

Latest version of the slides available at

http://cse.osu.edu/~subramon/mug20-mvapich2-tutorial.pdf

How to Boost the Performance of Your HPC/AI Applications with MVAPICH2 Libraries?

A Tutorial at MUG'20

by

The MVAPICH Team

The Ohio State University

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

Overview of the MVAPICH2 Project

- High Performance open-source MPI Library
- Support for multiple interconnects
 - InfiniBand, Omni-Path, Ethernet/iWARP, RDMA over Converged Ethernet (RoCE), and AWS EFA
- Support for multiple platforms
 - x86, OpenPOWER, ARM, Xeon-Phi, GPGPUs (NVIDIA and AMD (upcoming))
- Started in 2001, first open-source version demonstrated at SC '02
- Supports the latest MPI-3.1 standard
- <u>http://mvapich.cse.ohio-state.edu</u>
- Additional optimized versions for different systems/environments:
 - MVAPICH2-X (Advanced MPI + PGAS), since 2011
 - MVAPICH2-GDR with support for NVIDIA GPGPUs, since 2014
 - MVAPICH2-MIC with support for Intel Xeon-Phi, since 2014
 - MVAPICH2-Virt with virtualization support, since 2015
 - MVAPICH2-EA with support for Energy-Awareness, since 2015
 - MVAPICH2-Azure for Azure HPC IB instances, since 2019
 - MVAPICH2-X-AWS for AWS HPC+EFA instances, since 2019
- Tools:
 - OSU MPI Micro-Benchmarks (OMB), since 2003
 - OSU InfiniBand Network Analysis and Monitoring (INAM), since 2015



- Used by more than 3,100 organizations in 89 countries
- More than 820,000 (> 0.8 million) downloads from the OSU site directly
- Empowering many TOP500 clusters (June '20 ranking)
 - 4th, 10,649,600-core (Sunway TaihuLight) at NSC, Wuxi, China
 - 8th, 448, 448 cores (Frontera) at TACC
 - 12th, 391,680 cores (ABCI) in Japan
 - 18th, 570,020 cores (Nurion) in South Korea and many others
- Available with software stacks of many vendors and Linux Distros (RedHat, SuSE, OpenHPC, and Spack)
- Partner in the 8th ranked TACC Frontera system
- Empowering Top500 systems for more than 15 years

Architecture of MVAPICH2 Software Family (for HPC and DL)

High Performance Parallel Programming Models									
Message Passing Interface	PGAS	Hybrid MPI + X							
(MPI)	(UPC, OpenSHMEM, CAF, UPC++)	(MPI + PGAS + OpenMP/Cilk)							



* Upcoming

Production Quality Software Design, Development and Release

- Rigorous Q&A procedure before making a release
 - Exhaustive unit testing
 - Various test procedures on diverse range of platforms and interconnects
 - Test 19 different benchmarks and applications including, but not limited to
 - OMB, IMB, MPICH Test Suite, Intel Test Suite, NAS, Scalapak, and SPEC
 - Spend about 18,000 core hours per commit
 - Performance regression and tuning
 - Applications-based evaluation
 - Evaluation on large-scale systems
- All versions (alpha, beta, RC1 and RC2) go through the above testing

Automated Procedure for Testing Functionality

- Test OMB, IMB, MPICH Test Suite, Intel Test Suite, NAS, Scalapak, and SPEC
- Tests done for each build done build "buildbot" \bullet
- Test done for various different combinations of *environment variables* meant to \bullet trigger different communication paths in MVAPICH2

