



# Troubleshooting Guidelines for Installing and Using MVAPICH2 and MVAPICH2-X

MVAPICH2 User Group (MUG) Meeting

by

**Devendar Bureddy**

The Ohio State University

E-mail: [bureddy@cse.ohio-state.edu](mailto:bureddy@cse.ohio-state.edu)

<http://www.cse.ohio-state.edu/~bureddy>



# Outline

- User Resources
- Frequently reported issues and Common mistakes
- Useful Diagnostics
- Performance Troubleshooting
- Getting help and Bug report details

# User Resources

- [MVAPIVH2 Quick Start Guide](#)
- [MVAPICH2 User Guide](#)
  - Long and very detailed
  - FAQ
- [MVAPICH2 Web-Site](#)
  - [Overview](#) and [Features](#)
  - [Reference performance](#)
  - [Publications](#)
- [Mailing List](#) Support
  - mvapich-discuss
- [Mailing List Archives](#)
- All above resources accessible from: <http://mvapich.cse.ohio-state.edu/>

# Outline

- User Resources
- Frequently reported issues and Common mistakes
- Useful Diagnostics
- Performance Troubleshooting
- Getting help and Bug report details

## Frequently reported issues and Common mistakes

- Job Startup issues
- MPI\_Init and Other MPI errors
- Creation of CQ or QP failure
- Failed to register memory with Infiniband HCA
- Multicast group creation failed
- Infiniband setup issues
- MVAPICH2 over RoCE issues
- MPI + OpenMP, Multi-threaded MPI shows bad performance

# Job Startup issues

- **Symptoms**

- [mpirun\_rsh][child\_handler] Error in init phase, aborting! (0/2 mpispawn connections)

- **Cause**

- Host file is not correct
- SSH issues

- **Troubleshooting**

- Verify host file
- Password less ssh
- DNS or /etc/hosts

# MPI\_Init and Other MPI errors

- **Symptoms**

- “Fatal error in MPI\_Init:  
Other MPI error”

- **Cause**

- Could be because of multiple reasons

- **Troubleshooting**

- Reconfigure with `–enable-g=dbg –enable fast=none` to better understand the problem

[cli\_0]: aborting job:

Fatal error in MPI\_Init:

Other MPI error, error stack:

MPID\_Init\_thread(408).....:

MPID\_Init(308).....: channel initialization failed

MPIDI\_CH3\_Init(283).....:

MPIDI\_CH3I\_RDMA\_init(171)....:

rdma\_setup\_startup\_ring(389): cannot open hca device

# Creation of CQ or QP failure

## • Symptoms

- **libibverbs: Warning: RLIMIT\_MEMLOCK is 32768 bytes.**

This will severely limit memory registrations.

Other MPI error, error stack:

MPID\_Init\_thread(449).....:

MPID\_Init(365).....: channel initialization failed

MPIDI\_CH3\_Init(313).....:

MPIDI\_CH3I\_RDMA\_init(170)....:

rdma\_setup\_startup\_ring(416): **cannot create cq**

## • Cause

- Memory buffers used in verbs operations and ib context uses pinned memory
- Inability to pin the required memory

## • Troubleshooting

- Make sure enough memory set for “max locked memory” (limit -l)
- recommended “unlimited” on all compute nodes
- User guide section

- [http://mvapich.cse.ohio-state.edu/support/user\\_guide\\_mvapich2-2.0a.html#x1-1360009.4.3](http://mvapich.cse.ohio-state.edu/support/user_guide_mvapich2-2.0a.html#x1-1360009.4.3)



# Failed to register memory with InfiniBand HCA

- **Symptoms**

- “Cannot register vbuf region”
- “Abort: vbuf pool allocation failed”
- QP errors, node failures

