HPC hardware like Myrinet. There are two pure Java devices in Ibis. The first
17 called `TCPIbis` provides communication using the traditional `java.io` pack-
18 age. The second called `NIOIbis` uses the Java NIO package. Although `TCPIbis`
19 and `NIOIbis` provide blocking and non-blocking communication at the device
20 level, the higher-levels only use blocking versions of these methods. Neverthe-
21 less, MPJ/Ibis does not provide a multithreaded communication device, unlike
22 MPJ Express, key to harness the performance of multi-core processors.

23 F-MPJ[20] is our message-passing communication middleware that pro-
24 vides shared memory and high-speed networks (e.g., InfiniBand, Myrinet, and
25 SCI) communication support through the use of JFS. However, the use of
26 Java IO sockets in its communication device `iodev` limits scalability as the
27 progress engine of F-MPJ has to check every connection for incoming mes-
28 sages, unlike Java NIO sockets whose support is already implemented in the
29 `select` method.

30 MPJ Express project is currently the most active project in terms of adop-
31 tion by the HPC community, presence on academia and production environ-
32 ments, and available documentation. This project is also stable and publicly
33 available along with its source code at <http://mpj-express.org>. Therefore,
34 MPJ Express has been selected for the integration of the `ibvdev` middleware
35 in a production MPJ library.
36
37

38 4.2 MPJ Express Communication Devices Design 39

40 MPJ Express has a layered design that enables its incremental development
41 and provides the capability to update and swap layers in or out as required.
42 Thus, at runtime end users can opt to use a high performance proprietary net-
43 work device, or choose a pure Java device, based either on sockets or threads,
44 for portability.
45

46 Figure 3 illustrates an overview of the MPJ Express design and the differ-
47 ent levels of the software. From top to bottom, it can be seen that a message-
48 passing application in Java (MPJ application) calls MPJ Express point-to-
49 point and collective primitives. These primitives implement the MPJ commu-
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

nications API on top of the `xdev` layer, which has been designed as a pluggable architecture and provides a simple but powerful API. This design facilitates the development of new communication devices in order to provide custom implementations on top of specific native libraries and HPC hardware. Thus, `xdev` is portable as it presents a single API and provides efficient communication on different system configurations.

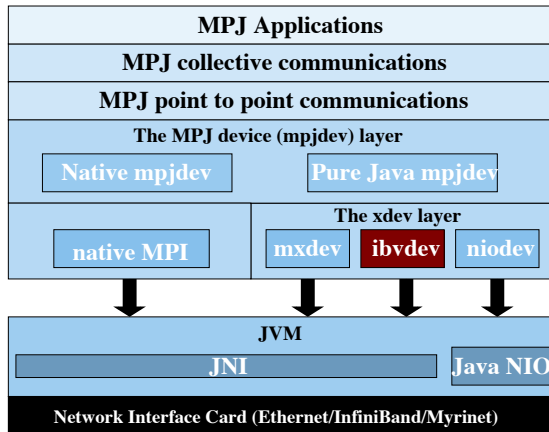


Fig. 3 Overview of the MPJ Express design including `ibvdev`

Figure 3 shows the three implementations of the `xdev` API for networked communication: `niodev` on Java NIO, and hence TCP/IP, and `mxdev` on Myrinet, as well as the developed `xdev` middleware for direct InfiniBand support, `ibvdev` (depicted in red).

4.2.1 `xdev` API Design

The `xdev` API, presented in Listing 1, has been designed with the goal of being simple and small, providing only basic communication methods, in order to ease the development of `xdev` devices. An `xdev` communication device is similar to the MPI communicator class, but with reduced functionality. The `init` method starts the communication device operation. The `id` method returns the identification (`ProcessID`) of the device. The `finish` method is the last method to be called and completes the device operation.

The `xdev` communication primitives only include point-to-point communication, both blocking (`send` and `recv`, like `MPI_Send` and `MPI_Recv`) and non-blocking (`isend` and `irecv`, like `MPI_Isend` and `MPI_Irecv`). Synchronous communications are also embraced (`srend` and `issend`). These communication methods use `PID` (`ProcessID`) objects instead of using ranks as arguments

to send and receive primitives. In fact, the `xxdev` layer is focused on providing basic communication methods and it does not deal with high level message-passing abstractions such as groups and communicators. Therefore, a PID object unequivocally identifies a device object.

```

1 public abstract class Device {
2     public static Device newInstance(String dev);
3     ProcessID [] init(String [] args);
4     ProcessID id();
5     void finish();
6
7     Request isend(Buffer buf, PID dest, int tag, int cntx);
8     void send(Buffer buf, PID dest, int tag, int cntx);
9     Request issend(Buffer buf, PID dest, int tag, int cntx);
10    void ssend(Buffer buf, PID dest, int tag, int cntx);
11    Status recv(Buffer buf, PID src, int tag, int cntx);
12    Request irecv(Buffer buf, PID src, int tag, int cntx, Status s);
13    Status probe(PID src, int tag, int cntx);
14    Status iprobe(PID src, int tag, int cntx);
15    Request peek();
16 }

```

Listing 1 API of the `xdev.Device` class

4.3 Communication Device Design

Figure 4 presents the overall design of the communication middleware, which consists of three distinct parts. The first is the definition of a new device, `ibvdev`, in the `xdev` layer of MPJ Express (1 in Figure 4). The analysis of the other high-speed network support in MPJ Express, the implementation of the `mxdev` device, reveals that it also uses native code via JNI to rely on the MX library, thus posing similar design issues as `ibvdev`. The MX library [23] provides a set of primitives similar to those needed to implement `xdev` interface, so there are a number of functions, such as `mx_isend`, `mx_issend`, `mx_irecv` and `mx_wait`, that are used in the JNI layer. Therefore, `mxdev` acts as a Java wrapper layer to MX library, so that the implementation of a method in `xdev` generally delegates directly in a native method that performs the requested operation in MX library. Nevertheless, the design of `mxdev` is not directly applicable to `ibvdev` since InfiniBand lacks an MX-style library that implements the functionality and operations that must be implemented in `xdev`. The available communication layer for `ibvdev` is the IBV API, which offers low-level methods for the management of the HCA InfiniBand card.

