

# iRCCE: A Non-blocking Communication Extension to the RCCE Communication Library for the Intel Single-Chip Cloud Computer

Carsten Clauss, Stefan Lankes, Jacek Galowicz, Thomas Bemmerl  
Chair for Operating Systems, RWTH Aachen University  
Kopernikusstr. 16, 52056 Aachen, Germany  
{clauss,lankes,galowicz,bemmerl}@lfbs.rwth-aachen.de

February 22, 2011

## 1 Introduction

The Single-Chip Cloud Computer (SCC) experimental processor [3] is a 48-core *concept vehicle* created by Intel Labs as a platform for many-core software research. The 48 cores are arranged in a 6x4 on-die mesh of tiles with two cores per tile. The SCC chip possesses four on-die memory controllers for addressing the external main memory. Additionally, each tile possesses a small amount of fast on-die memory that is also accessible to all other cores in a shared-memory manner. These special memory regions are the so-called *Message-Passing Buffers* (MPBs) of the SCC. The SCC's architecture does not provide any cache coherency between the cores, but rather offers a low-latency infrastructure in terms of these MPBs for explicit message-passing between the cores. Thus, the processor resembles a *Cluster-on-Chip* architecture with distributed but shared memory where each core can run its own operating system instance. Communication between processes hosted by different OS instances running on different cores can then be conducted either via a dedicated region of the off-die shared memory or via the fast on-die MPBs for increased performance. The RCCE communication library [4, 5] offers a customized message-passing interface for programming these SCC features by means of a simplified application programming interface (API). This message-passing environment, that is in turn based on a simpler one-sided communication mechanism (`RCCE.put/RCCE.get`), offers two-sided but *blocking* (often also referred to as *synchronous*) point-to-point communication functions (`RCCE.send/RCCE.recv`) as well as a set of collective communication operations (`RCCE.barrier`, `RCCE.bcast`, `RCCE.reduce`, ...). However, the lack of *non-blocking* point-to-point communication capabilities within the current RCCE library has driven us to extend RCCE by such asynchronous message-passing functions (`iRCCE.isend/iRCCE.irecv`, see Section 3). Furthermore, we have also improved the performance of some RCCE functions, as for example the blocking send and receive operations, by applying an assembler-coded and SCC-customized memory copy routine (`iRCCE.memcpy`, see Section 5). In order not to interfere with the current RCCE library and its future updates, we have placed our extensions into an additional auxiliary library with a separated namespace called *iRCCE*. In this manual, we detail the installation and the usage of these iRCCE extensions (see Section 2 and Section 8) and we present some performance comparisons between the current RCCE release and the improved iRCCE functions (see Section 7).

## 2 Getting Started

### 2.1 Installation Guide

1. Download, configure and build the common RCCE library with the *non-gory* interface.
2. Download<sup>1</sup> iRCCE as a TAR-File and unpack it within the desired installation folder.  
(e.g. within the RCCE root directory by calling `tar -xvzf iRCCE.tar.gz`)
3. Change into the iRCCE directory.
4. Type `./configure <path-to-rcce>` (= path to the RCCE installation).  
(e.g. `./configure ..` when you have unpacked iRCCE within the RCCE root directory)
5. Call `make`. All iRCCE related functions will then be added to the common RCCE library.<sup>2</sup>

### 2.2 Overview of the Basic iRCCE Functions

#### Library Initialization Function:

- `iRCCE_init();`

#### Non-Blocking Send and Receive Functions:<sup>3</sup>

- `int iRCCE_isend(char *, size_t, int, iRCCE_SEND_REQUEST *);`
- `int iRCCE_isend_test(iRCCE_SEND_REQUEST *, int *);`
- `int iRCCE_isend_wait(iRCCE_SEND_REQUEST *);`
- `int iRCCE_isend_push(void);`
- `int iRCCE_irecv(char *, size_t, int, iRCCE_RECV_REQUEST *);`
- `int iRCCE_irecv_test(iRCCE_RECV_REQUEST *, int *);`
- `int iRCCE_irecv_wait(iRCCE_RECV_REQUEST *);`
- `int iRCCE_irecv_push(void);`

#### Blocking but Improved Send and Receive Functions:

- `int iRCCE_send(char *, size_t, int);`
- `int iRCCE_recv(char *, size_t, int);`

#### Optimized Put and Get Functions:

- `int iRCCE_put(t_vcharp, t_vcharp, int, int);`
- `int iRCCE_get(t_vcharp, t_vcharp, int, int);`

---

<sup>1</sup>The iRCCE package can be found on the website of the Intel Many-core Application Research Community (MARC) [1]

<sup>2</sup>That in turn means that you do not have to link against iRCCE but just against the extended RCCE library.

<sup>3</sup>For convenience, the the following function names are mirrored into the common RCCE namespace, too.

## 2.3 Overview of the Source Tree of iRCCE

At the top level of the iRCCE directory structure, there are the following folders and files:

- `README.txt` gives a short introduction into iRCCE and its installation
- `Makefile.in` is a template for a later `Makefile` to be generated by the `configure` script
- `configure` should be called initially during the installation process with the path to the RCCE library as the first argument (see Section 2.1)
- `include` is a folder containing the `iRCCE.h` header file that must be included by all applications using iRCCE
- `src` is the folder that contains all iRCCE-related C source files
- `apps` is a folder that contains a simple iRCCE application example as well as some simple benchmark tools that are all described below
- `doc` is the folder containing this document as a PDF file as well as its  $\text{\LaTeX}$  sources

## 2.4 Description of the Application Example and the Benchmark Tools

In order to build the application example and the benchmark tools, just change into the `apps` folder after the installation of iRCCE and call `make all`. The following iRCCE executables will then be built:

- **pingpong** A simple Ping-Pong benchmark that utilizes the improved blocking send and receive functions of iRCCE. The benchmark reports *round-trip time* / 2 as well as *bandwidth* for ascending message sizes by measuring and averaging the time for  $n$  repetitions of the Ping-Pong pattern that is shown in Figure 1. Optional arguments that can be passed to the executable are the number  $n$  of repetitions (default is 10000) and the maximum of the messages sizes (default is 1024) that should be tested.
- **pingping** This benchmark performs a common variation of the Ping-Pong pattern that can only be realized by means of non-blocking communication functions. In contrast to the Ping-Pong benchmark, the time is measured under the aggravating circumstance that the outgoing message is interleaved with an incoming one [2]. The benchmark reports the pure average *ping time* for  $n$  repetitions for ascending message sizes. The difference between both patterns and between their clocked benchmark times can be made clear by comparing Figure 1 with Figure 2.
- **spam** This application enqueues a line of user-chosen length of requests and processes it afterwards. Hence, this is just a stress test for simulating a scenario with a huge amount of outstanding non-blocking communication requests.
- **madmonkey** Often cited in discussions about evolution, the infinite monkey theorem states that a monkey typing random keys on a typewriter will certainly write any given text after a (long) period time (Shakespeare's works are often mentioned here). While parallel computer architectures are very suitable for simulating this scenario, this application demonstrates the use of non-blocking send/receive calls in situations where no one can tell when sends/receives shall occur. This application needs a word to be found as program argument. Try for example "iRCCE"; the search for this word should not take much longer than one minute when running the program on two cores with one Master and one Slave process.

Each of these executables can just be started like a *normal* RCCE application. For example like:  
> `rcceun -nue 2 -f rc.hosts pingping 1000 65536`

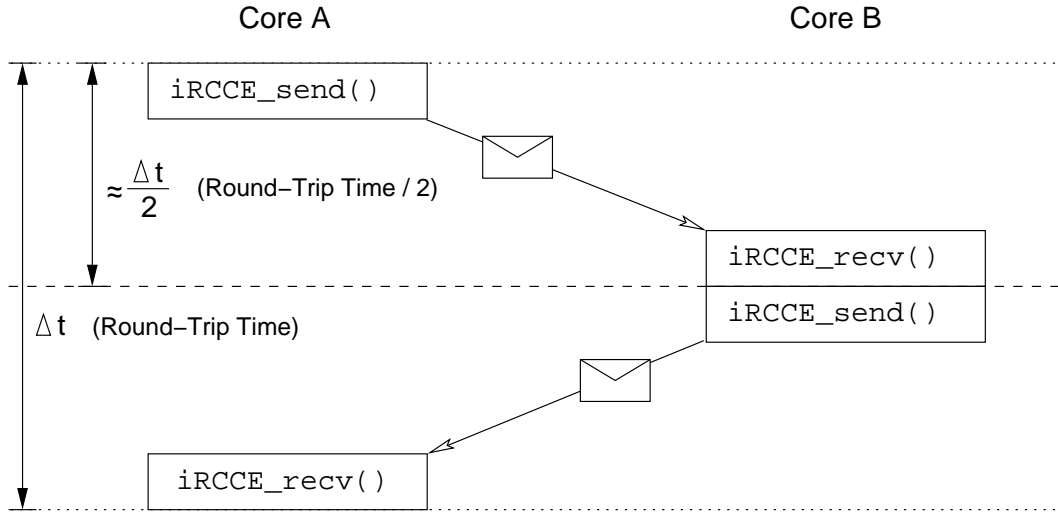


Figure 1: Communication Pattern of the iRCCE *Ping-Pong* Benchmark

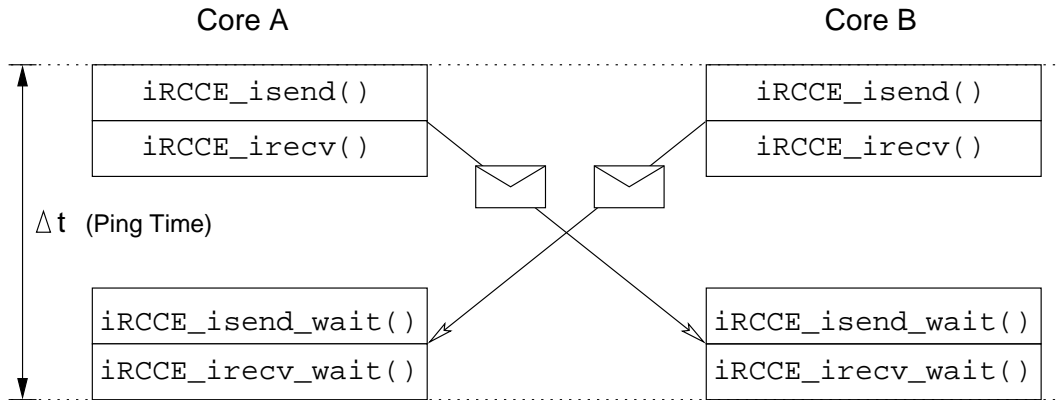


Figure 2: Communication Pattern of the iRCCE *Ping-Ping* Benchmark

## 3 Non-blocking Communication

### 3.1 Interleaved Communication and Computation

In many situations, performance can be improved by interleaving communication and computation. This is particularly true in situations where the communication progress stalls due to a temporary lack of communication resources like intermediate buffers. In such a situation, a blocking communication function has to wait (either actively by polling or passively by sleeping) until the needed resource becomes available again. An alternative mechanism in such a situation is to use non-blocking communication functions that do not block but return back to the application immediately if the communication progress can temporarily not be fostered. Of course, it is up to the application to exploit the interim time until the communication progress can be pushed on again. This can be done either by processing computational tasks or by pushing on with other communication requests. But this in turn means that the application must check repeatedly by itself whether the communication is still stuck or not. For this purpose, non-blocking communication functions usually pass back a so-call *request handle* which can then be used by additional *push*, *test* or *wait* functions to ensure the communication's progress and eventually its completion.