- **Cause**

- Limited registered (pinned) memory

- **Troubleshooting**

- OFED parameters : `log_num_mtt`, `log_mtt_per_seg`
- $\text{max\_reg\_mem} = (2^{\text{log\_num\_mtt}}) * (2^{\text{log\_mtts\_per\_seg}}) * \text{PAGE\_SIZE}$
- Some OFED default values are too low (< 2GB)
- clusters with large physical memory (> 64)
- **Recommendation** : increase `log_num_mtt` value
  - $\text{max\_reg\_mem} = (2^{24}) * (2^1) * (4 \text{ kB}) = 128 \text{ GB}$
- User guide section
  - [http://mvapich.cse.ohio-state.edu/support/user\\_guide\\_mvapich2-2.0a.html#x1-1130009.1.1](http://mvapich.cse.ohio-state.edu/support/user_guide_mvapich2-2.0a.html#x1-1130009.1.1)

# Multicast group creation failed

- **Symptoms**

- [host1:mpi\_rank\_0][create\_2level\_comm]

**Warning: Multicast group setup failed. Not using any multicast features**

- **Cause**

- Umad device permission
- OpenSM issues

- **Troubleshooting**

- Check umad device user permissions

```
$ ls -l /dev/infiniband/umad0
```

```
crw-rw-rw- 1 root root 231, 0 Aug  9 02:04 /dev/infiniband/umad0
```

- Slow opensm response

- MV2\_MCAST\_COMM\_INIT\_TIMEOUT

- Maximum multicast groups reached ( very unlikely). Check opensm logs

- User guide section

- [http://mvapich.cse.ohio-state.edu/support/user\\_guide\\_mvapich2-2.0a.html#x1-620006.9](http://mvapich.cse.ohio-state.edu/support/user_guide_mvapich2-2.0a.html#x1-620006.9)

# InfiniBand setup issues

- **Symptoms**

- [0->6150] send desc error, wc\_opcode=0  
[0->6150] **wc.status=12**, wc\_opcode=0, vbuf->phead->type=25 = XXXX  
[4979] Abort: [] Got completion with error 12, vendor code=0x81, **dest rank=6150**
- wc.status : **12 (IBV\_WC\_RETRY\_EXC\_ERR), 13 (IBV\_WC\_RNR\_RETRY\_EXC\_ERR)**

- **Cause**

- Bad QP attributes
- Loose cable, bad HCA or a bad switch blade
- Remote side is in a bad state
- Heavy congestion in the network

- **Troubleshooting**

- MV2\_DEFAULT\_RETRY\_COUNT
- Map src, dest ranks to host file and check those specific nodes

# MVAPICH2 over RoCE issues

- **Symptoms**

- Intermittent hangs

- **Cause**

- Most likely setup issues

- **Troubleshooting**

- Requires loss-less Ethernet fabric
- Configure Ethernet switch to treat RoCE traffic as loss-less
- Create a separate VLAN interface
- All VLAN interfaces appear as additional GID index
- Select non-default GID index with MV2\_DEFAULT\_GID\_INDEX
- Use VLAN IP addresses in */etc/mv2.conf* in RDMA CM mode
- User guide section
  - [http://mvapich.cse.ohio-state.edu/support/user\\_guide\\_mvapich2-2.0a.html#x1-380005.2.7](http://mvapich.cse.ohio-state.edu/support/user_guide_mvapich2-2.0a.html#x1-380005.2.7)

# MPI + OpenMP , Multi-threaded MPI shows bad performance

- **Symptoms**

Poor performance, hangs

- **Cause**

CPU affinity enabled by default

All OpenMP , pthreads in the application process bind to same core

- **Troubleshooting**

- Turn off affinity

`MV2_ENABLE_AFFINITY = 0`

- Choose binding level

`MV2_CPU_BINDING_LEVEL=socket`

- User guide section

- [http://mvapich.cse.ohio-state.edu/support/user\\_guide\\_mvapich2-2.0a.html#x1-550006.5](http://mvapich.cse.ohio-state.edu/support/user_guide_mvapich2-2.0a.html#x1-550006.5)

# Outline

- User Resources
- Frequently reported issues and Common mistakes
- **Useful Diagnostics**
- Performance Troubleshooting
- Getting help and Bug report details

## Useful Diagnostics

- What parameters are being used by my job?
- Where is the segmentation fault?
- What is the peak memory used by my app?
- Is process binding working as expected?