Therefore, an MX-like library has been defined in order to provide `ibvdev` with a set of communication primitives with message-passing semantics on InfiniBand, to ease the development of the `xdev` communication device. This library has been denominated IBV eXpress (IBVX) (2 in Figure 4). With this design, a native communication library has been implemented on top of IBV to provide basic message-passing communication primitives to higher level layers

(either Java or non Java). Thus, the new communication device `ibvdev` can rely on IBVX through JNI. The design of this layer allow the access to IBVX from MPJ Express through its `ibvdev` device (3 in Figure 4).

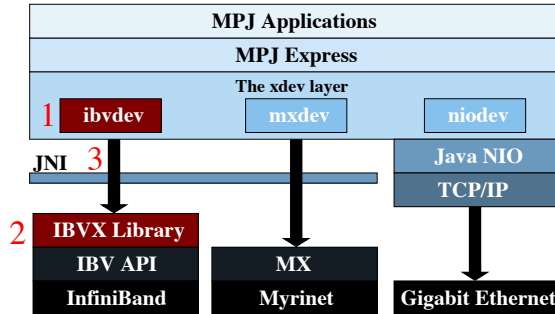


Fig. 4 Overall design of the communication library

4.3.1 IBV eXpress Library Design

The IBVX library is a scalable and high performance low-level C message-passing middleware for communication on InfiniBand systems. It has been designed using the same approach as `xdev` communication devices. In fact, there is a mapping of `xdev` methods to IBVX functions, except for methods `id` used for process identification, and `getSendOverhead` and `getRecvOverhead`, which are available only at the Java level as give information about the buffer handling. The IBVX API is presented in Listing 2. Like `xdev` API, IBVX includes only point-to-point communication, both blocking and non-blocking, and also synchronous communication support. In order to support non-blocking operations IBVX implements `IBV_Wait` and `IBV_Test` functions, which handle non-blocking operation requests.

```

1 IBV_Init(char **pNames, int *pList, int nProcs, int rank, int psl);
2 IBV_Finalize();
3 IBV_Isend(void *buf, int size, int dst, int tag, int ctx, Request *r);
4 IBV_Issend(void *buf, int size, int dst, int tag, int ctx, Request *r);
5 IBV_Irecv(void *buf, int size, int src, int tag, int ctx, Request *r);
6 IBV_Send(void *buf, int size, int dst, int tag, int ctx);
7 IBV_Ssend(void *buf, int size, int dst, int tag, int ctx);
8 IBV_Recv(void *buf, int size, int src, int tag, int ctx, Status *s);
9 IBV_Wait(Request *request, Status *status);
10 IBV_Test(Request *request, Status *status);
11 IBV_Iprobe(int src, int tag, int context, Status *status);
12 IBV_Probe(int src, int tag, int context, Status *status);
13 Request* IBV_Peek();

```

Listing 2 Public interface of the IBV eXpress library

4.3.2 `ibvdev` JNI Layer Design

The design of the JNI layer of `ibvdev` is quite straightforward as it acts as a thin wrapper over IBVX. Thus, each native method of `ibvdev` delegates on a native IBVX function through JNI, implementing a series of three steps: (1) get Java objects associated parameters required for calling the corresponding library function in IBVX; (2) call IBVX function; and (3) save the results in the appropriate attributes of the Java objects involved in the communication. As general rules in the implementation of the JNI layer it has been extensively used the caching of object references, thus minimizing the overhead associated with the JNI calls.

4.4 IBV eXpress Library Implementation

IBVX library implements non-blocking low-level communication primitives (see Listing 2) on top of IBV API. The first decision is the transport service used to create the queue pairs. Not all transports services support RDMA operations (see Table 1), whose support is desirable, so these transport services (UC and UD) are discarded.

Moreover, for RD and XRC transport services is not applicable the InfiniBand end-to-end flow control and this requires the development of a specific flow control software layer, which can add significant overhead if the implementation is not efficient. Therefore, the RC transport service has been selected as it provides reliability, delivery order, data loss detection and error detection.

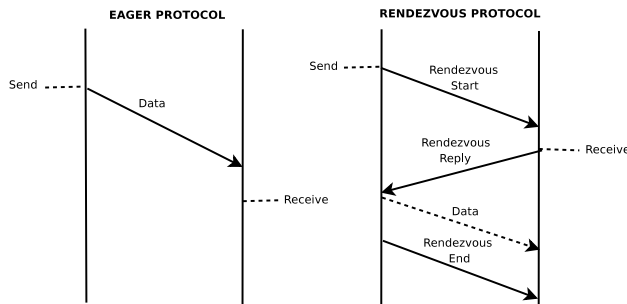
IBVX implements all communication operations as non-blocking communication primitives. Then, blocking communication support is implemented as a non-blocking primitive followed by an `IBV_Wait` call. Therefore, the basic set of functions implemented consists of `IBV_Init`, `IBV_Finalize` and non-blocking communication functions (`IBV_Isend`, `IBV_Issend`, `IBV_Irecv`), and the function that checks the completion of a non-blocking operation (`IBV_Test`). Thus, the operation that waits for the completion of a non-blocking operation (`IBV_Wait`) has been implemented following a strategy of polling (busy loop) as a continuous loop calling `IBV_Test` until the test is positive (thus minimizing latency). Blocking communication functions (`IBV_Send`, `IBV_Ssend` and `IBV_Recv`) have been implemented by a call to its corresponding non-blocking function followed by an `IBV_Wait` call. Moreover, the probe operation, which checks for incoming messages without actual receipt of any of them, has been also implemented in the non-blocking version `IBV_Iprobe`, where the blocking version (`IBV_Probe`) relies on the non-blocking operation completion.