### 3.2 An Exemplary Scenario

For example, a non-blocking *receive function* will just check if the respective incoming message is already completely available. If this is not the case (e.g. the sender has not even started the message transfer), the function just records all needed parameters (like the source and the length of the expected message as well as the address of the respective receive buffer) within the request handle and returns it immediately to the application level. The program can then perform some other application-related calculations, provided that these are independent from the data of the pending message. However, during these calculations, the program has to call a specific *push function* repeatedly in order to ensure that the communication progress is being fostered in an interleaved manner. Afterwards, the program has to call an additional *test* or *wait function* in order to check whether the message transfer has been completed. This approach should be illustrated by the following pseudo-code:

```
# Initialize the non-blocking receive request:
CALL non_blocking_receive(OUT request_handle)

# Perform the interleaved computation and communication:
# (the calculation task can be divided into n subtasks)
FOR i = 1 TO n
    CALL partial_calculation(IN i)
    CALL push_communication(VOID)
NEXT i = i + 1

# Ensure that the communication has yet been completed:
CALL wait_for_communication_completion(IN request_handle)
```

Of course, not only the receiving of a message can be conducted in such a non-blocking and interleaved manner, but also the sending of a message can be processed in a similar way. The considerations for non-blocking send requests are quite the same as stated above for the receive request and even the programming patterns are quite analogous. Therefore, we do not present an exemplary non-blocking send scenario at this point.

### 3.3 Interleaved vs. Overlapped Communication

On systems, where the communication can be conducted autonomously by a communication controller (e.g. a DMA engine) or by a communication thread (e.g. on a multi-core CPU), not only an interleaving but rather a true *overlapping* of communication and computation can be achieved. The main difference between interleaved and overlapped communication is that interleaving implies a concurrent but still serialized processing, whereas a true overlapping results in a parallel processing of communication and computation. Thus, the first approach helps to hide wait states and to break up message dependencies (see Section 3.5 for an example), but on the other hand it requires an explicit switch-over between computation and communication by frequently calling the push function. In contrast to this, the second approach can achieve real parallelism and thus does not need an explicit pushing for progress.<sup>4</sup> However, since there is no other asynchronous hardware (like a communication controller) available for the cores on the SCC, the iRCCE library (currently) just implements the first approach.

### 3.4 The Non-Blocking Communication Approach of iRCCE

As stated above, the iRCCE library implements an interleaving mechanism for non-blocking communication operations. This is achieved by using the standard RCCE communication functions<sup>5</sup> as a template for their iRCCE counterparts, but instead of waiting for a progress flag<sup>6</sup> to be set, the non-blocking iRCCE functions just test the respective flag<sup>7</sup> and return immediately if a progress cannot be made right away. In order to resume the communication later on, the current progress states are stored in request handles that are of type `iRCCE_SEND/RECV_REQUEST`<sup>8</sup>. These handles are set up and returned by the respective non-blocking send/receive functions (`iRCCE_isend()/iRCCE_irecv()`) that have to be called in order to initiate a non-blocking communication request. Subsequently, the communication progress can be pushed on by calling `iRCCE_isend/irecv.push()` repeatedly. Again, also these functions return immediately if the communication progress can temporarily not be fostered. Finally, the completion of a once pending communication request must be ensured by a call to the `iRCCE_isend/irecv_test()` or to the `iRCCE_isend/irecv_wait()` function with the respective request handle as function parameter.

**Before the completion of a non-blocking operation is not ensured by a call to these functions, neither the respective receive buffer is guaranteed to be valid (it is likely that the message has yet not arrived in the receive buffer) nor the respective send buffer is allowed to be modified (it is likely that the message has yet not been copied out of the send buffer).**

In order to handle multiple outstanding communication requests, the iRCCE library implements a queuing mechanism. This is necessary, because it is possible to initiate subsequent non-blocking send requests, for example, to the same destination even before the first message gets started to be transferred. Thus, without a queuing mechanism, the order of the messages could not be observed. The iRCCE library implements this queuing mechanism in terms of single-linked lists with *one*<sup>9</sup> send and 48 receive queues per core. If the first element in such a list is not NULL (that means that a previous request is still in progress) the new request is just added at the end of the list. Otherwise, the `iRCCE_isend/irecv()` function tries to complete the new request directly. If such a completion cannot be achieved, the request becomes the head of the respective list and the function returns immediately.

---

<sup>4</sup>Nevertheless, it should be mentioned that the synchronization between a communication and a computation thread, as well as the latency for setting up a communication controller like a DMA engine, can expose an enormous overhead that can even nullify the performance improvements gained by the parallel progress, especially for short messages.

<sup>5</sup>These are the send and receive functions (`RCCE_send()/RCCE_recv()`) of the *non-gory* interface of RCCE.

<sup>6</sup>See the `RCCE_wait_until()` calls of `RCCE_send/recv_general()` in `RCCE_send/recv.c`.

<sup>7</sup>See the `iRCCE_test_flag()` calls of `iRCCE_push_send/recv_request()` in `iRCCE_isend/irecv.c`.

<sup>8</sup>These are just C structs that are defined in `iRCCE.h`

<sup>9</sup>Consequently, send requests (even to different remote ranks) are always processed in the order of their respective `iRCCE_isend()` calls; whereas receive requests are handled in the order of the message arrival.

According to this, if one wants to foster the communication progress irrespectively from a specific request, just the *first* pending request in the respective queue needs to be pushed; and that is exactly what `iRCCE_isend/irecv_push()` does. In contrast to this, `iRCCE_isend/irecv_test()` only checks and pushes that request that is passed as the function argument. However, these test functions can alternatively be called with `NULL` as the argument indicating that the respective queue as a whole is meant. And likewise, when calling `iRCCE_isend/irecv_wait()` with `NULL` as the argument, the completion of the whole pending queue is waited for. Furthermore, even `iRCCE_isend/irecv()` can be called with `NULL` as the request argument. In such a case, a subsequent `wait()` call will be issued internally.