# What parameters are being used by my job?

- **MV2\_SHOW\_ENV\_INFO**

- Show values of the run time parameters
- 1 ( short list), 2 (full list)

- **Example**

```
$ mpirun_rsh -np 2 -hostfile hfile MV2_SHOW_ENV_INFO=1 ./exec
```

```
MVAPICH2-2.0a Parameters
```

```
-----  
PROCESSOR ARCH NAME      : MV2_ARCH_INTEL_XEON_E5_2680_16  
HCA NAME                  : MV2_HCA_MLX_CX_FDR  
HETEROGENEOUS             : NO  
MV2_VBUF_TOTAL_SIZE      : 17408  
MV2_IBA_EAGER_THRESHOLD  : 17408  
MV2_RDMA_FAST_PATH_BUF_SIZE : 5120  
MV2_EAGERSIZE_1SC        : 8192  
MV2_PUT_FALLBACK_THRESHOLD : 8192  
MV2_GET_FALLBACK_THRESHOLD : 0  
MV2_SMP_EAGERSIZE        : 8193  
MV2_SMPI_LENGTH_QUEUE    : 524288  
MV2_SMP_NUM_SEND_BUFFER  : 16  
MV2_SMP_BATCH_SIZE       : 8  
-----
```



## What parameters are being used by my job? (contd.)

- MPI-T
  - Initial support added in upcoming MVAPICH2 2.0a release
  - Several variables exposed with this interface to the tools
    - Memory allocation and usage information
    - Different collective algorithm invocation counters
    - Shared-memory usage tracing
    - UD retransmission count
    - Progress polling counters
    - Expected and unexpected receive queue matching attempts
    - Many more planned ..
  - Several control variables that can be set/tuned and runtime

## What parameters are being used by my job? (contd.)

- MVAPICH2 control variables as seen by a sample tool developed at LLNL

Variable	VRB	Class	Type	Bind	R/O	CNT	ATM
posted_recvq_length	U/D-2	LEVEL	UINT	n/a	YES	YES	NO
unexpected_recvq_length	U/D-2	LEVEL	UINT	n/a	YES	YES	NO
posted_recvq_match_attempts	U/D-2	COUNTER	UNKNOW	n/a	NO	YES	NO
unexpected_recvq_match_attempts	U/D-2	COUNTER	UNKNOW	n/a	NO	YES	NO
time_failed_matching_postedq	U/D-2	TIMER	DOUBLE	n/a	NO	YES	NO
time_matching_unexpectedq	U/D-2	TIMER	DOUBLE	n/a	NO	YES	NO
unexpected_recvq_buffer_size	U/D-2	LEVEL	UNKNOW	n/a	YES	YES	NO
mem_allocated	U/B-1	LEVEL	ULLONG	n/a	YES	YES	NO
mem_allocated	U/B-1	HIGHWAT	ULLONG	n/a	YES	YES	NO
mv2_progress_poll_count	D/B-7	COUNTER	ULONG	n/a	NO	NO	NO
mv2_rdma_ud_retransmit_count	D/B-7	COUNTER	ULONG	n/a	YES	NO	NO
coll_bcast_binomial	U/B-1	COUNTER	ULLONG	n/a	YES	YES	NO
coll_bcast_scatter_doubling_allgather	U/B-1	COUNTER	ULLONG	n/a	YES	YES	NO
coll_bcast_scatter_ring_allgather	U/B-1	COUNTER	ULLONG	n/a	YES	YES	NO
mv2_num_2level_comm_requests	U/D-2	COUNTER	ULONG	n/a	YES	YES	NO
mv2_num_2level_comm_success	U/D-2	COUNTER	ULONG	n/a	YES	YES	NO
mv2_num_shmem_coll_calls	T/B-4	COUNTER	ULONG	n/a	YES	YES	NO
mv2_coll_bcast_binomial	T/B-4	COUNTER	ULLONG	n/a	YES	YES	NO
mv2_coll_bcast_scatter_doubling_allgather	T/B-4	COUNTER	ULLONG	n/a	YES	YES	NO
mv2_coll_bcast_scatter_ring_allgather	T/B-4	COUNTER	ULLONG	n/a	YES	YES	NO
mv2_coll_bcast_scatter_ring_allgather_shm	T/B-4	COUNTER	ULLONG	n/a	YES	YES	NO
mv2_coll_bcast_shmem	T/B-4	COUNTER	ULLONG	n/a	YES	YES	NO
mv2_coll_bcast_knomial_internode	T/B-4	COUNTER	ULLONG	n/a	YES	YES	NO
mv2_coll_bcast_knomial_intranode	T/B-4	COUNTER	ULLONG	n/a	YES	YES	NO
mv2_coll_bcast_mcast_internode	T/B-4	COUNTER	ULLONG	n/a	YES	YES	NO
mv2_coll_bcast_pipeline	T/B-4	COUNTER	ULLONG	n/a	YES	YES	NO