4.4.1 IBV eXpress Communication Protocols

Message-passing libraries usually implement two different communication protocols:

- 1 **1. Eager Protocol:** the sender process eagerly sends the entire message to
2 the receiver. In order to achieve this, the receiver needs to provide a suf-
3 ficient number of buffers to handle incoming messages. This protocol has
4 minimal startup overheads and is used to implement low latency message
5 passing communication for smaller messages (typically < 128 KB, config-
6 urable threshold).
7
- 8 **2. Rendezvous Protocol:** this protocol negotiates (via control messages)
9 the buffer availability at the receiver side before the message is actually
10 transferred. This protocol is used for transferring large messages (typically
11 > 128 KB), whenever the sender is not sure whether the receiver actually
12 has enough buffer space to hold the entire message.
13

14 Figure 5 presents graphically the operation of Eager and Rendezvous pro-
15 tocols.



16 **Fig. 5** MPI Eager and Rendezvous Protocols

31 4.4.2 Message Format

32 The presence of control messages in the operation of the rendezvous protocol
33 and the need for a receiving process to uniquely distinguish a message, has
34 forced the introduction of message header before the actual data payload.
35 Thus, a message is defined as the union of a header of 20 bytes (starting from
36 the beginning), which is followed by the data payload, as shown in Figure 6.

37 The header consists of 5 fields of 4 bytes each representing in this order:
38 the process rank that sends the message, the destination process rank, the
39 tag or label of the message, the context to which it belongs, and the type of
40 message. All header fields are integers, for all types of messages.

41 4.4.3 Eager Protocol

42 The overhead of data copies is small for short messages, like eager protocol
43 and control messages, which are eagerly push through the network to achieve
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

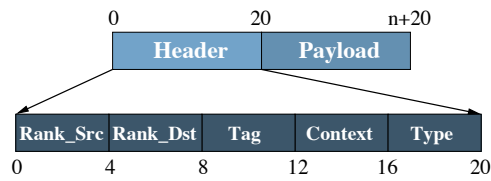


Fig. 6 Message format in IBVX library

the lowest latency. This operation matches with the semantic of InfiniBand send/receive communication.

In `IBV_Init` a reliable connection is set up between every two processes. For a single process, the send and receive queues of all connections are associated with a single CQ (Completion Queue). Through this CQ, the completion of all send and RDMA operations can be detected at the sender side. The completion of receive operations (or arrival of incoming messages) can also be detected through the CQ (see Figure 7).

The InfiniBand Architecture requires the pinning of buffers previous to the communication, thus they must be registered with the hardware. In the eager protocol implementation (shown in Figure 7), the buffer pinning and unpinning overhead is avoided by using a pool of pre-pinned, fixed size buffers for communication. For sending an eager data message, the data is copied to one of the buffers first and sent out from this buffer to the send queue (1 in Figure 7). At the receiver side, a number of buffers from the pool are pre-posted (2 in Figure 7). After the message is received, the payload is copied to the destination buffer (3 in Figure 7). The communication of control messages also uses this buffer pool as they are actually sent using the eager protocol.

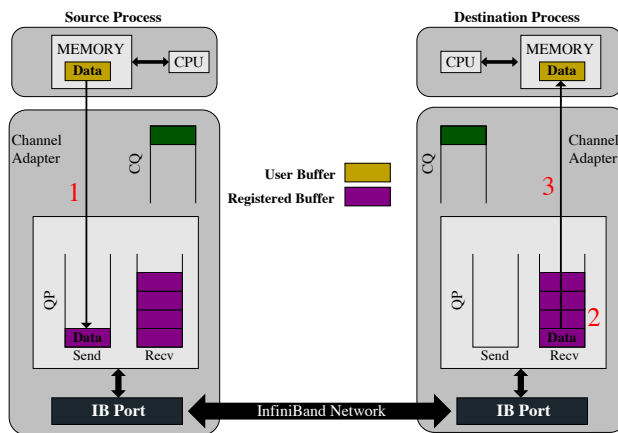


Fig. 7 Eager protocol implementation in IBVX

4.4.4 Rendezvous Protocol

When transferring large messages it is extremely beneficial to avoid extra data copies. A zero-copy Rendezvous protocol implementation can be achieved by using RDMA operations. The Rendezvous protocol negotiates the buffer availability at the receiver side. However, the actual data can be transferred either by using RDMA Write or RDMA Read. RDMA Write-based approaches can totally eliminate intermediate copies and efficiently transfer large messages. RDMA Read based approaches can enable both zero copy and computation and communication overlap. Similar approaches have been widely used for implementing MPI communications over different interconnects [24] [25].

The RDMA Write-based protocol is illustrated in Figure 8 (right). In this implementation, the buffers are pinned down in memory and the buffer addresses are exchanged via control messages. The sending process first sends a control message to the receiver (RNDZ_START). The receiver replies to the sender using another control message (RNDZ_REPLY). This reply message contains the receiving application's buffer information along with the remote key to access that memory region. The sending process then sends the large message directly to the receiver's application buffer by using RDMA Write (DATA). Finally, the sending process issues another control message (RNDZ_END) which indicates to the receiver that the message has been placed in the application buffer.

IBVX uses a *progress engine* to discover incoming messages and to make progress on outstanding sends. As can be seen in Figure 8, the RDMA Write based Rendezvous Protocol generates multiple control messages which have to be discovered by the *progress engine*. Since the *progress engine* operation is based on polling, it requires a call to the IBVX library.

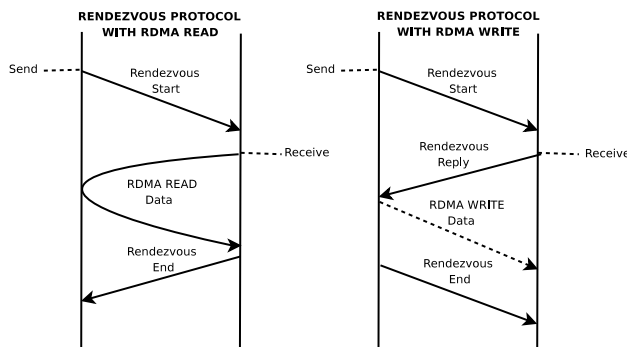


Fig. 8 Rendezvous Protocol alternatives

RDMA Read operation presents a small number of control messages and thus a reduced set of I/O bus transactions. In addition, since the receiver

can progress independently of the sender (once the RNDZ_START message is sent), the sender does not need to call any IBVX progress, the data transfer proceeds with RDMA Read without direct control of the sender.