### 3.5 An Example Code: The Ping-Ping Pattern

In this section, we want to detail the communication kernel of the Ping-Ping benchmark that is part of the iRCCE distribution (see Section 2.4). The Ping-Ping pattern (see Figure 2) is an example for common communication patterns where messages are exchanged in a *symmetric* manner. In contrast to the notorious Ping-Pong pattern (see Figure 1), the Ping-Ping pattern is symmetric in this respect that both participating processes do exactly the same: (1) initiate a send request, (2) initiate a receive request and then (3) wait for their completion:

```
iRCCE_SEND_REQUEST send_request;
iRCCE_RECV_REQUEST recv_request;

char send_buffer[length];
char recv_buffer[length];

int my_rank      = RCCE_ue();
int remote_rank = (my_rank + 1) % 2;

RCCE_barrier(&RCCE_COMM_WORLD);

timer = RCCE_wtime();

for(round=0; round < numrounds+1; round++)
{
    /* (1) send PING via non-blocking send: */
    iRCCE_isend(send_buffer, length, remote_rank, &send_request);

    /* (2) receive PING via non-blocking recv: */
    iRCCE_irecv(recv_buffer, length, remote_rank, &recv_request);

    /* (3) wait for completion: */
    iRCCE_isend_wait(&send_request);
    iRCCE_irecv_wait(&recv_request);
}

timer = RCCE_wtime() - timer;
```

As one can see, without non-blocking communication functions, this pattern could not be realized because the symmetric blocking send calls would both stuck in anticipation of the matching but subsequent receive calls.

## 4 Extended Range of Functions

Besides the yet presented iRCCE functions for handling non-blocking communication, we have also added some more higher-level functions for a more convenient handling of outstanding non-blocking requests. This additional part of the iRCCE API should be detailed in this section.

### 4.1 Functions for Canceling Requests

First of all, we have added functions for canceling already enqueued but not yet started send and receive requests (`iRCCE_isend_cancel()`/`iRCCE_irecv_cancel()`). By means of these functions, it can be *attempted* to remove such a request from the respective waiting queue until it becomes the head of it. However, if the request has already become the first and the actual communication has already been started, a subsequent canceling is no longer possible. In such a case, the cancel function returns the information that the requested withdrawal has failed. But this in turn means that the started request has to be matched by a respective send or receive call on the remote side.<sup>10</sup> Therefore, the application programmer should be clear in one's mind about the fact that a cancel call may fail and that the application's communication pattern must be designed in such a way that all started non-blocking requests are resolved.

### 4.2 Functions for Handling Multiple Outstanding Requests

In many cases one wants request some amount of non-blocking send/receive operations, then do calculations in some kind of loop and from time to time test if the send/receive operations have succeeded. For bigger amounts of messages one will certainly write a loop for testing all those requests. Those loops are an unnecessary and uncomfortable duplication of code, so we wrote functions to handle this in a more elegant manner. `iRCCE_test_all()` will traverse a linked list of send/receive requests and return success if there have been no pending requests, while `iRCCE_wait_all()` is not returning until there are no pending requests left. One only have to initialize a data structure of type `iRCCE_WAIT_LIST` and add requests to it first. Analogous to this, `iRCCE_test_any()`/`iRCCE_wait_any()` can be used for testing or waiting for completion of *any* pending operation in the list.

### 4.3 Overview of the Additional iRCCE Functions

- `int iRCCE_isend_cancel(iRCCE_SEND_REQUEST*, int*);`
- `int iRCCE_irecv_cancel(iRCCE_RECV_REQUEST*, int*);`
- `void iRCCE_init_wait_list(iRCCE_WAIT_LIST*);`
- `void iRCCE_add_to_wait_list(iRCCE_WAIT_LIST*, iRCCE_SEND_REQUEST*, iRCCE_RECV_REQUEST*);`
- `int iRCCE_test_all(iRCCE_WAIT_LIST*, int*);`
- `int iRCCE_wait_all(iRCCE_WAIT_LIST*);`
- `int iRCCE_test_any(iRCCE_WAIT_LIST*, iRCCE_SEND_REQUEST**, iRCCE_RECV_REQUEST**);`
- `int iRCCE_wait_any(iRCCE_WAIT_LIST*, iRCCE_SEND_REQUEST**, iRCCE_RECV_REQUEST**);`

---

<sup>10</sup>This is because otherwise the pending request will wait eternally at the head of the queue and will hinder all following requests from being processed.



## 5 Other Improvements of iRCCE to RCCE

### 5.1 SCC-optimized Memory Copy Functions

On the current SCC architecture, a write access does *not* perform a cache line fill even if the cache line is not present (*cache on read not on write*). Therefore, the core writes in this case directly to the main memory and, of course, this is quite expensive in terms of time. However, if a present cache line is modified, these changes are not directly applied to main memory but just to the cache. Thus, when touching a cache line by a read access before modifying its content, subsequent write accesses to this line do not imply a write through to the main memory. The following SCC-customized functions copy memory from the on-die buffers (MPB) to an off-die region while pre-fetching the cache lines of the destination: `iRCCE_put()` / `iRCCE_memcpy_put()`. Therefore, these functions avoid the above mentioned SCC-specific bad behavior of write misses. In turn, if the destination is located in on-die memory (MPB), classical pre-fetching techniques<sup>11</sup> are used by the following functions in order to increase the copy performance: `iRCCE_get()` / `iRCCE_memcpy_get()`. All these functions are internally used by the iRCCE communication functions; but they can also be utilized outside of the library.

### 5.2 Pipelined Send and Receive Functions

The common RCCE send and receive functions<sup>12</sup> use the core-local MPB space for sending messages to remote cores. As one may know, these local buffers are 8KByte MPB space per core. If a message to be sent is bigger than the available local MPB space, it must be divided into chunks that are then sent piecewise. For this purpose, the common RCCE functions determine a chunk size that is equal to the amount of the available local MPB space. However, when using the whole available local MPB space for one single message chunk, the message transport becomes necessarily a serialized process: (1a) sender puts chunk into the MPB, (1b) sender signalizes the chunk's arrival, (2a) receiver gets chunk from the MPB, (2b) receiver signalizes the chunk's removal (see Figure 3). This processing scheme must then be performed iteratively until the whole message has been transferred.