# Where is the segmentation fault?

- **MV2\_DEBUG\_SHOW\_BACKTRACE**

- Shows backtrace with debug builds (--enable-g=dbg, --enable-fast=none)

- **Example**

- segmentation fault report with out much information

```
[host1:mpi_rank_0][error_sighandler] Caught error: Segmentation fault (signal 11)
```

- `mpirun_rsh -np 2 --hostfile hfile MV2_DEBUG_SHOW_BACKTRACE=1 ./exec`

```
[error_sighandler] Caught error: Segmentation fault (signal 11)
```

```
[print_backtrace] 0: libmpich.so.10(print_backtrace+0x22) [0x2af447e29d9a]
```

```
[print_backtrace] 1: libmpich.so.10(error_sighandler+0x7c) [0x2af447e29ef2]
```

```
[print_backtrace] 2: libmpich.so.10(allocate_vbufs+0x71) [0x2af447de6d9f]
```

```
[print_backtrace] 3: libmpich.so.10(rdma_iba_allocate_memory+0x101) [0x2af447dd5ca2]
```

```
[print_backtrace] 4: libmpich.so.10(MPIDI_CH3I_RDMA_init+0x1569) [0x2af447dce9f1]
```

```
[print_backtrace] 5: libmpich.so.10(MPIDI_CH3_Init+0x406) [0x2af447da32f4]
```

```
[print_backtrace] 6: libmpich.so.10(MPID_Init+0x31f) [0x2af447d8a91b]
```

```
[print_backtrace] 7: libmpich.so.10(MPIR_Init_thread+0x3e0) [0x2af447f90aca]
```

```
[print_backtrace] 8: libmpich.so.10(MPI_Init+0x1de) [0x2af447f8f645]
```

```
[print_backtrace] 9: ./mpi_hello() [0x400746]
```

# What is the peak memory used by my app?

- **MV2\_DEBUG\_MEM\_USAGE\_VERBOSE**

- Show memory usage statistics
- 1 ( rank 0 usage), 2 ( all ranks)

- **Example**

```
$ mpirun_rsh -np 2 -hostfile hfile MV2_DEBUG_MEM_USAGE_VERBOSE=1 ./exec
```

```
[mv2_print_mem_usage]      VmPeak:      79508 kB  VmHWM:      16340 kB
```

```
[mv2_print_vbuf_usage_usage]  RC VBUFs:512  UD VBUFs:0 TOT MEM:8828 kB
```

# Is process binding working as expected?

- **MV2\_SHOW\_CPU\_BINDING**

- Display CPU binding information
- Launcher independent

- **Examples**

- **MV2\_SHOW\_CPU\_BINDING=1 MV2\_CPU\_BINDING\_POLICY=scatter**

```
-----CPU AFFINITY-----  
RANK:0 CPU_SET: 0  
RANK:1 CPU_SET: 8
```

- **MV2\_SHOW\_CPU\_BINDING=1 MV2\_CPU\_BINDING\_POLICY=core**

```
-----CPU AFFINITY-----  
RANK:0 CPU_SET: 0  
RANK:1 CPU_SET: 1
```