Results Grid for (QA-PATCHES	/master	with 513	d83 with	n test list	runs						Results	or mvapich2	🖉 Testing - 🔳 Lists - 🔢 Grid 🔒	History 🖀 Latest	🕈 Run
Branch Filter														mvapich2-QA-PATC	CHES/master	ger
next-ea ruhel	er-v master-x ne la/nextadr next-ads	xt-gdr mv2	2x-mv2gdr-merg	ge-gar m	v2x-mv2gdr-m	ierge-x	next-v mast	er-ea mvzx	-mv2gdr-merg	je master-g	ar	mvapich2 / m	apich2-git / QA	A-PATCHES/master / Branch Grid / 513d83 /	Rev Grid / 597858	
												View Results	View -			
Revision Filter	All 513d83 be	173f 7a053	37 9074b3	33c936	9cff8b	ae118f	211aa5 479	e88 a119e4	l.							
												status:		passed		
Cluster Filter A	All nowlab ri	gordon sta	ampede ri2	ł hpcac	ibmfrs00	talapas	talapas-In1					Branch		O&-DATCHES/master		
												branch.		Contraction and an and a second se		
unted Runs	Total Runs	Test List C	Count	Succ	ess Rate		Lost Rate	Failure Ra	ite	Running Ra	te	Revision:		513d83216fda61a70d041e33390900c0cbe53	544	
unted Runs	Total Runs	Test List C	Count	Succ. 70.06	ess Rate %		Lost Rate	Failure Ra	ite	Running Ra	te	Revision: Channel:		513d83216fda61a70d041e33390900c0cbe53	544	
unted Runs	Total Runs 971	Test List (1399	Count	Succ 70.06	ess Rate %		Lost Rate	Failure Ra	ite	Running Ra	ite	Revision: Channel: Group:		513d83216lda61a70d041e33390900c0cbe53 gen2 basic_1	544	
ounted Runs	Total Runs 971	Test List (1399	Count	Succ 70.06	ess Rate %		Lost Rate 10.49%	Failure Ra	ite	Running Ra	ite	Revision: Channel: Group: Type:		613d83216fda61a70d041e33390900c0cbe53 gen2 basic_1 mpich2	544	
unted Runs	Total Runs 971	Test List (1399	Count	5ucc 70.06	ess Rate %		Lost Rate 10.49%	Failure Ra 19.34%	ite	Running Ra	ite	Revision: Channel: Group: Type: Cluster:		613d83216/ta61a70d041e33390900c0cbe63i gen2 basic_1 mpich2 ri		
Orenne - Danae -	Total Runs 971	Test List (1399	Count	Succ 70.06	ess Rate % gen2	mitech	Lost Rate 10.49%	Failure Ra	Ite	Running Ra	te	Revision: Channel: Group: Type: Cluster: Results ID:		615852168561463550960c0cbe535 gen2 basic_1 rrpich2 rl 55027858	44	
unted Runs 1 Groups : / Types	Total Runs 971 compilation	Test List (1399	Count	Succ 70.06	ess Rate % gen2 intel	mpibench	Lost Rate 10.49%	Failure Ra 19:34%	nas	Running Ra 0.1% nas btio	scalapack	Revision: Channel: Group: Type: Cluster: Results ID: Builder Loca	tion:	51585216/14/14/23350900/c/cbe53 gen2 basic_1 rrpich2 rl 55027658 /homenunbol/mapich2/install/QA-PATCHE5ii	nasteribasici 613d63216/da	61a700
unted Runs 1 Groups :/ Types	Total Runs 971 compilation	Test List (1399 imb	imb4	Succ 70.06	ess Rate % gen2 intel	mpibench	Lost Rate 10.49% mpich2	Failure Ra 19.34%	nas	Running Ra 0.1% nas btio	scalapack	Revision: Channel: Group: Type: Cluster: Results ID: Builder Loci Running Loc	tion: ation:	51383216861873001163359990000066533 gen2 basic_1 nmpich2 n 50027858 /homehunbotimvapich2insta8QA-PATCHESK /homehunbotimvapich2insta8QA-PATCHESK	naster/basic/613d63216/da	161a70 xe53b4
