



An Efficient
Method for
Stream
Semantics
over RDMA

MacArthur
and Russell

Background

RDMA vs. TCP
Related Work
UNH EXS

Dynamic
Protocol

Motivation
Overview
Scenario

Performance
Evaluation

Simple
Distance

Conclusions

An Efficient Method for Stream Semantics over RDMA

Patrick MacArthur <pio3@cs.unh.edu>
Robert D. Russell <rdr@cs.unh.edu>

Department of Computer Science
University of New Hampshire
Durham, NH 03824-3591, USA

28th IEEE International Parallel & Distributed Processing
Symposium (IEEE IPDPS 2014)
May 21, 2014



Acknowledgements

An Efficient
Method for
Stream
Semantics
over RDMA

MacArthur
and Russell

Background

RDMA vs. TCP
Related Work
UNH EXS

Dynamic
Protocol

Motivation
Overview
Scenario

Performance
Evaluation

Simple
Distance

Conclusions

The authors would like to thank the University of New Hampshire InterOperability Laboratory for the use of their RDMA cluster for the development, maintenance, and testing of UNH EXS. We would also like to thank the UNH-IOL and Ixia for the use of an Anue network emulator for performance testing.

This material is based upon work supported by the National Science Foundation under Grant No. OCI-1127228 and under the National Science Foundation Graduate Research Fellowship Program under award number DGE-0913620.



An Efficient
Method for
Stream
Semantics
over RDMA

MacArthur
and Russell

Background

RDMA vs. TCP
Related Work
UNH EXS

Dynamic Protocol

Motivation
Overview
Scenario

Performance Evaluation

Simple
Distance

Conclusions

Background



Differences Between RDMA and TCP Sockets

An Efficient
Method for
Stream
Semantics
over RDMA

MacArthur
and Russell

Background

RDMA vs. TCP

Related Work
UNH EXS

Dynamic
Protocol

Motivation
Overview
Scenario

Performance
Evaluation

Simple
Distance

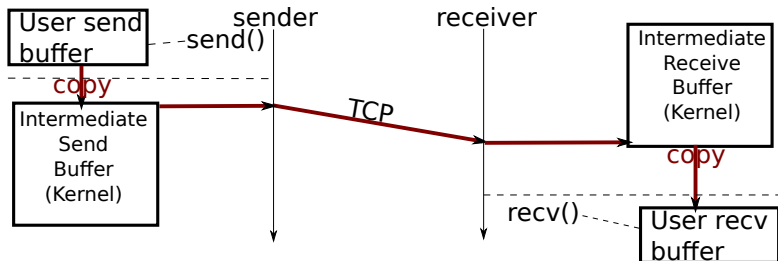
Conclusions

TCP (Transmission Control Protocol) Sockets

- Kernel involvement in all data transfers
- Buffered in kernel-space on both sides of connection
- Byte-stream oriented protocol
- Synchronous programming interface

RDMA (Remote Direct Memory Access)

- “Kernel bypass”: data transfers with no OS involvement
- “Zero-copy”: Direct virtual memory to virtual memory transfers
- Message-oriented
- Asynchronous programming interface