- **MV2\_SHOW\_CPU\_BINDING=1 MV2\_CPU\_BINDING\_POLICY=scatter MV2\_CPU\_BINDING\_LEVEL=socket**

```
-----CPU AFFINITY-----  
RANK:0 CPU_SET: 0 1 2 3 4 5 6 7  
RANK:1 CPU_SET: 8 9 10 11 12 13 14 15
```

- **MV2\_SHOW\_CPU\_BINDING=1 MV2\_CPU\_BINDING\_POLICY=bunch MV2\_CPU\_BINDING\_LEVEL=socket**

```
-----CPU AFFINITY-----  
RANK:0 CPU_SET: 0 1 2 3 4 5 6 7  
RANK:1 CPU_SET: 0 1 2 3 4 5 6 7
```

# Outline

- User Resources
- Useful Diagnostics
- Frequently reported issues and Common mistakes
- **Performance Troubleshooting**
- Getting help and Bug report details

# Performance Trouble shooting

- Check “active\_speed” in “ibv\_devinfo -v” output
- Check OFED memory registration limits (log\_num\_mtt, log\_mtt\_per\_seg)
- Increase registration cache size
  - MV2\_NDREG\_ENTRIES, MV2\_NDREG\_ENTRIES\_MAX
- Are huge pages configured?
- SMP copy block size : MV2\_SMP\_SEND\_BUF\_SIZE
- Small message performance
  - RDMA fast path thresholds
    - MV2\_NUM\_RDMA\_BUFFER, MV2\_RDMA\_FAST\_PATH\_BUF\_SIZE
  - Eager thresholds
    - MV2\_IBA\_EAGER\_THRESHOLD, MV2\_VBUF\_TOTAL\_SIZE
- Large message performance
  - RNDV protocols : MV2\_RNDV\_PROTOCOL
- Collectives
  - Try different algorithms , change algorithm specific parameters
  - More in later talks

# Outline

- User Resources
- Frequently reported issues and common mistakes
- Useful Diagnostics
- Performance Troubleshooting
- **Getting help and Bug report details**



## Getting Help

- Check the [MVAPICH2 FAQ](#)
- Check the [Mailing List Archives](#)
- Basic System Diagnostics
  - `ibv_devinfo` - at least one port should be `PORT_ACTIVE`
  - `ulimit -l` - should be “unlimited” on all compute nodes
  - host resolution: DNS or `/etc/hosts`
  - password-less ssh login
  - run IB perf tests for all the message sizes(-a option)
    - `ib_send_lat`, `ib_send_bw`
  - run system program (like `hostname`) and MPI hello world program

## Getting Help (contd.)

- More diagnostics
  - Already fixed issue: always try with latest release
  - Regression: verifying with previous release
  - Application issue: verify with other MPI libraries
  - Launcher issue: verifying with multiple launchers (mpirun\_rsh, mpiexec.hydra)
  - Debug mode
  - Compiler optimization issues: try with different compiler

# Submitting Bug Report

- Subscribe to mvapich-discuss and send problem report
- Include as much information as possible
- Run-time issues
  - Config flags (“mpiname -a” output)
  - Exact command used to run the application
  - Run-time parameters in the environment
  - Standalone reproducer program
  - Information about the IB network
    - OFED version
    - ibv\_devinfo
  - Remote system access

## Submitting Bug Report (contd.)

- Build and Installation issues
  - MVAPICH2 version
  - Compiler version
  - Platform details ( OS, kernel version..etc)
  - Configure flags
  - Attach Config.log file
  - Attach configure, make and make install step output
    - `./configure {-flags} 2>&1 | tee config.out`
    - `Make 2>&1 | tee make.out`
    - `Make install 2>&1 | tee install.out`

# Web Pointers

NOWLAB Web Page

<http://nowlab.cse.ohio-state.edu>

MVAPICH Web Page

<http://mvapich.cse.ohio-state.edu>