Therefore, a smarter approach is not to use the whole local MPB as one big chunk but to divide it into two smaller chunks. This is because in this case sender and receiver can work on the MPB simultaneously in a pipelined and parallelized manner: (1a) sender puts message chunk A into the MPB, (1b) sender signalizes the arrival of chunk A, (2a) sender puts message chunk B into the MPB; and *meanwhile*, the receiver can remove chunk A from the MPB, (2b) sender/receiver signalize the arrival/removal of chunk B/A; and so on (see Figure 4). The differences between both approaches can be made clear by comparing Figure 3 and Figure 4.

This pipelining approach is implemented within iRCCE in terms of the blocking `iRCCE_send()` and `iRCCE_recv()` functions. It is important that each call of these functions is matched with a call of its respective counterpart on the remote side. That means that a message that is sent via `iRCCE_send()` must necessarily be received via `iRCCE_recv()`; and vice versa. This is because the pipelining requires that both the sender and the receiver cooperate *synchronously* according to the pipelining technique. That this technique can really help to improve the communication performance especially for larger messages can be proved by applying appropriate benchmark tool, as for example shown here in Section 7. However, one should remember that this technique (at least in its current implementation) just takes effect for messages that are bigger than the local MPB space of 8KByte.

---

<sup>11</sup>See `include/scc_memcpy.h` for more details.

<sup>12</sup>as well as the non-blocking and the blocking send and receive functions of iRCCE

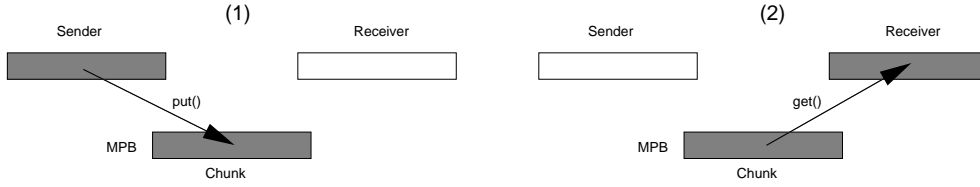


Figure 3: Serialized Transport of Message Chunks through the MPB Space

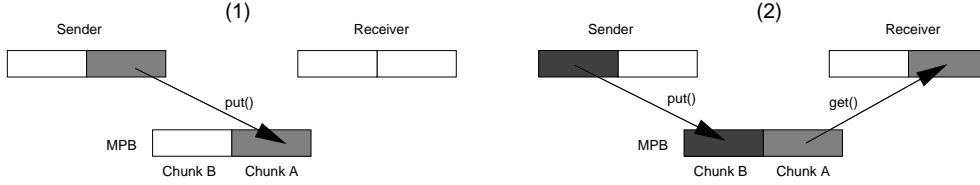


Figure 4: Pipelined Transport of Message Chunks through the MPB Space

## 6 Known Problems and Issues

- **Differences between MPI and iRCCE Semantics**

Although some function names of iRCCE are quite similar to their counterparts of the Message Passing Interface Standard (MPI) [6], usage and semantics of both APIs differ in detail. So, for example, the test function of iRCCE only checks and pushes that non-blocking request that is passed via the respective function call. In contrast to this, the MPI test function also triggers the so-called progress engine that ensures that all waiting requests are checked for progress irrespectively from the passed request of the function call.

- **Matching of Non-Blocking Requests with Blocking Function Calls**

Although not extensively tested, a call e.g. of the blocking RCCE send function (`RCCE_send()`) should match with a call of its non-blocking iRCCE counterpart (`iRCCE_irecv()`) on the remote side; and vice versa. *But beware!* This is true for matching non-blocking iRCCE calls with blocking RCCE calls, but not for mixing non-blocking calls with the blocking but *pipelined* functions of iRCCE, as for example `iRCCE_recv()`.<sup>13</sup> Moreover, calling a blocking RCCE function after initiating a still pending non-blocking iRCCE function can cause deadlocks because of the missing `push()` calls!<sup>14</sup>

- **Mixing of Non-Blocking Function Calls with Collective Operations**

Do not use collective communication operations (like `RCCE_bcast` or `RCCE_barrier`) when there are still outstanding non-blocking communication requests! This is because since RCCE does not use message tags like MPI, iRCCE cannot distinguish internally between collective and point-to-point requests. Thus, an overlapping of non-blocking transfers with collective communication patterns can lead to message mismatches and deadlocks.

- **Allocating MPB Space or Flags during Non-Blocking Communication**

Do not allocate new flags or new MPB regions when there are still outstanding non-blocking communication requests! This is because the amount and the position of the MPB space used for the asynchronous data transfers is recorded at creation time of the respective non-blocking communication request and internally used throughout the request's completion.

<sup>13</sup>This is because the pipelining requires that both the sender and the receiver cooperate *synchronously* according to the pipelining technique (see Section 5.2).

<sup>14</sup>Better use `iRCCE_isend/irecv()` with `NULL` as the request argument in this case because the internal `wait()` call avoids such deadlocks (see Section 3.4).

## 7 Performance Results and Comparisons

In this section, we want to present some performance results. In Figure 5, one can see the Ping-Pong bandwidth<sup>15</sup> (measure with the iRCCE Ping-Pong benchmark, see Section 2.4) for different message sizes and different optimization approaches. All measurements were done between the cores rck00 and rck01 with network and memory running at 800MHz and a core frequency of 533MHz (Tile533\_Mesh800\_DDR800). The used libraries were standard RCCE (V 1.0.13, `big_flags`, `nongory`) and our iRCCE. The figure shows a performance comparison between the *standard RCCE* functions (`RCCE_send/recv()`), the utilization of *improved memory copy* functions (`iRCCE_memcpy_put/get()`) and the the applying of *Pipelining* (`iRCCE_send/recv()`) for long messages (*length* > 8192). As one can see, iRCCE outperforms RCCE with respect to the communication bandwidth especially in case of larger messages. However, the average latency for a 1Byte message still remains the same since the improved memory copy function does not take effect for messages smaller than a cache line.