The Rendezvous Protocol over RDMA Read is also illustrated in Figure 8 (left). Here the sending process begins with the RNDZ_START message, which has embedded the virtual address and memory handle information of the message buffer to be sent. Thus, upon the receipt of this RNDZ_START message all the information about the application buffer is available to the receiving process, and no RNDZ_REPLY message needs to be sent any more. Upon its discovery, the receiving process issues the DATA message over RDMA Read. When the operation has been completed, it informs the sending process by a RNDZ_END message. This approach, although simple, poses several design challenges that has to be addressed before directly utilize RDMA Read:

- Limited Outstanding RDMA Reads: The number of outstanding RDMA Reads on any QP is a fixed number (typically 8 or 16), decided during the QP creation.
- Issuing RNDZ_END Message: According to InfiniBand specification [4], Send or RDMA Write transactions, are not guaranteed to finish in order with outstanding RDMA Reads.

For these reasons, the Rendezvous protocol has been implemented with RDMA Write operation, in order to benefit from a more productive development.

4.4.5 Cache of Registered Buffers

In Rendezvous protocol, data buffers are pinned on-the-fly. However, the buffer pinning and unpinning overhead can be reduced by using the pin-down cache technique [26]. The idea is to maintain a cache of registered buffers. When a buffer is first registered, it is put into the cache. When the buffer is unregistered, the actual unregister operation is not carried out and the buffer stays in the cache. Thus the next time when the buffer needs to be registered, we need not to do anything because it is already in the cache. The effectiveness of Pin-down Cache depends on how often the application reuses its buffers. If the reuse rate is high, most of the buffer registration and de-registration operations can be avoided.

4.5 JNI Layer Implementation Details

The JNI layer is a wrapper for IBVX library, in order to make it accessible from Java. Therefore, it implements the functions that the `javah` utility generated in terms of native operations contained in communication device Java classes. The development of this layer must take into account the design of the MPJ Express buffering layer [27]. The use of this buffering layer incurs a copying overhead, that can be significant for large messages, and is considered

a performance bottleneck for MPJ Express, so the handling of this layer has to be implemented efficiently.

The core class of the buffering layer used for packing and unpacking data is `mpjbuf.Buffer`. This class provides two storage options: static and dynamic. Implementations of static storage use the interface `mpjbuf.RawBuffer`. It is possible to have alternative implementations of the static section depending on the actual raw storage medium. In addition, it also contains an attribute of type `byte[]` that represents the dynamic section of the message. Figure 9 shows two implementations of the `mpjbuf.RawBuffer` interface. The first, `mpjbuf.NIOBuffer` is an implementation based on `ByteBuffer`s. The second, `mpjbuf.NativeBuffer` is an implementation for the native MPI device, which allocates memory in the native C code. Figure 9 shows the primary buffering classes in the `mpjbuf` API.

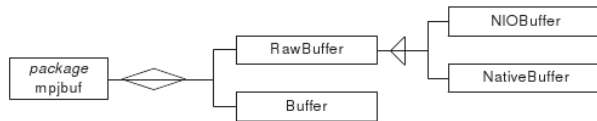


Fig. 9 Primary Buffering Classes in `mpjbuf`

Regarding `mpjbuf.Buffer` class design, it is necessary to handle at the JNI layer a second call to the IBVX library when communicating a buffer with data in the two sections (static and dynamic). To support this operation efficiently the first 4 bytes of the static buffer indicate the size of the dynamic part of the buffer. Thus, the overhead of this protocol in terms of buffering space, returned by `getSendOverhead` and `getRecvOverhead` methods, is 4 bytes. These methods, implemented for every MPJ Express communication device, are used to express the extra space needed in the static buffer to implement the buffering layer support, and they are profusely used when handling the buffer contents.

5 Performance Evaluation

This section presents a performance evaluation of the developed communication device `ibvdev`, compared to native MPI libraries (MVAPICH and OpenMPI) and the MPJ Express communications devices `niodev` over InfiniBand (using IPoIB) and `smpdev` for shared memory communication. This evaluation consists of a microbenchmarking of point-to-point data transfers (Subsection 5.2) and collective communications (Subsection 5.3), as well as an analysis of the impact on the overall performance of the use of the developed library on several representative MPJ codes (Subsection 5.4).

5.1 Experimental Configuration

The evaluation of `ibvdev` has been carried out in a cluster which consists of 8 nodes, each of them with 8 GB of RAM and 2 Intel Xeon E5520 quad-core Nehalem processors. Although each node has 8 cores, the HyperThreading (HT) is enabled so it is possible to run 16 processes per node concurrently. The interconnection networks are InfiniBand (16 Gbps of maximum theoretical bandwidth), with OFED driver 1.5, and Gigabit Ethernet (1 Gbps). The OS is Linux CentOS 5.3 with kernel 2.6.18 and the JVM is Sun JDK 1.6.0_13. The evaluated MPJ implementation is MPJ Express [28] version 0.36 (labeled MPJE in graphs) and the evaluated MPI implementations are MVAPICH [25] v1.2.0 and OpenMPI [24] v1.3.3. The PSL (Protocol Switch Limit) MPJ Express attribute, the threshold between eager and rendezvous send protocols, has been set to 128 KB message size for all the benchmarks. F-MPJ and MPJ/Ibis results are not shown for clarity purposes, apart from the fact that `ibvdev` is only integrated in MPJ Express. However, as they are sockets-based implementations, its performance is similar to `niodev` results.

5.2 Point-to-point Micro-benchmarking

In order to micro-benchmark MPJ point-to-point and collectives primitives performance our own micro-benchmark suite [29], similar to Intel MPI Benchmarks used for MPI libraries, has been used due to the lack of suitable micro-benchmarks for MPJ evaluation. Here, the results shown are the half of the round-trip time of a pingpong test or its corresponding bandwidth. The transferred data are byte arrays, avoiding the serialization overhead that would distort the analysis of the results.