1 Groups () Types collectives allgathery	Total Runs 971 compilation	Test List (1399	imb-4	Succ 70.06 imb4 cuda	ess Rate % gen2 intel 513683 21 22 4	mpibench	Lost Rate 10.49% mpich2 513d53	Failure Ra 19.34% mpich2 cyclic 513653	nas Nija	Running Ra 0.1% nas btio	scalapack	Revision: Channet: Group: Type: Cluater: Results ID: Builder Loc: Running Loc	tion: ation:	513832168461a70001e3359090000e653 gen2 basic_1 rrpich2 r 50027658 /home1utbot/imapich2/install QA-PATCHESs /home1utbot/imapich2/install QA-PATCHESs /homenutbot/imapich2/install QA-PATCHESs	nasteribasic/613d63216/da 651a70d041e33390900c0ci 651a70d041e33390900c0ci	161a70 5e53b4 5e53b4
unted Runs 1 Groops ; / Types collectives alignitherv	Total Runs 971 compilation NoA	Test List (1399	imb4	Succ 70.06	ess Rate % gen2 intel 513683 21 22 4 6 513683	mpibench	Lost Rate 10.49%	mpich2 cyclic	nas NeA	Running Ra 0.1% nas btio	scalapack 513a83	Revision: Channet: Group: Type: Cluster: Results ID: Builder Loc: Running Loc Log Locatio Results Loc	tion: ation: 1: tion:	513832168661a70001e33590900c0ce533 gen2 basic_1 rrpich2 rl 50027858 /homenunbothmapich2/installCA-PATCHESi /homenunbothmapich2/installCA-PATCHESi /homenunbothmapich2/install30321656 /homenunbothmapich2/insulfs1308321656	naster/basic/61365321653 6112704041e33360900c0ci	161a700 be53b4
Inted Runs	Total Runs 971 compilation Not	Test List (1399 imb	Lount	Succ 70.06	ess Rate % gen2 intel 513000 21 22 4 513000 21 513000 21 4 513000 21 22 22 22 22	mpibench NIA	Lost Rate 10.49% mpich2 513453 513453	Failure Ra 19.34% mpich2 cyclic \$13683	nas NVA	Running Ra 0.1% nas blo NA	scalapack \$13d33 \$13d33	Revision: Channel: Group: Type: Cluster: Results ID: Builder Loc: Running Lo: Log Locatio Results Loc Owner(s):	tion: ation: 1: tion:	513832168618/12/00-11433590900-00-be53 gen2 basic_1 rrgin12 r 55027558 /homenunbothmapich2/installCA-PATCHESis /homenunbothmapich2/installCA-PATCHESis /homenunbothmapich2/installSA21668 /homenunbothmapich2/installS321668	nasteribasic0 154652 1664 6 1 a7040 4 1e33390500.00d	i61a70 be53b4
1 Croops // Types - Collectives alignitier 1 Collectives alignitier 1	Total Runs 971 compilation NoA NA	Test List (1399 imb 513653	51363	Succ 70.06 imb4 cuda N/A	ess Rate % gen2 intel 513003 21 22 4 6 513003 21 22 4 6	mpibench N/A	Lost Rate 10.49% snpich2 513483	Failure Ra 19.34% mplch2 cycle \$13883	nas NA NA	Running Ra 0.1% nas bilo NA NA	scalapack \$13883	Revision: Channel: Group: Type: Cluster: Results ID: Builder Loc: Running Lo Log Locatio Results Loc Owner(s): Hosts:	tion: ation: :: tion:	5138321686181270041e33590900c0cbe53 gen2 basic_1 mpich2 n 55027585 //omenunbotimapich2/installCA-PATCHESis //omenunbotimapich2/installCA-PATCHESis //omenunbotimapich2/installCA-PATCHESis //omenunbotimapich2/installSA121686 //omenunbotimapich2/installSA121686 //omenunbotimapich2/installSA121686	nasteribasic0 154532 1654 6 5 a70404 te33390500.00d	161a70x be53b4 be53b4