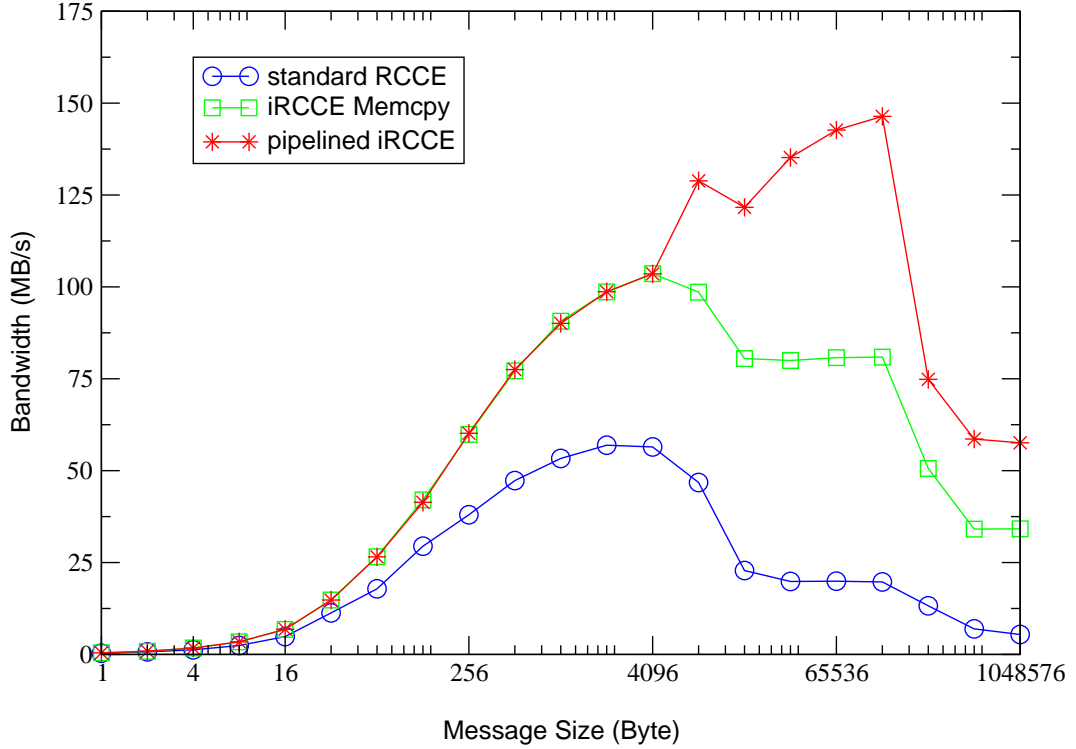


Figure 5: Some Results of the *iRCCE Ping-Pong Benchmark*: Performance comparison in terms of communication bandwidth between the *standard RCCE* functions (`RCCE_send/recv()`), the utilization of *improved Memory-Copy* functions (`iRCCE_memcpy_put/get()`) and the the applying of *Pipelining* (`iRCCE_send/recv()`) for long messages (*length* > 8192)

<sup>15</sup>See [5] for a detailed analysis of the bandwidth curve progression and its correlation with cache level sizes.

## 8 The iRCCE Application Programming Interface

### 8.1 Library Initialization Function

#### **int iRCCE\_init(void)**

This function must be invoked before any other iRCCE function is called. The function initializes crucial data structures of the iRCCE library as for example the queues for pending communication requests. The return value should always be `iRCCE_SUCCESS`.

### 8.2 Functions for Non-Blocking Sending

#### **int iRCCE\_isend(char \*buffer, size\_t length, int dest, iRCCE\_SEND\_REQUEST \*request)**

This function initiates a request for a non-blocking send operation. The function initializes a request handle of type `iRCCE_SEND_REQUEST` that must be passed as a pointer and that is used later on for checking for the operation's completion. The function returns `iRCCE_SUCCESS` in that case that the request could be finished already within this function call. However, usually the function returns `iRCCE_PENDING` or `iRCCE_RESERVED` indicating that the communication has been started but not yet finished or that there are prior requests pending in the send queue and that the request is being reserved. For more details about this function see Section 3.

<b>buffer</b>	starting address of the message to be sent
<b>length</b>	length of the outgoing message in bytes
<b>dest</b>	rank (ID) of the target process (UE)
<b>request</b>	request handle that is used later on for checking for completion (if set to <code>NULL</code> , a subsequent <code>iRCCE_isend_wait()</code> call will be issued internally)

#### **int iRCCE\_isend\_test(iRCCE\_SEND\_REQUEST \*request, int \*flag)**

This function checks whether an initiated request for a non-blocking send operation is already completed. The function checks and pushes only that request that is passed as the request handle argument. However, if `NULL` is passed as the argument, the function checks for the completion of the whole pending send queue. The function is non-blocking and the return values and their meanings are the same as for `iRCCE_isend()`. Therefore, the argument `flag` can also be omitted by passing `NULL` instead.

<b>request</b>	request handle returned by <code>iRCCE_isend()</code>
<b>flag</b>	flag that indicates whether the message has been sent (1) or not (0)

#### **int iRCCE\_isend\_wait(iRCCE\_SEND\_REQUEST \*request)**

This function can be used for waiting for the completion of a pending send request. Since this function blocks until the respective request is finished, the function also pushes the pending receive queue as well as the pending send requests that are enqueued prior to that request in order to avoid deadlocks. However, if `NULL` is passed as the request handle argument, the completion of the whole pending send queue is waited for. The return value should always be `iRCCE_SUCCESS`.