RDMA WRITE Data Transfer

An Efficient Method for Stream Semantics over RDMA

MacArthur and Russell

Background

RDMA vs. TCP

Related Work

UNH EXS

Dynamic Protocol

Motivation

Overview

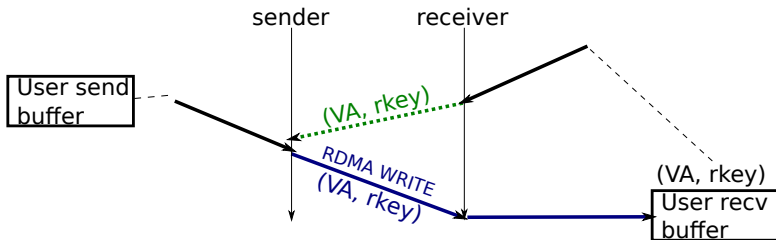
Scenario

Performance Evaluation

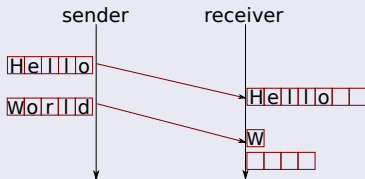
Simple

Distance

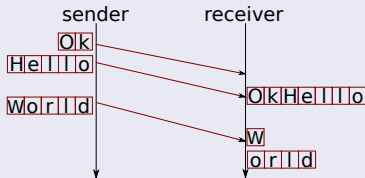
Conclusions



Message Transfer (UDP)



Byte Stream Transfer (TCP)





Prior Implementations of Sockets over RDMA

An Efficient
Method for
Stream
Semantics
over RDMA

MacArthur
and Russell

Background

RDMA vs. TCP

Related Work

UNH EXS

Dynamic
Protocol

Motivation
Overview
Scenario

Performance
Evaluation

Simple
Distance

Conclusions

- Sockets Direct Protocol (SDP) (defined by InfiniBand specification [InfiniBand 2011])
 - BCopy (buffering on both sides)
 - ZCopy (zero-copy, send() blocks) [Goldenberg 2005]
 - AZ-SDP (asynchronous, zero-copy, seefault handler) [Balaji 2006]
- uStream (asynchronous but not zero-copy) [Lin 2009]



Current Implementations of Sockets over RDMA

An Efficient
Method for
Stream
Semantics
over RDMA

MacArthur
and Russell

Background

RDMA vs. TCP

Related Work

UNH EXS

Dynamic
Protocol

Motivation

Overview

Scenario

Performance
Evaluation

Simple

Distance

Conclusions

- SMC-R (100% compatibility with TCP/IP and sockets)
- rsockets (high-performance sockets replacement)
[Hefty 2012]
- UNH EXS (extended sockets) [ISC 2005, Russell 2009]



UNH EXS (Extended Sockets)

An Efficient
Method for
Stream
Semantics
over RDMA

MacArthur
and Russell

Background

RDMA vs. TCP

Related Work

UNH EXS

Dynamic Protocol

Motivation

Overview

Scenario

Performance Evaluation

Simple

Distance

Conclusions

- Based on ES-API (Extended Sockets API) published by the Open Group [ISC 2005]
- Extensions to sockets API to provide asynchronous, zero-copy transfers
 - Memory registration (`exs_mregister()`, `exs_mderegister()`)
 - Event queues for completion of asynchronous events (`exs_qcreate()`, `exs_qdequeue()`, `exs_qdelete()`)
 - Asynchronous operations (`exs_send()`, `exs_recv()`, `exs_accept()`, `exs_connect()`)
- UNH EXS supports `SOCK_SEQPACKET` (reliable message-oriented) and `SOCK_STREAM` (reliable stream-oriented) modes



Example

Example asynchronous send operation

```
exs_mhandle_t mh = exs_mregister(buf, bufsize, EXS_ACCESS_READ);
exs_qhandle_t qh = exs_qcreate(10);

if (exs_send(fd, buf, bufsize, 0, mh, 0, qh) < 0) {
    perror("Could not start send operation");
    /* bail out */
}

/* do work in parallel with data transfer */

exs_event_t ev;
if (exs_qdequeue(qh, &ev, 1, NULL) < 0) {
    perror("Could not get send completion event");
    /* bail out */
}

fprintf(stderr, "Send of %d/%d bytes complete with errno=%d\n",
        bufsize, ev.exs_evt_union.exs_evt_xfer.exs_evt_length,
        ev.exs_evt_errno);
```