Figures 10 and 11 show point-to-point latencies (for short messages) and bandwidths (for long messages) on InfiniBand and shared memory, respectively. The `ibvdev` middleware obtains significant point-to-point performance benefits, thus obtaining 11 μ s start-up latency and up to 7.2 Gbps bandwidth. The threshold between eager and rendezvous send protocols can be observed in the bandwidth graph at 128 KB, which confirms the efficiency of the implementation of the zero-copy rendezvous protocol with RDMA Write for `ibvdev`. These results outperforms significantly `niodev` over InfiniBand, limited to 65 μ s start-up latency and below 3 Gbps bandwidth.

Compared to native MPI libraries, `ibvdev` obtains a similar bandwidth than MVAPICH (7 Gbps) in this testbed, surpassing it even at several points (e.g., 32 KB, 256 KB and 512 KB message sizes). Nevertheless, OpenMPI shows the best performance from 32 KB message size, obtaining up to 9.2 Gbps bandwidth. As for latency, `ibvdev` obtains better results than MVAPICH (13 μ s) and only slightly worse from OpenMPI (10 μ s), again the best performer.

Regarding shared memory communication performance, `ibvdev` obtains much better start-up latency, 6 μ s, than the multithreading `smpdev` middleware, which achieves 17 μ s, which means that `ibvdev` has implemented a highly

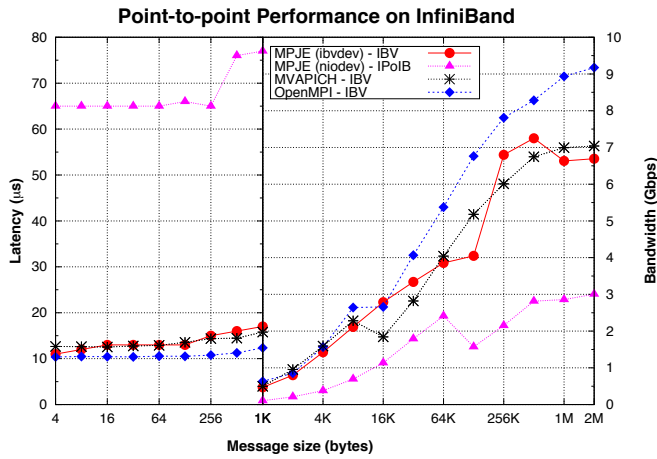


Fig. 10 Message-passing point-to-point performance on InfiniBand

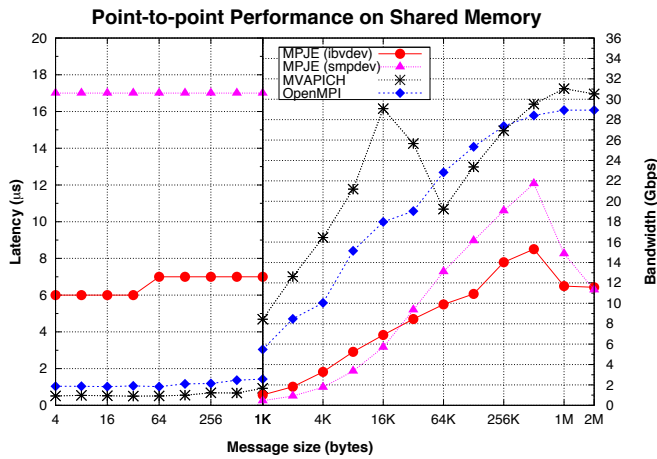


Fig. 11 Message-passing point-to-point performance on shared memory

efficient communication protocol and that `smpdev` presents poor start-up latency, caused by an excess of synchronizations. The native MPI libraries are again the best performers obtaining $0.5 \mu\text{s}$ and $1 \mu\text{s}$ for `MVAPICH` and `OpenMPI`, respectively, due to their efficient communications support on shared memory. Regarding bandwidth, MPJ devices are far from native MPI libraries, obtaining worse performance (15.3 Gbps and 22 Gbps for `ibvdev` and `smpdev`, respectively).

5.3 Collective Primitives Micro-benchmarking

Figure 12 presents the aggregated bandwidth for representative MPJ data movement operations (broadcast and allgather), and computational operations (reduce and allreduce double precision sum operations) with 128 processes. The aggregated bandwidth metric has been selected as it takes into account the global amount of data transferred. The `niodev` allgather results could not be taken due to flaws in the implementation that hanged its operation. In addition to `ibvdev`, `niodev` and MPI communications it has been evaluated the performance of multithreaded versions of the MPJ collective operations, running only one process per node, and 16 threads within each process. Thus, instead of running 128 processes on the cluster, only 8 processes are being used, taking advantage of intra-node communications through multithreading. This hybrid support of network and multithreading communications is one of the main advantages of Java middleware for scalable and efficient communication on clusters of multi-core processors.

The results confirm that `ibvdev` outperforms significantly `niodev`, achieving up to one order of magnitude higher performance, although generally the performance benefit is 2 or 3 times better. Moreover, both `ibvdev` and `niodev` take advantage of the multithreaded collectives. With respect to the MPI libraries, `ibvdev` achieves better performance than MPI collectives for short messages, up to 16 KB - 256 KB, thanks to the exploitation of multithreading in collectives implementation and the use of a high PSL (128 KB), whereas MPI libraries use smaller PSL (8 KB). However, for longer messages the MPI collectives achieve much better performance due to the use of better collective algorithms, and the use of pipelined transfers.

5.4 Kernel/Application Performance Analysis

The impact of `ibvdev` on the scalability of Java parallel codes has been analyzed using the NAS Parallel Benchmarks (NPB) implementation for MPJ (NPB-MPJ) [30], selected as the NPB are probably the benchmarks most commonly used in the evaluation of languages, libraries and middleware for HPC. In fact, there are implementations of the NPB for MPI, OpenMP and hybrid MPI/OpenMP.

Four representative NPB codes have been evaluated: CG (Conjugate Gradient), FT (Fourier Transform), IS (Integer Sort) and MG (Multi-Grid). Moreover, the `jGadget` [31] cosmology simulation application has also been analyzed. These MPJ codes have been selected as they show very poor scalability with MPJ Express over InfiniBand. Hence, these are target codes for the evaluation of the impact on performance of the use of `ibvdev` in MPJ Express. The results have been obtained using up to 64 processes instead of 128, due to memory constraints on the cluster.