<b>request</b>	request handle returned by <code>iRCCE_isend()</code>
----------------	---

### **int iRCCE\_isend\_push(void)**

This function pushes the progress of non-blocking communication requests that are enqueued in the pending send queue. The function is non-blocking and the return value is `iRCCE_PENDING` in case of still pending send requests or `iRCCE_SUCCESS` if the send queue is empty.

## **8.3 Functions for Non-Blocking Receiving**

### **int iRCCE\_irecv(char \*buffer, size\_t length, int source, iRCCE\_RECV\_REQUEST \*request)**

This function initiates a request for a non-blocking receive operation. The function initializes a request handle of type `iRCCE_RECV_REQUEST` that must be passed as a pointer and that is used later on for checking for the operation's completion. The function returns `iRCCE_SUCCESS` in that case that the request could be finished already within this function call. However, usually the function returns `iRCCE_PENDING` or `iRCCE_RESERVED` indicating that the communication has been started but not yet finished or that there are prior requests pending in the receive queue and that the request is being reserved. For more details about this function see Section 3.

<code>buffer</code>	starting address of the receive buffer
<code>length</code>	length of the expected message in bytes
<code>source</code>	rank (ID) of the source process (UE)
<code>request</code>	request handle (if set to <code>NULL</code> , <code>iRCCE_irecv_wait()</code> will be called internally)

### **int iRCCE\_irecv\_test(iRCCE\_RECV\_REQUEST \*request, int \*flag)**

This function checks whether an initiated request for a non-blocking receive operation is already completed. The function checks and pushes only that request that is passed as the request handle argument. However, if `NULL` is passed as the argument, the function checks for the completion of the whole pending receive queue. The function is non-blocking and the return values and their meanings are the same as for `iRCCE_irecv()`. Therefore, the argument `flag` can also be omitted by passing `NULL` instead.

<code>request</code>	request handle returned by <code>iRCCE_irecv()</code>
<code>flag</code>	flag that indicates whether the message has been received (1) or not (0)

### **int iRCCE\_irecv\_wait(iRCCE\_RECV\_REQUEST \*request)**

This function can be used for waiting for the completion of a pending receive request. Since this function blocks until the respective request is finished, the function also pushes the pending send queue as well as the pending receive requests that are enqueued prior to that request in order to avoid deadlocks. However, if `NULL` is passed as the request handle argument, the completion of the whole pending receive queue is waited for. The return value should always be `iRCCE_SUCCESS`.

<code>request</code>	request handle returned by <code>iRCCE_irecv()</code>
----------------------	---

### **int iRCCE\_irecv\_push(void)**

This function pushes the progress of non-blocking communication requests that are enqueued in the pending receive queues. The function is non-blocking and the return value is `iRCCE_PENDING` in case of still pending receive requests or `iRCCE_SUCCESS` if the receive queue is empty.

## 8.4 Blocking but Pipelined Communication Functions

### **int iRCCE\_send(char \*buffer, size\_t length, int dest)**

This function is quite similar to the blocking `RCCE_send()` function of RCCE. The main difference is that this function uses an SCC-optimized memory copy routine and that a pipeline technique for larger messages is used (see also Section 5 and Section 7). The threshold value for the message size, when pipelining should take effect, is 8KByte per default and the function call must be matched by a remote call of `iRCCE_recv()`.

<b>buffer</b>	starting address of the message to be sent
<b>length</b>	length of the outgoing message in bytes
<b>dest</b>	rank (ID) of the target process (UE)

### **int iRCCE\_recv(char \*buffer, size\_t length, int source)**

This function is quite similar to the blocking `RCCE_recv()` function of RCCE. The main difference is that this function uses an SCC-optimized memory copy routine and that a pipeline technique for larger messages is used (see also Section 5 and Section 7). The threshold value for the message size, when pipelining should take effect, is 8KByte per default and the function call must be matched by a remote call of `iRCCE_send()`.

<b>buffer</b>	starting address of the receive buffer
<b>length</b>	length of the expected message in bytes
<b>source</b>	rank (ID) of the source process (UE)

## 8.5 SCC-customized Put/Get and Mem-Copy Functions

### **int iRCCE\_put(t\_vcharp target, t\_vcharp source, int size, int rank)**

This is the SCC-optimized version of the `RCCE_put()` function (see Section 5.1 for more details). The function copies the contents of the buffer pointed to by **source** into the MPB location pointed to by **target**. The data type `t_vcharp` is similar to `volatile char*` and defined in `RCCE.h`.

<b>target</b>	an MPB address that will be converted to the appropriate address on the UE
<b>source</b>	start address of the data in private memory of the calling UE
<b>size</b>	size of the data to be put into the MPB in bytes
<b>rank</b>	rank (ID) of the target process (UE)

### **int iRCCE\_get(t\_vcharp target, t\_vcharp source, int size, int rank)**

This is the SCC-optimized version of the `RCCE_get()` function (see Section 5.1 for more details). The function copies the contents of the MPB location pointed to by **source** into the buffer pointed to by **target**. The data type `t_vcharp` is similar to `volatile char*` and defined in `RCCE.h`.

<b>target</b>	address of the destination buffer in private memory of the calling UE
<b>source</b>	an offset that, when combined with the remote rank, points to the source MPB
<b>size</b>	size of the data to be gotten from the MPB in bytes
<b>rank</b>	rank (ID) of the source process (UE)

**void\* iRCCE\_memcpy\_put(void\* dest, const void\* src, size\_t num)**

This function copies `num` bytes from memory area `src` to memory area `dest` in an SCC-optimized manner (see also Section 5.1). It can be used instead of the common `memcpy()` routine of `<string.h>`.

<code>dest</code>	start address of the destination memory area
<code>src</code>	start address of the source memory area
<code>num</code>	number of bytes to be copied from source to destination

**void\* iRCCE\_memcpy\_get(void\* dest, const void\* src, size\_t num)**

This function copies `num` bytes from memory area `src` to memory area `dest` in an SCC-optimized manner (see also Section 5.1). It can be used instead of the common `memcpy()` routine of `<string.h>`.