An Efficient
Method for
Stream
Semantics
over RDMA

MacArthur
and Russell

Background

RDMA vs. TCP
Related Work
UNH EXS

Dynamic
Protocol

Motivation
Overview
Scenario

Performance
Evaluation

Simple
Distance

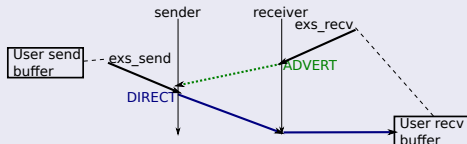
Conclusions

Dynamic Protocol

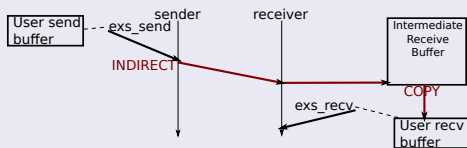


- Provide familiar **byte-stream** abstraction over RDMA
- **Dynamically** combine best aspects of zero-copy RDMA and fast send response benefit of TCP-style buffering.
- Deliver user data from sender to receiver **in order** with **no errors**
- Implementation of sender and receiver should be as **independent** as possible
- Work well over large **distance**
- Work **automatically** without user intervention

Direct Transfer (Streams and Messages)

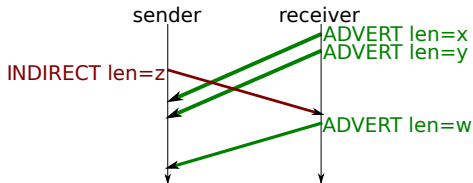


Indirect Transfer (Streams only)



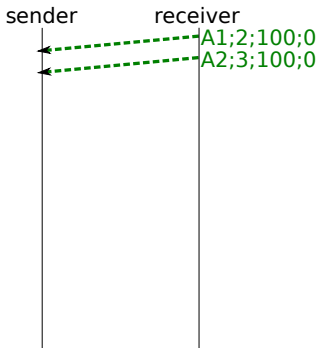
Key idea: UNH EXS automatically uses direct or indirect transfer based on current conditions

- `exs_recv(fd, buf, n, ...)` may actually receive between 1 and n bytes (**not** known in advance by receiver)
- `exs_send(fd, buf, n, ...)` will transfer n bytes in absence of network errors (due to asynchronous nature of EXS)
- Advertisements may arrive late (i.e., after sender has already sent corresponding data indirectly)
- Sender must distinguish between “fresh” and “stale” advertisements

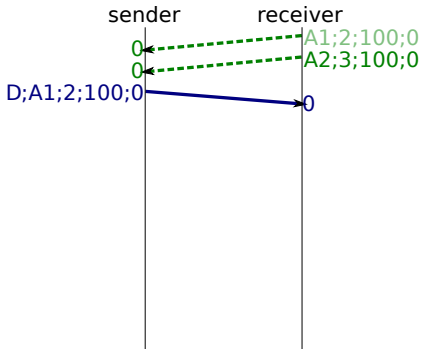




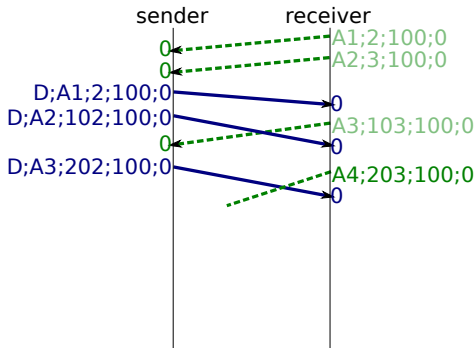
- Advertisement includes estimated **byte sequence number** and **phase number**
- Byte sequence number is estimated sequence number of first byte in stream satisfying associated `exs_recv()` request.
 - Each advertisement at receiver increases estimate by 1 (minimum transfer length)
 - Each transfer completion at receiver increases estimate by $n - 1$ (adjusts for actual number of bytes transferred)
 - Exact value known at sender
- Phase number is monotonically nondecreasing identifier of transfer phase
 - Incremented on each switch between direct and indirect transfers on both sides
 - At sender, always \geq that of last processed advertisement