Figure 13 shows the NPB-MPJ CG, IS, FT and MG results, respectively, for the Class C workload in terms of MOPS (Millions of Operations Per Sec-

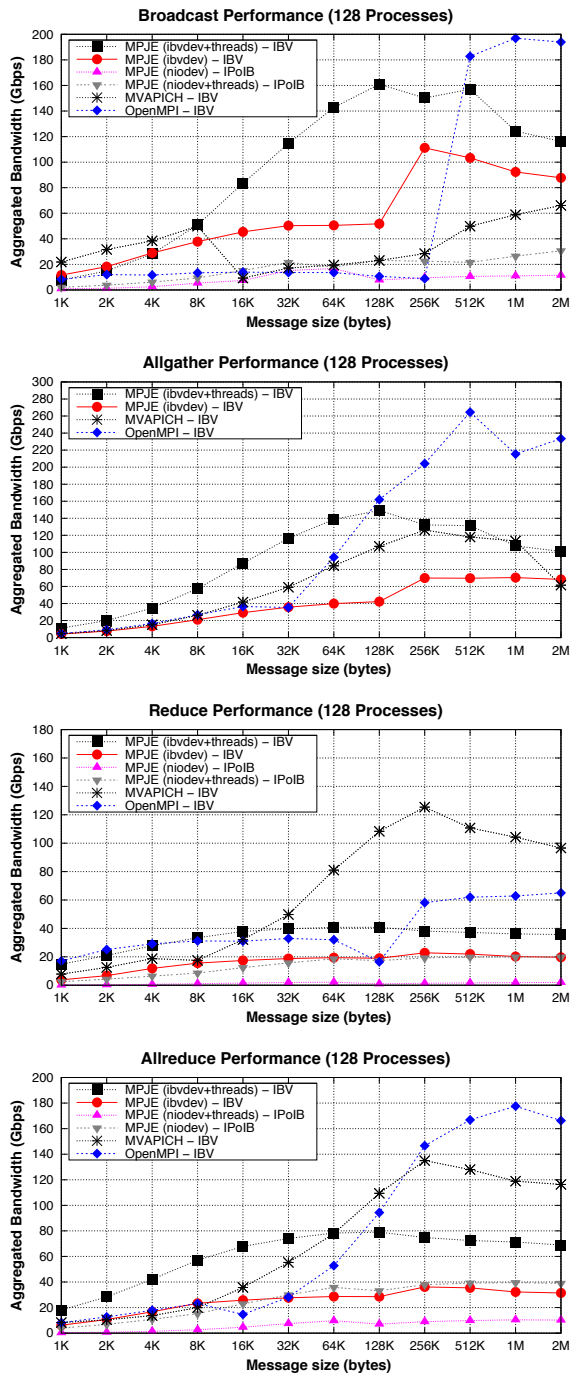


Fig. 12 Message-passing collective primitives performance

1 ond) (left) and its corresponding scalability, in terms of speedup (right). For
2 CG kernel, `ibvdev` doubles the performance of the `niodev` device over In-
3 finiBand, with almost 9000 MOPS compared to less than 4000 MOPS on 64
4 processes. With respect to IS kernel, the results for `niodev` over InfiniBand
5 show a significant slowdown with 64 processes, not taking advantage of the
6 use of 64 processes, while `ibvdev` keeps on scaling and gets up to 650 MOPS,
7 significantly outperforming the `niodev` results. Regarding FT, `ibvdev` also
8 doubles the performance of the `niodev` device over InfiniBand, with around
9 17000 MOPS compared to less than 8000 MOPS. Finally, the impact of `ibvdev`
10 on MG is smaller than for the remaining codes as this NPB is less communica-
11 tion intensive, as obtains relatively good speedups, even with `niodev` (speedup
12 of 30 with 64 processes).
13

14 The performance comparison of `ibvdev` against MPI libraries has two dif-
15 ferent analysis, depending on the metric used. If we take into account the
16 MOPS achieved, MPI benchmarks obtain always the best performance, around
17 a 50% higher than `ibvdev` results. The poorer performance of NPB-MPJ can
18 be attributed to the lower performance of the JVM compared to native com-
19 pilers. However, if we have a look at the speedups, `ibvdev` outperforms MPI
20 for FT and MG, while obtains slightly lower scalability for CG and IS, which
21 suggests that `ibvdev` implements a highly efficient communication support,
22 even comparable to MPI libraries, and that the use of efficient communication
23 libraries can bridge the gap between Java and natively compiled languages
24 provided that an efficient communication support is made available.
25

26 The jGadget application is the MPJ implementation of Gadget [32], a
27 popular cosmology simulation code initially implemented in C and parallelized
28 using MPI that is used to study a large variety of problems like colliding and
29 merging galaxies or the formation of large-scale structures. This application
30 has been selected for the performance evaluation of `ibvdev`, measuring its
31 performance using up to 64 processes instead of 128, due to memory constraints
32 on the cluster (each Java process is using its own JVM).
33

34 Figure 14 presents the performance results of jGadget running a two mil-
35 lion particles cluster formation simulation. As jGadget is a communication-
36 intensive application, with important collective operations overhead, only mod-
37 est speedups are obtained. Here `ibvdev` can take advantage of the use of 64 pro-
38 cesses (speedup above 22), whereas `niodev` over IPoIB remains with a speedup
39 of 16. Regarding MPI results, OpenMPI and MVAPICH achieve around 45%
40 higher speedup than `ibvdev` on 64 processes, which suggests that this middle-
41 ware is bridging the gap between Java and natively compiled applications in
42 HPC.
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