<code>dest</code>	start address of the destination memory area
<code>src</code>	start address of the source memory area
<code>num</code>	number of bytes to be copied from source to destination

## 8.6 Cancel Functions for Non-blocking Requests

**int iRCCE\_isend\_cancel(iRCCE\_SEND\_REQUEST \*request, int \*flag)**

This function tries to cancel a yet not finished non-blocking send request. If the request has not yet been started of being processed, the canceling should be successful, otherwise this request can no longer be canceled. The returned flag value indicates whether the canceling was successful or not. For more details about this function see Section 4.1.

<code>request</code>	request that should be removed from the waiting send request queue
<code>flag</code>	flag that indicates whether the canceling was successful (1) or not (0)

**int iRCCE\_irecv\_cancel(iRCCE\_RECV\_REQUEST \*request, int \*flag)**

This function tries to cancel a yet not finished non-blocking receive request. If the request has not yet been started of being processed, the canceling should be successful, otherwise this request can no longer be canceled. The returned flag value indicates whether the canceling was successful or not. For more details about this function see Section 4.1.

<code>request</code>	request that should be removed from the waiting receive request queue
<code>flag</code>	flag that indicates whether the canceling was successful (1) or not (0)

## 8.7 Functions for Handling Multiple Outstanding Requests

**void iRCCE\_init\_wait\_list(iRCCE\_WAIT\_LIST\* wait\_list)**

This function initializes a wait-list and must be called before adding pending communication requests to the respective wait-list. A wait-list of type `iRCCE_WAIT_LIST` can handle both send and receive requests.

<code>wait_list</code>	a pointer to an object of the opaque iRCCE data type <code>iRCCE_WAIT_LIST</code>
------------------------	---

**IRCCE\_add\_to\_wait\_list(IRCCE\_WAIT\_LIST\* wait\_list, IRCCE\_SEND\_REQUEST \*send\_request, IRCCE\_RECV\_REQUEST \*recv\_request)**

This function adds a pending send and/or receive request to a wait-list that must be initialized before by `IRCCE_init_wait_list()`. Either the send request argument or the receive request argument can be set to NULL if there is not request of the respective type to be added.

<code>wait_list</code>	a pointer to the respective wait-list object
<code>send_request</code>	a pointer to a pending send request to be added to the wait-list (can be NULL)
<code>recv_request</code>	a pointer to a pending receive request to be added to the wait-list (can be NULL)

**IRCCE\_test\_all(IRCCE\_WAIT\_LIST\* wait\_list, int \*flag)**

This function tests for completion of *all* requests in the passed wait-list. The flag is set to 1, if all respective requests are finished, or to 0 otherwise. See Section 4.2 for more details about this function.

<code>wait_list</code>	a pointer to the respective wait-list object
<code>flag</code>	flag that indicates whether the wait-list is processed (1) or not (0)

**IRCCE\_wait\_all(IRCCE\_WAIT\_LIST\* wait\_list)**

This function just waits for completion of *all* requests in the passed wait-list. See Section 4.2 for more details about this function.

<code>wait_list</code>	a pointer to the respective wait-list object
------------------------	--

**IRCCE\_test\_any(IRCCE\_WAIT\_LIST\* wait\_list, IRCCE\_SEND\_REQUEST \*\*send\_request, IRCCE\_RECV\_REQUEST \*\*recv\_request)**

This function tests for completion of *any* request in the passed wait-list. It returns a pointer to that finished request or NULL otherwise. In case of a finished send request, `send_request` points to this request and `recv_request` is set to NULL; and vice versa in case of a finished receive request.

<code>wait_list</code>	a pointer to the respective wait-list object
<code>send_request</code>	returned pointer to a finished send request (or NULL)
<code>recv_request</code>	returned pointer to a finished receive request (or NULL)

**IRCCE\_wait\_any(IRCCE\_WAIT\_LIST\* wait\_list, IRCCE\_SEND\_REQUEST \*\*send\_request, IRCCE\_RECV\_REQUEST \*\*recv\_request)**

This function just waits for completion of *any* request in the passed wait-list. In case of a finished send request, `send_request` points to this request and `recv_request` is set to NULL; and vice versa in case of a finished receive request. The rank (ID) of the respective sender/receiver can afterwards be determined via `send_request->dest/recv_request->source`.

<code>wait_list</code>	a pointer to the respective wait-list object
<code>send_request</code>	returned pointer to a finished send request (or NULL)
<code>recv_request</code>	returned pointer to a finished receive request (or NULL)



## Acknowledgment

The research and development of the iRCCE library was supported by Intel Corporation. We would like to thank especially Ulrich Hoffmann, Michael Konow and Michael Riepen of Intel Braunschweig for their help and guidance as well as Carsten Scholtes of the University of Bayreuth for his useful feedback.

## References

- [1] Intel Many-core Applications Research Community. <http://communities.intel.com/community/marc>.
- [2] Intel Corporation. *Intel MPI Benchmarks – Users Guide and Methodology Description*, 2006. Version 3.0.
- [3] Intel Corporation. *SCC External Architecture Specification (EAS)*, July 2010. Revision 0.98.
- [4] T. Mattson and R. van der Wijngaart. *RCCE: a Small Library for Many-Core Communication*. Intel Corporation, May 2010. Software 1.0-release.
- [5] T. Mattson, R. van der Wijngaart, M. Riepen, T. Lehnig, P. Brett, W. Haas, P. Kennedy, J. Howard, S. Vangal, N. Borkar, G. Ruhl, and S. Dighe. The 48-core SCC Processor: The Programmer’s View. In *Proceedings of the 2010 ACM/IEEE Conference on Supercomputing (SC10)*, New Orleans, LA, USA, November 2010.
- [6] Message Passing Interface Forum. *MPI: A Message-Passing Interface Standard*. High-Performance Computing Center Stuttgart (HLRS), September 2009. Version 2.2.