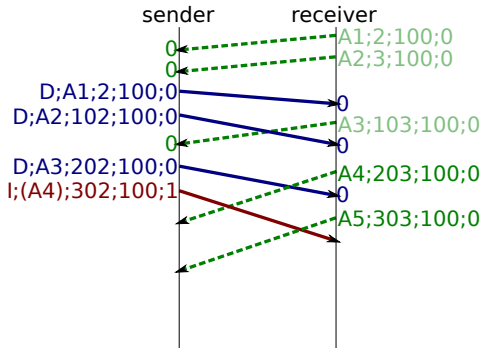
At start, sender has sent 2 bytes to receiver directly. Receiver issues two `exs_recv` requests.



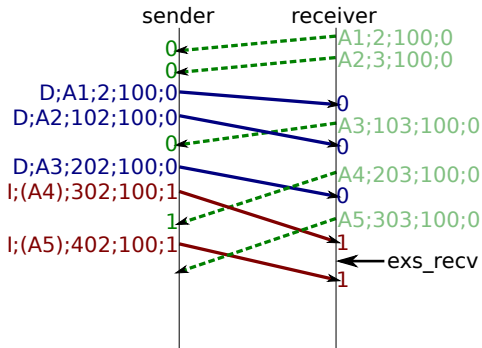
Sender issues single `exs_send` request.



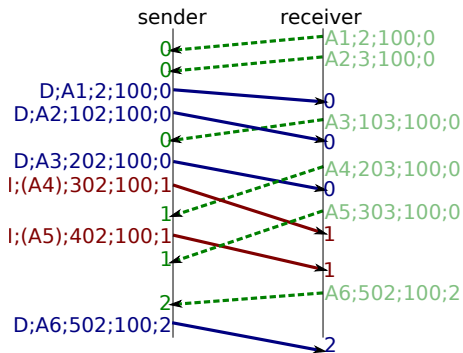
Sender and receiver continue to issue `exs_send` and `exs_rcv` requests.



Advertisement A4 is delayed, so sender sends **indirect** transfer.



Receiver **holds off** on sending off advertisements until all pending advertisements satisfied.



Receiver starts sending advertisements again once caught up. Sender matches advertisement once sequence number in ADVERT matches and phase number in ADVERT is greater.



An Efficient
Method for
Stream
Semantics
over RDMA

MacArthur
and Russell

Background

RDMA vs. TCP
Related Work
UNH EXS

Dynamic
Protocol

Motivation
Overview
Scenario

Performance
Evaluation

Simple
Distance

Conclusions

Performance Evaluation



Performance Evaluation

An Efficient Method for Stream Semantics over RDMA

MacArthur and Russell

Background

RDMA vs. TCP
Related Work
UNH EXS

Dynamic Protocol

Motivation
Overview
Scenario

Performance Evaluation

Simple
Distance

Conclusions

- Comparison of dynamic protocol with baseline protocols
 - **Direct-only** protocol: sender always waits for ADVERT (receiver never copies)
 - **Indirect-only** protocol receiver never sends ADVERTs (receiver always copies)
- Measure throughput, CPU usage, and percent of direct sends (average across entire run)



Simple Tests

An Efficient Method for Stream Semantics over RDMA

MacArthur and Russell

Background

RDMA vs. TCP
Related Work
UNH EXS

Dynamic Protocol

Motivation
Overview
Scenario

Performance Evaluation

Simple
Distance

Conclusions

- Systems used: Intel Xeon 2.40 GHz E5-2609 CPUs, 64 GB RAM, PCI-e Gen 3
- HCAs: Mellanox ConnectX-3 54.54 Gbps FDR InfiniBand HCAs
- Connection: Mellanox SX6036 FDR InfiniBand switch