1 Crops / Types - Collectives aligather 1 Collectives aligather 1	Total Runs 971 compilation No. No.	Test List (1399 imb 51363 51363	51363	Succ 70.06	ess Rate % gen2 intel 513683 21 22 4 513683 21 21 22 4 6 513683 21 4 6 513683 21	mpibench NIA NIA	Lost Rate 10.49% mpick2 51343 51343	Failure Ra 19.34% mpich2 cyclic \$13883	nas NeA NeA	Running Ra 0.1% nas bilo NA NA	scalapack \$1383 \$1363	Revision: Channel: Group: Type: Cluster: Results ID: Builder Loc: Running Lo: Log Lostio Results Loc Owner(s): Hosts: Start Time:	tion: ation: 1: tion:	51383216861812/00-11433350900c/cbe53 gen2 basic_1 mpin2 n 55027568 Anomenunbotimwapich2installQA-PATCHE50 Anomenunbotimwapich2installQA-PATCHE50 Anomenunbotimwapich2installQA-PATCHE50 Anomenunbotimwapich2installQA-PATCHE50 Anomenunbotimwapich2installQA-PATCHE50 Anomenunbotimwapich2installQA-PATCHE50 Anomenunbotimwapich2installQA-PATCHE50 Anomenunbotimwapich2installQA-PATCHE50 Jan 5, 2020, 6,24 p.m.	144 naster/basic/6 158532 1683 16 1 a70404 1e03380500.cbc/ 6 1 a70404 1e03380500.cbc/	+61a70 be53b4 be53b4
Ceoups :/ Types → Collectives aligather 1 Collectives aligather 2	Total Runs 971 971 NoA NoA NoA	Test List (1399 imb 51363 51363 51363	51363 51363	Succ 70.06 imb-t cuda N/A N/A	ess Rate % gen2 inter 21 22 4 513603 21 22 4 6 513603 21 22 4 6 513603 21 22 22 4 6	mpibench NA NA	Lost Rate 10.49% meich2 51363 51363 51363	Failure Ra 19.34% mpich2 cyclic \$13683 \$13683	118 1183 1183 1183 1183 1183 1183 1183	Running Ra 0.1% mas bitio NVA NVA	te scalapack 513683 513683	Revision: Channel: Group: Type: Cluster: Results ID: Builder Loc: Running Loc Log Locatio Results Loc Owner(s): Hosts: Start Time: End Time:	tion: ation: 1: tion:	5138321686132700416335909000066533 gen2 basic_1 mpin2 n 55027568 homenunbotimwapich2installGA.PATCHESIs homenunbotimwapich2installGA.PATCHESIs homenunbotimwapich2installGA.PATCHESIs homenunbotimwapich2installGA.PATCHESIs homenunbotimwapich2installGA.PATCHESIs homenunbotimwapich2installGA.PATCHESIs homenunbotimwapich2installGA.PATCHESIs homenunbotimwapich2installGA.PATCHESIs homenunbotimwapich2installGA.PATCHESIs homenunbotimwapich2installGA.PATCHESIs homenunbotimwapich2installGA.PATCHESIs homenunbotimwapich2installGA.PATCHESIs homenunbotimwapich2installGA.PATCHESIs homenunbotimwapich2installGA.PATCHESIs homenunbotimwapich2installGA.PATCHESIs homenunbotimwapich2installGA.PATCHESIs homenunbotimwapich2installGA.PATCHESIs homenunbotimwapich2installGA.PATCHESIs homenunbotimwapich2installGA.PATCHESIs homenunbotimwapich2installGA.PATCHESis homenunbotimwapich2installGA.PA	naster/basic/61365321653 161120304163330500c/bci 611270404163330500c/bci	161a70d be53b44