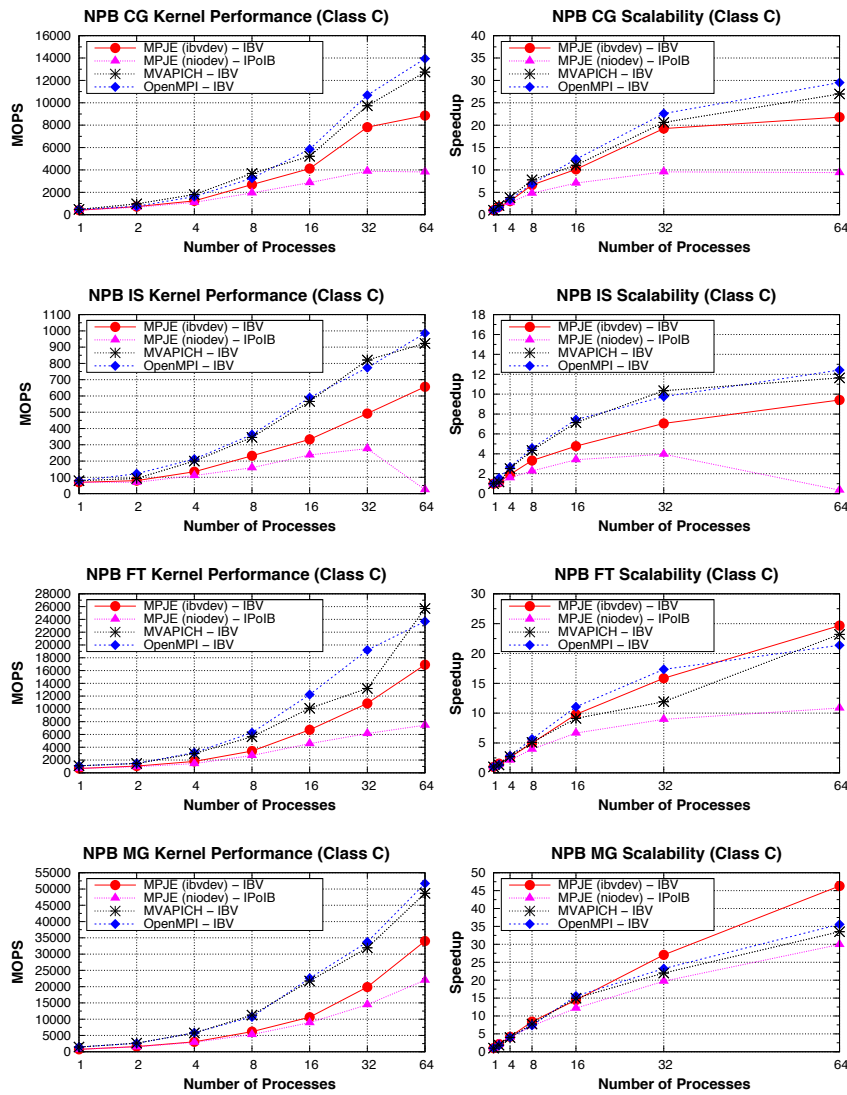


Fig. 13 Performance and scalability of NPB-MPI/MPJ codes

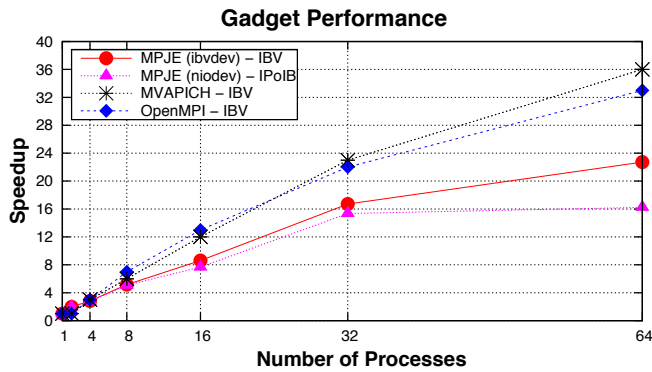


Fig. 14 Scalability of MPI/MPJ Gadget

6 Conclusions

This paper has presented `ibvdev`, a scalable and efficient low-level Java message-passing device for communication on InfiniBand systems. The increase in the number of cores per system demands languages with built-in multithreading and networking support, such as Java, as well as scalable and efficient communication middleware that can take advantage of multi-core systems. The developed device transparently provides Java message-passing applications with efficient performance on InfiniBand thanks to its direct support on IBV and the efficient and scalable implementation of a lightweight communication protocol which is able to take advantage of RDMA over InfiniBand.

The performance evaluation of `ibvdev` on an InfiniBand multi-core cluster has shown that this middleware obtains start-up latencies and bandwidths similar to MPI performance results, obtaining in fact up to 85% start-up latency reduction and twice the bandwidth compared to previous Java middleware on InfiniBand. Additionally, the impact of `ibvdev` on message-passing collective operations is significant, achieving up to one order of magnitude performance increases compared to previous Java solutions, especially when taking advantage of shared memory intra-process (multithreading) communication. The analysis of the impact of the use of `ibvdev` on MPJ applications shows significant performance increase compare to sockets-based middleware (`niodev`), which helps to bridge the gap between Java and natively compiled codes in HPC. To sum up, the efficiency of this middleware, which is even competitive with MPI point-to-point transfers, increments the scalability of communications intensive Java applications, especially in combination with the native multithreading support of Java.

Acknowledgements This work was funded by the Ministry of Science and Innovation of Spain under Project TIN2010-16735, and by the Xunta de Galicia under the Consolidation Program of Competitive Research Groups and Galician Network of High Performance Computing.