Dynamic Protocol Throughput Comparison

outstanding sends = $\frac{1}{2}$ outstanding receives

An Efficient Method for Stream Semantics over RDMA

MacArthur and Russell

Background

RDMA vs. TCP
Related Work
UNH EXS

Dynamic Protocol

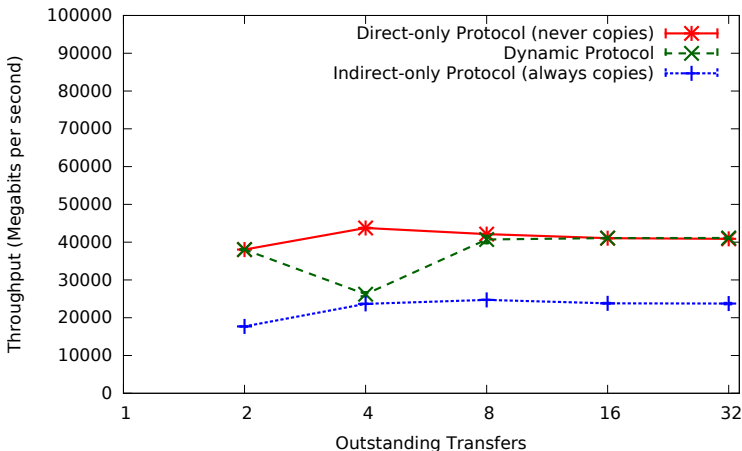
Motivation
Overview
Scenario

Performance Evaluation

Simple
Distance

Conclusions

Higher is better





Dynamic Protocol CPU Comparison

outstanding sends = $\frac{1}{2}$ outstanding receives

An Efficient Method for Stream Semantics over RDMA

MacArthur and Russell

Background

RDMA vs. TCP
Related Work
UNH EXS

Dynamic Protocol

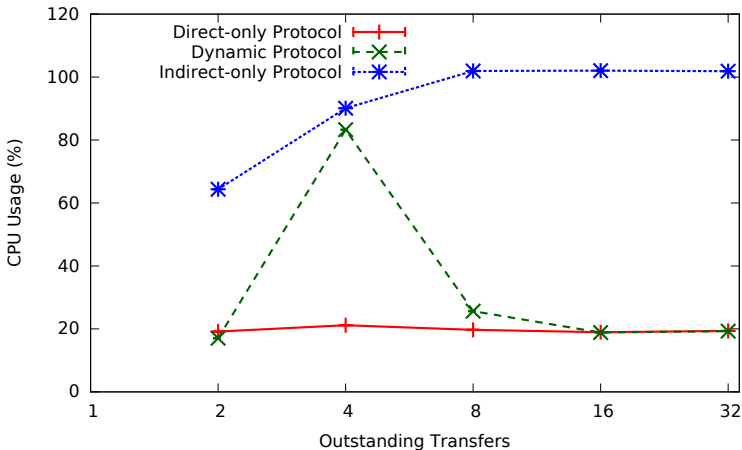
Motivation
Overview
Scenario

Performance Evaluation

Simple
Distance

Conclusions

Lower is better





Dynamic Protocol Throughput

Outstanding receives fixed at 32

An Efficient Method for Stream Semantics over RDMA

MacArthur and Russell

Background

RDMA vs. TCP
Related Work
UNH EXS