Summary of all tests for one commit

Summary of an individual test

-PATCHES/master gen2 mpich2 basic 1

Details of individual combinations in one test

myapich2 myapich2-git OA-PATCHES/master Branch Grid 513d83 Rev Grid 597858 Result

CFLAG	8
Part	ENV
1	IVV2_DEBUG_SHOW_BACKTRACE+1 IVV2_CKPT_USE_AGGREGATION=0 MV2_USE_U0_HYBRID=0 MV2_CN_DEMAND_THRESHOLD=1 MV2_USE_UD_ZCOPY=0 PRIVIL RSH=1 MV2_USE_R0MA_CM=1 MV2_CPU_BINDING_LEVEL=SOCKET MV2_CPU_BINDING_POLICY=SCATTER.MV2_USE_BITORIC COMM_SUT1=1 MV2_ COMM_SEVIL_THRESHOLD=1 MV2_USE_PRIORITY_FACTOResH_UD_UBRANY_PATH=logtportox0/26_S0m/optioned_25_bit logtporte_01 0_branes_2011_132innux/completenbinteR4-logtIntel/2017completes_und_branes_2011.1125innux/institutioner84-logtCattes.01.182-logtcata6_01864/opticata6_01
2	MV2_DEBUG_SHOW_BACKTRACE+1W2_OKPLUSE_ADGREANTONO MV2_USE_UD_UHRIDHO MV2_ON_DEMAND_IHRESHICH MV2_USE_UD_ZOOPVO- RINN, RSH-HW2_CPL BANDAL CHULLS-SOCKET MV2_CPL BINDING-FOLVEXATTER MV2_USE_BITONG_COMM. SHLT1 MV2_BITONLCOMM.SHLT1 OLDH 1MV2_BIMP_PRIORITY_FACTORH-LD_UBRAVY_PATH-oppposobub/2 6 onto-/optionor2 4/to-optionor2 4/to optionor2 4/
3	MV2_DEBUG_SHOW_BACKTRACE1 MV2_OWT_USE_ADBREATIONO MV2_USE_U0_THISRID-MV2_ON_EDAMAD_THEESHOLD-MV2_USE_U0_ZOOPH-0 EX_TAC1-USE_MIRPUR, ISH-IH VV2_CUP UNIXING_EVEL-SOCKT IVX_CUP UNIXING_FOUND-SACTER MV2_USE_BITOINC_COMM_SEVET-MV2_DIS M_SEVET_THESHOLD-MV2_SMP_PRIORITY_ACTOR-64 L0_LIBMAY_PATH-teptpreduct/2.6.00% job/sptbs/20.5.00% job/sptbs
4	MVD_DEBUG_SHOW_BACKTRACE+1 MVD_CHXT_USE_ACREGATIONED MV2_USE_UD_THRRITHE MV2_USE_UD_ZCXXY+VU EUE_MMXR_INSH+1 MV2_USE_DCX 1 MVD_CH2_IDMOIL_SEVEL-SCXXET MV2_CH2_IDMOIDS_POLCH*SCXTER.MV2_USE_IDTXXCC_CMM_SPLT+1 MV2_USE_DCX_CMM_SPLT+1 THRESHCLD SMP_PRORTY_FACTOR+4 LD_UBRAVY_XH+Noptprotox26 S00%/pdbband_268/sozpcx01 01864/ uptimized17/complex_and_banks_2017.1 132/muuch binnelse1.aptimized17/complex_and_banks_2017.1 132/muuch/s01466 5 0166-9/dbband_5066_D042-666.abt/s014-666.abt/s
5	MVD_DEBIO_SHOW_BACKTRACE+1 MVD_OVEF_USE_ADGREGATION 0 MV2_USE_UD_IVHRID=0 MVD_USE_UD_ZCOCY+0 USE_URANN, SH+1 MVD_OVE_UBE LEVEL+SOCKET MVD_CFU_DIBIONO_DUCHSACTER MVD_USE_BTONC_COMA SHIT-1 MVD_UST_DTONC_COMA SHIT-1 MVD_USE_PRORITY R=64 LD_UBBACY_PATH-optonebuZ2.6 000/pathoon2.56tt; coptos0.5 0.01664/ vpbtuse207/complete_and_planete_2017.1328/muu/complex/biorde54/gottese pplane_and_pathone_2017.1328/muu/tabatheel64 vpbtuse2.6 1006/cpbtuse2.6 000/
6	IAV2_DEBUG_SHOW_BACKTRAGE+1AV2_CAPT_USE_ADGREGATION-014V2_SUP_USE_UIACC+0.MV2_USE_UD_HIRBID=014V2-CV_DEMAND_THRESHLDI- SE_UD_ZCOP*+0USE_UIARUN_RSH+1AV2_USE_ADDA_CA+1AV2_CPU_BINDING_LEVEL+SCORT_HIV2_CPU_BINDING_POLITY-ATT-HOPMONT_ HIV2_USE_STANDAC_COMM_SPT1+1AV2_STORC_COMM_SPT1_THRESHLDI-UI-HIV2_SUP_POLITY_AFTCH-RSHLDI-UU_BINDING_POLITY-ATT-HOPMONT_SATT AUX2_USE_STANDAC_COMM_SPT1+1AV2_STORC_COMM_SPT1_THRESHLDI-UI-HIV2_SUP_POLITY_AFTCH-RSHLDI_UU_BINDING_POLITY-ATT-HOPMONT_SATT AUX2_USE_STANDAC_COMM_SPT1+1AV2_USE_ADDAG_STT1_THRESHLDI-UI-HIV2_SUP_POLITY_AFTCH-RSHLD_UU_BINDING AUX2_STANDAC_COMM_SPT1+1AV2_USE_ADDAG_STITT-HISTORICAUCH-HIV2_SUP_POLITY_AFTCH-RSHLD_UU_BINDING AUX2_STANDACCASAT_USE_ADDAG_ADDAG_STANDACCASATURAT_STANDACCASATURAT_COMM_SPT1-FITCH-RSHLDI-UU_BINDING AUX2_STANDACCASAT_USE_ADDAG_ADDAG_STANDACCASATURAT_STANTACCASATURAT_STANTACCASA