References

1. G.L. Taboada, J. Touriño, and R. Doallo. Java for High Performance Computing: Assessment of Current Research and Practice. In *Proc. 7th Intl. Conf. on Principles and Practice of Programming in Java (PPPJ'09)*, pages 30–39, Calgary, Canada, 2009.
2. B. Blount and S. Chatterjee. An Evaluation of Java for Numerical Computing. *Scientific Programming*, 7(2):97–110, 1999.
3. A. Shafi, B. Carpenter, and M. Baker. Nested Parallelism for Multi-core HPC Systems using Java. *Journal of Parallel and Distributed Computing*, 69(6):532–545, 2009.
4. InfiniBand Trade Association. Infiniband Architecture Specification Volume 1, Release 1.2.1, 2004. <http://www.infinibandta.org/> [Last visited: April 2011].
5. M. Baker, B. Carpenter, and A. Shafi. A Pluggable Architecture for High-Performance Java Messaging. *IEEE Distributed Systems Online*, 6(10):1–4, 2005.
6. IETF Draft. IP over IB. <http://www.ietf.org/old/2009/ids.by.wg/ipoib.html> [Last visited: April 2011].
7. Z. Hongwei, H. Wan, H. Jizhong, H. Jin, and Z. Lisheng. A Performance Study of Java Communication Stacks over InfiniBand and Gigabit Ethernet. In *Proc. 4th IFIP Intl. Conf. Network and Parallel Computing (NPC'07)*, pages 602–607, Dalian, China, 2007.
8. R. Veldema, R.F.H. Hofman, R. Bhoedjang, and H.E. Bal. Run-time Optimizations for a Java DSM Implementation. *Concurrency and Computation: Practice and Experience*, 15(3-5):299–316, 2003.
9. R.d.R. Righi, P.O.A. Navaux, M.C. Cera, and M. Pasin. Asynchronous Communication in Java over Infiniband and DECK. In *Proc. 17th Intl. Symp. on Computer Architecture and High Performance Computing (SBAC-PAD'05)*, pages 176–183, Rio de Janeiro, Brazil, 2005.
10. G.L. Taboada, J. Touriño, and R. Doallo. Java Fast Sockets: Enabling High-speed Java Communications on High Performance Clusters. *Computer Communications*, 31(17):4049–4059, 2008.
11. H. Wan, Z. Hongwei, H. Jin, H. Jizhong, and Z. Lisheng. Jdib: Java Applications Interface to Unshackle the Communication Capabilities of InfiniBand Networks. In *Proc. 4th IFIP Intl. Conf. Network and Parallel Computing (NPC'07)*, pages 596–601, Dalian, China, 2007.
12. H. Wan, H. Jizhong, H. Jin, Z. Lisheng, and L. Yao. Enabling RDMA Capability of InfiniBand Network for Java Applications. In *Proc. 4th IFIP Intl. Conf. on Networking, Architecture, and Storage (NAS'08)*, pages 187–188, Chongqing, China, 2008.
13. L. Yao, H. Jizhong, G. Jinjun, and H. Xubin. uStream: A User-Level Stream Protocol over Infiniband. In *Proc. 15th Intl. Conf. on Parallel and Distributed Systems (ICPADS'09)*, pages 65–71, Shenzhen, China, 2009.
14. D. Dunning, G. Regnier, G. McAlpine, D. Cameron, B. Shubert, F. Berry, A.M. Merritt, E. Gronke, and C. Dodd. The Virtual Interface Architecture. *IEEE Micro*, 18(2):66–76, 1998.
15. Compaq, Intel and Microsoft Corporations. The Virtual Interface Architecture Specification, Version 1.0, 1997.
16. OpenFabrics Alliance Website. <http://www.openfabrics.org/> [Last visited: April 2011].
17. B. Carpenter, G. Fox, S.-H. Ko, S. Lim. mpiJava 1.2: API specification. <http://www.hpjava.org/reports/mpiJava-spec/mpiJava-spec/mpiJava-spec.html> [Last visited: April 2011].
18. The mpiJava project. <http://www.hpjava.org/mpiJava.html> [Last visited: April 2011].
19. M. Bornemann, R. V. van Nieuwpoort, and T. Kielmann. MPJ/Ibis: a Flexible and Efficient Message Passing Platform for Java. In *Proc. 12th European PVM/MPI Users' Group Meeting (EuroPVM/MPI'05)*, pages 217–224, Sorrento, Italy, 2005.

20. G.L. Taboada, J. Touriño, and R. Doallo. F-MPJ: Scalable Java Message-passing Communications on Parallel Systems. *Journal of Supercomputing*, (In press).
21. Java Grande Forum. <http://www.javagrande.org> [Last visited: April 2011].
22. R. V. van Nieuwpoort, J. Maassen, G. Wrzesinska, R. Hofman, C. Jacobs, T. Kielmann, and H.E. Bal. Ibis: a Flexible and Efficient Java-based Grid Programming Environment. *Concurrency and Computation: Practice and Experience*, 17(7-8):1079–1107, 2005.
23. MX User's Guide. <http://www.myri.com/scs/MX/doc/mx.pdf> [Last visited: April 2011].
24. Open Source High Performance MPI Library. <http://www.open-mpi.org/> [Last visited: April 2011].
25. MPI over InfiniBand, 10GigE/iWARP and RDMAoE. <http://mvapich.cse.ohio-state.edu/> [Last visited: April 2011].
26. H. Tezuka, F. O'Carroll, A. Hori, and Y. Ishikaw. Pin-down Cache: A Virtual Memory Management Technique for Zero-copy Communication. In *Proc. 12th. Intl. Parallel Processing Symposium / 9th Symposium on Parallel and Distributed Processing (IPPS/SPDP'98)*, pages 308–314, Orlando, FL, USA, 1998.
27. M. Baker, B. Carpenter, and A. Shafi. A Buffering Layer to Support Derived Types and Proprietary Networks for Java HPC. *Scalable Computing Practice and Experience*, 8(4):343–358, 2007.
28. MPJ Express Website. <http://mpj-express.org/> [Last visited: April 2011].
29. G.L. Taboada, J. Touriño, and R. Doallo. Performance Analysis of Java Message-Passing Libraries on Fast Ethernet, Myrinet and SCI Clusters. In *Proc. 5th IEEE Intl. Conf. on Cluster Computing (CLUSTER'03)*, pages 118–126, Hong Kong, China, 2003.
30. D.A. Mallón, G.L. Taboada, J. Touriño, and R. Doallo. NPB-MPJ: NAS Parallel Benchmarks Implementation for Message-Passing in Java. In *Proc. 17th Euromicro Intl. Conf. on Parallel, Distributed, and Network-Based Processing (PDP'09)*, pages 181–190, Weimar, Germany, 2009.
31. M. Baker, B. Carpenter, and A. Shafi. MPJ Express meets Gadget: Towards a Java code for cosmological simulations. In *13th European PVM/MPI Users' Group Meeting (EuroPVM/MPI'06)*, pages 358–365, Bonn, Germany, 2006.
32. V. Springel. The Cosmological Simulation Code GADGET-2. *Monthly Notices of the Royal Astronomical Society*, 364(4):1105–1134, 2005.