Dynamic Protocol

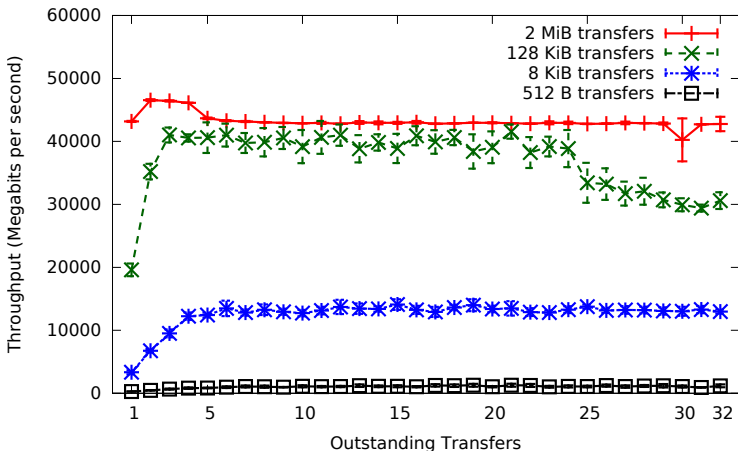
Motivation
Overview
Scenario

Performance Evaluation

Simple
Distance

Conclusions

Higher is better





Dynamic Protocol Percent Direct Sends

Outstanding receives fixed at 32

An Efficient Method for Stream Semantics over RDMA

MacArthur and Russell

Background

RDMA vs. TCP
Related Work
UNH EXS

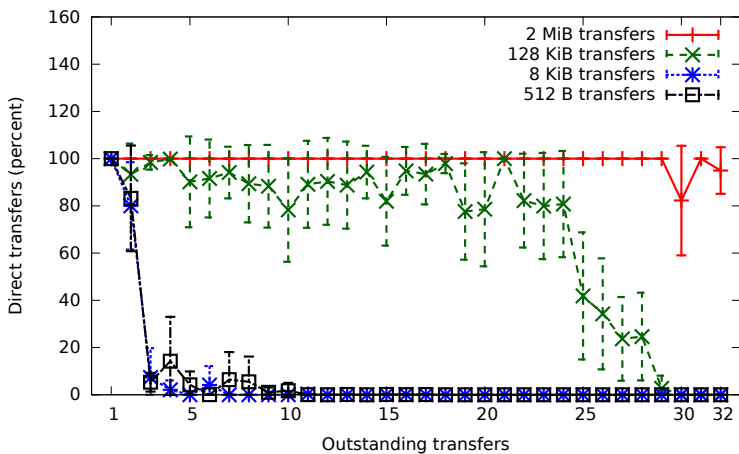
Dynamic Protocol

Motivation
Overview
Scenario

Performance Evaluation

Simple
Distance

Conclusions





An Efficient
Method for
Stream
Semantics
over RDMA

MacArthur
and Russell

Background

RDMA vs. TCP
Related Work
UNH EXS

Dynamic
Protocol

Motivation
Overview
Scenario

Performance
Evaluation

Simple
Distance

Conclusions

- Systems used: Intel Xeon 2.93 GHz X-5670 CPUs, 64 GB RAM, PCI-e Gen 2
- HCAs: Mellanox ConnectX-2 10 Gbps RoCE HCAs
- Distance simulated using Ixia ANUE Network Emulator introducing 48 ms round-trip delay