7	IND_DEBUG_SHOW_BACKTRAGET IND_CART_USE_ADGREGATION-DIAV2_SUP_USE_UNICE-DIAV2_USE_UD_HINDIGE/DIAV2_USE_PRI_BRANKENT INV2_USE SE_UD_ZCOPY-O USE_INFRUM_RSH-HINZ_COPJ_BINONG_LEVEL-SOCIET INV2_CPU_BINONG_PUBLY-SACTER INV2_USE_PRI_BRANKENT INV2_USE_ DOML_SPLTT_UN_BITOINC_COMM_SIV_2 SIV_THRESHOUTT INV2_SUP_PRIORITY_FACTOR-ELD_USERXP1_PATH-optprotoxit/26_Status/protect DE4/aptitet/2017/complexs_und_granes_2017.1152/InvavicempletBinIndE4/aptitet/2017/complexs_und_Bhanes_2017.1152/InvavireNtBinIndE4/aptitet/2017/complexs_und_granes_2017.1152/InvavireNtBinIndE4/aptitet/2017/complexs_und_granes_2017.1152/InvavireNtBinIndE4/aptitet/2017/complexs_und_granes_2017.1152/InvavireNtBinIndE4/aptitet/2017/complexs_und_granes_2017.1152/InvavireNtBinIndE4/aptitet/2017/complexs_und_granes_2017.1152/InvavireNtBinIndE4/aptitet/2017/complexs_und_granes_2017.1152/InvavireNtBinIndE4/aptitet/2017/complexs_und_granes_2017.1152/InvavireNtBinIndE4/aptitet/2017/complexs_und_granes_2017.1152/InvavireNtBinIndE4/aptitet/2017/complexs_und_granes_2017.1152/InvavireNtBinIndE4/aptitet/2017/complexs_und_granes_2017.1152/InvavireNtBinIndE4/aptitet/2017/complexs_und_granes_2017.1152/InvavireNtBinIndE4/aptitet/2017/complexs_und_granes_2017.1152/InvavireNtBinIndE4/aptitet/2017/complexs_und_granes_2017.1152/InvavireNtBinIndE4/aptitet/2017/complexs_und_granes_2017/1152/InvavireNtBinIndE4/aptitet/2017/complexs_und_granes_2017/1152/InvavireNtBinIndE4/aptitet/2017/complexs_und_granes_2017/1152/InvavireNtBinIndE4/aptitet/2017/complexs_und_granes_2017/1152/InvavireNtBinIndE4/aptitet/2017/complexs_und_granes_2017/1152/InvavireNtBinIndE4/aptitet/2017/complexs_und_granes_2017/1152/InvavireNtBinIndE4/aptitet/2017/complexs_und_granes_2017/1152/InvavireNtBinIndE4/aptitet/2017/complexs_und_granes_2017/1152/InvavireNtBinIndE4/aptitet/2017/complexs_und_granes_2017/1152/InvavireNtBinIndE4/aptitet/2017/complexs_und_granes_2017/1152/InvavireNtBinIndE4/aptitet/2017/complexs_und_granes_2017/1152/InvavireNtBinIndE4/aptitet/2017/complexs_und_granes_2017/1152/InvavireNtBinIndE4/a
8	INZ_DEBUG_SHOW_BACKTRACE+1 INZ_CKPT_USE_AGGREGATION-0 INZ_SMP_USE_LIMIC2+0 INZ_USE_UD_HYBRD=0 INZ_ON_DBMAND_THRESHOLD+1 SE_UD_ZCOP*+0 USE_IMPRUN_RSH+1 INZ_USE_R0MA_CM=1 INZ_OPU_BINING_LEVE-ISOCKET INZ_OPU_BINING_POLI*>CATTER INZ_USE_BTONU M_SPLT1=1 INZ_DTONIC_COMM_RSPLT_THRESHOLD+1 INZ_SMP_R0R0TH_TACTOR+10_LID_RARY_TATH-redptore00028_0 INIto hybrbiox02 Statio hybrbiox02