Throughput Over Distance

Outstanding sends and receives fixed at 32

An Efficient Method for Stream Semantics over RDMA

MacArthur and Russell

Background

RDMA vs. TCP

Related Work

UNH EXS

Dynamic Protocol

Motivation

Overview

Scenario

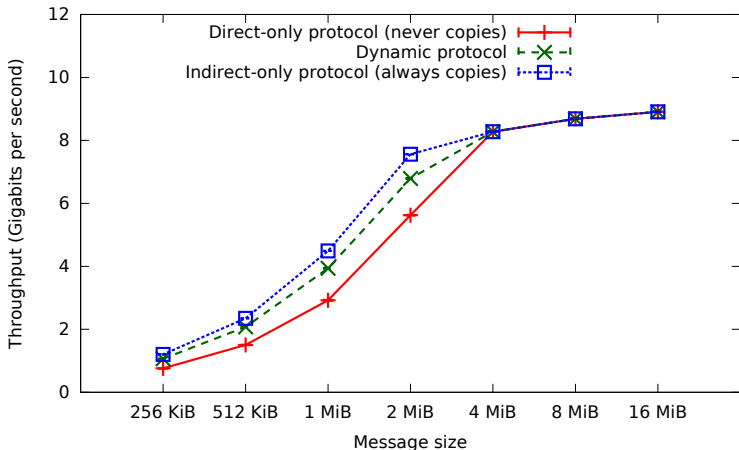
Performance Evaluation

Simple

Distance

Conclusions

Higher is better





An Efficient
Method for
Stream
Semantics
over RDMA

MacArthur
and Russell

Background

RDMA vs. TCP
Related Work
UNH EXS

Dynamic
Protocol

Motivation
Overview
Scenario

Performance
Evaluation

Simple
Distance

Conclusions

Conclusions



Contributions

An Efficient
Method for
Stream
Semantics
over RDMA

MacArthur
and Russell

Background

RDMA vs. TCP
Related Work
UNH EXS

Dynamic
Protocol

Motivation
Overview
Scenario

Performance
Evaluation

Simple
Distance

Conclusions

- Design of a **set of algorithms** for dynamically choosing between direct and indirect transfers in a **byte-stream protocol**
- Proof of **correctness**
- **Implementation** and testing of these algorithms within UNH EXS
- **Performance** evaluation



Thanks!

An Efficient
Method for
Stream
Semantics
over RDMA

MacArthur
and Russell

Background

RDMA vs. TCP
Related Work
UNH EXS

Dynamic
Protocol

Motivation
Overview
Scenario

Performance
Evaluation

Simple
Distance

Conclusions

Questions?

Patrick MacArthur <pio3@cs.unh.edu>

Robert D. Russell <rdr@cs.unh.edu>

<https://www.iol.unh.edu/services/research/unh-exs>




An Efficient
Method for
Stream
Semantics
over RDMA


MacArthur
and Russell

References

Backup

References

- 

Infiniband Trade Association,
“Supplement to Infiniband Architecture Specification
Volume 1, Release 1.2.1: Annex A4: Sockets Direct
Protocol (SDP),”
Oct. 2011.
- 

D. Goldenberg, M. Kagan, R. Ravid, and M. S. Tsirkin,
“Zero copy sockets direct protocol over
Infiniband—preliminary implementation and performance
analysis,”
*in High Performance Interconnects, 2005. Proceedings.
13th Symposium on.* IEEE, 2005, pp. 128–137.



P. Balaji, S. Bhagvat, H.-W. Jin, and D. K. Panda,
“Asynchronous zero-copy communication for synchronous
sockets in the sockets direct protocol (SDP) over
InfiniBand,”

*in Parallel and Distributed Processing Symposium, 2006.
IPDPS 2006. 20th International.* IEEE, 2006, pp. 8–pp.



Y. Lin, J. Han, J. Gao, and X. He,
“uStream: a user-level stream protocol over InfiniBand,”
*in Parallel and Distributed Systems (ICPADS), 2009 15th
International Conference on.* IEEE, 2009, pp. 65–71.



S. Hefty,
Rsockets.

Available: [https://www.openfabrics.org/
ofa-documents/doc_download/495-rsockets.html](https://www.openfabrics.org/ofa-documents/doc_download/495-rsockets.html)



Interconnect Software Consortium in association with the
Open Group,
“Extended Sockets API (ES-API) Issue 1.0,”
Jan. 2005.



———,
“A General-purpose API for iWARP and InfiniBand,”
*in the First Workshop on Data Center Converged and
Virtual Ethernet Switching (DC-CAVES)*, Sep. 2009.



Interconnect Software Consortium in association with the
Open Group,
“Extended Sockets API (ES-API) Issue 1.0,”
Jan. 2005.



An Efficient
Method for
Stream
Semantics
over RDMA

MacArthur
and Russell

References

Backup

Backup



Issue: Implementing “Zero-copy” Sockets

An Efficient
Method for
Stream
Semantics
over RDMA

MacArthur
and Russell

References

Backup

- Memory to be used for RDMA must (currently) be **registered**
- Sockets programmers make assumptions about memory used in I/O operations
 - Memory can be reused/freed as soon as `send()` returns
 - Alignment of messages does not matter
- The reality when using RDMA:
 - `ibv_post_send()` only queues send operation; buffer still in use when it returns
 - Not respecting adapter's natural alignment can cause severe performance degradation, especially on FDR adapters

```

1: while  $\neg$ EMPTY( $q_A$ ) do
2:    $A \leftarrow$  HEAD( $q_A$ )
3:   if PHASE_IS_INDIRECT( $P_S$ )  $\wedge$  ( $P_A < P_S \vee S_A < S_S$ )
   then
4:     if  $P_S < P_A$  then
5:        $P_S \leftarrow$  NEXT_PHASE( $P_A$ )
6:     end if
7:     throw away ADVERT  $A$ 
8:   else
9:     if PHASE_IS_INDIRECT( $P_S$ ) then
10:       $P_S \leftarrow P_A$   $\triangleright P_S$  is now direct
11:    end if
12:     $S_S \leftarrow S_S + l_w$ 
13:    send direct transfer

```

```

14:         return
15:     end if
16: end while
17: if  $\neg \text{FULL}(b_s)$  then
18:     if  $\text{PHASE\_IS\_DIRECT}(P_s)$  then
19:          $P_s \leftarrow \text{NEXT\_PHASE}(P_s)$             $\triangleright P_s$  is now indirect
20:     end if
21:      $S_s \leftarrow S_s + l_w$ 
22:      $b_s \leftarrow b_s - l_w$ 
23:     send indirect transfer
24:     return
25: end if

```

```

1: if  $b_r > 0 \vee k_a > 0 \vee k_b > 0$  then
2:   do not send ADVERT
3:   return
4: end if
5: if PHASE_IS_INDIRECT( $P_r$ ) then
6:    $P_r \leftarrow \text{NEXT\_PHASE}(P_r)$            ▷  $P_r$  is now direct
7: end if
8:  $P_A \leftarrow P_r$ 
9:  $S_A \leftarrow S'_r$ 
10: if MSG_WAITALL is set then
11:    $S'_r \leftarrow S'_r + l_r$ 
12: else
13:    $S'_r \leftarrow S'_r + 1$ 
14: end if

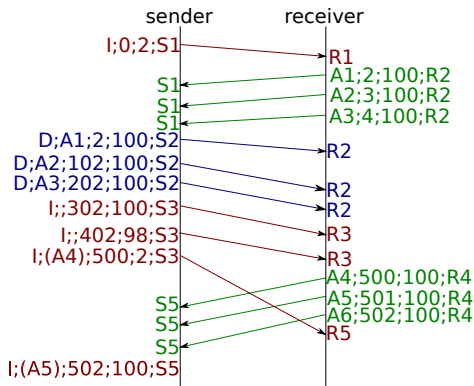
```

```

1: if incoming transfer is direct then
2:    $S_r \leftarrow S_r + I_w$ 
3:   if MSG_WAITALL was not set then
4:      $S'_r \leftarrow S'_r + I_w - 1$ 
5:   end if
6:   do normal processing
7: else                                ▷ incoming transfer is indirect
8:   if PHASE_IS_DIRECT( $P_r$ ) then
9:      $P_r \leftarrow \text{NEXT\_PHASE}(P_r)$ 
10:  end if
11:  do normal processing
12: end if

```

- 1: copy data from stream buffer
- 2: send ACK to sender notifying of freed space
- 3: $b_r \leftarrow b_r - l_c$
- 4: $S_r \leftarrow S_r + l_c$
- 5: **if** advert sent and MSG_WAITALL was *not* set **then**
- 6: $S'_r \leftarrow S'_r + l_c - 1$
- 7: **end if**



If sender phase were not incremented after processing A4-A6, would incorrectly match A6.