5

Scripts to Determine Performance Regression

- Automated method to identify performance regression between different commits
- Tests different MPI primitives
 - Point-to-point; Collectives; RMA
- Works with different
 - Job Launchers/Schedulers
 - SLURM, PBS/Torque, JSM
 - Works with different interconnects
- Works on multiple HPC systems
- Works on CPU-based and GPU-based systems

Performance regression of mvapich2-2.3rc2-x-3e5551 and mvapich2-masterx-2950c8 on FRONTERA (cascadelake architecture) Thu Aug 15 09:23:48 CDT 2019

OLD_TUNEVAR= NEW_TUNEVAR= Legend Dark Green : Performance of mvapich2-masterx-2950c8 is more than 5 % better than mvapich2-2.3rc2-x-3e5551

Inter-node

Light Green : Performance of mvapich2-masterx-2950c8 is less than 5 % better than mvapich2-2.3rc2-x-3e5551

Grey : Performance of mvapich2-masterx-2950c8 is same as mvapich2-2.3rc2-x-3e5551

Light Red : Performance of mvapich2-masterx-2950c8 is less than 5 % worse compared to mvapich2-2.3rc2-x-3e5551

 $Dark \ Red: Performance \ of \ mvapich 2-master x-2950c8 \ is \ more \ than \ 5 \ \% \ worse \ compared \ to \ mvapich 2-2.3rc2-x-3e5551$

2 14.79 14.52 -1.96 2 8.41 8.64 2.96 2 2.30 2.30 0.96 2 2	4 29.56 29.05 -1 % 4 17.11 17.43 1 % 4 2.30 2.30 0 %	8 58.98 58.02 -1 % 8 34.13 34.84 2 % 8 2.31 2 2 1	16 117.73 116.10 -1 % 16 68.41 69.74 1 % 16 16	32 219.70 221.33 0 % 32 127.38 132.60 3 %	64 438.70 442.21 0 % 64 253.67 265.58	128 862.87 852.53 -1 % 128 494.66	256 1821.56 1795.21 -1 % 256 1051 14	512 3505.96 3539.99 0 % 512	1K 6728.97 6714.72 0 % 1K	2K 11697.93 11634.75 0 % 2K	4K 15297.26 15186.84 0 %	8K 18873.10 18919.87 0 %
2 8.41 8.64 2.96 2.30 2.30 0.96 2 2	4 17.11 17.43 1% 4 2.30 2.30 0%	8 34.13 34.84 2.% 8 2.31 2.31	16 68.41 69.74 1 % 16	32 127.38 132.60 3 %	64 253.67 265.58	128 494.66	256 1051-14	512	1K	2K	4K	9L
2 2.30 2.30 0 %	4 2.30 2.30 0 %	8 2.31 2.31	16		4 %	509.46 2 %	1080.49 2 %	1987.26 2079.78 4 %	3801.76 3978.30 4 %	6357.54 6438.37 1 %	8422.30 8478.65 0 %	9972.75 10031.73 0 %
2		0%	2.47 2.50 -1 %	32 2.51 2.53 0 %	64 2.91 2.93 0 %	128 2.97 2.99 0 %	256 3.09 3.11 0 %	512 3.33 3.35 0 %	1K 3.82 3.85 0 %	2K 4.79 4.81 0 %	4K 7.03 6.96 0 %	8K 10.76 10.79 0 %
1.97 1.96 0 %	4 1.96 1.95 0 %	8 1.97 1.95 1 %	16 1.97 1.95 1 %	32 1.97 1.95 1 %	64 1.97 1.95 1 %	128 2.01 2.01 0 %	256 2.07 2.07 0 %	512 2.11 2.13 0 %	1K 2.23 2.24 0 %	2K 2.51 2.50 0 %	4K 3.07 3.12 -1 %	8K 3.79 3.76 0 %
2 1.58 1.59 0 %	4 1.59 1.58 0 %	8 1.59 1.58 0 %	16 1.59 1.59 0 %	32 1.63 1.63 0 %	64 1.62 1.62 0 %	128 1.64 1.64 0 %	256 1.98 1.97 0 %	512 2.03 2.02 0 %	1K 2.12 2.11 0 %	2K 2.33 2.32 0 %	4K 2.85 2.87 0 %	8K 3.58 3.62 -1 %
1 1.12 1.12 0 %	2 1.12 1.12 0 %	4 1.12 1.12 0 %	8 1.12 1.12 0 %	16 1.16 1.16 0 %	32 1.16 1.16 0 %	64 1.18 1.18 0 %	128 1.23 1.23 0 %	256 1.64 1.63 0 %	512 1.73 1.72 0 %	1K 1.89 1.90 0 %	2K 2.27 2.26 0 %	4K 3.30 3.16 4 %
2 8.84 12.41 28 %	4 17.67 24.90 29 %	8 35.26 50.43 30 %	16 69.95 102.34 31 %	32 139.10 209.77 33 %	64 272.80 441.25 38 %	128 527.61 906.82 41 %	256 1014.76 1649.39 38 %	512 1906.89 2984.43 36 %	1K 3512.06 5576.41 37 %	2K 5983.62 7606.62 21 %	4K 8808.00 2974.96 -196 %	8K 12102.21 11577.63 -4 %
2 8.65 11.37 23 %	4 18.55 22.70 18 %	8 37.12 45.17 17 %	16 70.05 88.18 20 %	32 140.56 169.93 17 %	64 277.59 361.68 23 %	128 535.34 759.80 29 %	256 992.70 1471.40 32 %	512 1834.85 2513.06 26 %	1K 3403.25 3894.56 12 %	2K 5539.04 5642.22 1 %	4K 7134.90 7440.25 4 %	8K 9246.59 9152.05 -1 %
2 8.72 4652819.64 2.55 5672494.10 17 %	4 4660101.58 5656039.25 17 %	8 4651994.45 5642174.30 17 %	16 4408316.86 5511754.90 20 %	32 4386219.34 5312786.53 17 %	64 4343210.34 5624117.68 22 %	128 4181690.38 5945467.49 29 %	256 3961422.35 5710061.67 30 %	512 3582486.69 4904271.86 26 %	1K 3358575.31 3815154.75 11 %	2K 2709562.74 2757052.57 1 %	4K 1751091.52 1816130.44 3 %	8K 1128687.3 1117241.9 -1 %
82	0 * 3 2 * 3 1 * 5 0 * 5 1 * 5 0 * 5 2 * 6 1 * 12 0 * 5 2 * 6 8 * 8 1 * 2 0 * 5 2 * 6 8 * 6 1 * 2 2 * 6 8 * 6 1 * 3 2 * 6 8 * 6 1 * 3 2 * 6 2 * 7 2	0 0 0 0 1 1 3 4 50 1 50 1 1.59 1.58 0 9 0 9 1	0 0 0 0 1 1 0 2 4 5 1 5 0 5 0 1	$\begin{array}{c c c c c c c c c c c c c c c c c c c $	$\begin{array}{ c c c c c c c c c c c c c c c c c c c$	$\begin{array}{c c c c c c c c c c c c c c c c c c c $	$\begin{array}{ c c c c c c c c c c c c c c c c c c c$	$ \begin{array}{c c c c c c c c c c c c c c c c c c c $	$ \begin{array}{ c c c c c c c c c c c c c c c c c c c$	$ \begin{array}{ c c c c c c c c c c c c c c c c c c c$	$ \begin{array}{ c c c c c c c c c c c c c c c c c c c$	0 0 0 0 0 1

osu_allgather	1	2	4	8	16	32	64	128	256	512	1K	2K	4K	8K	16K	32K	64K	128K	256K	512K	1M
	21.83	23.72	22.11	24.40	28.29	38.99	28.70	47.85	66.45	105.64	914.13	2330.43	726.91	1197.84	2460.79	3996.50	6274.14	13385.82	20382.71	43029.04	102436.87
	17.34	16.16	17.28	17.87	19.89	21.65	26.24	43.90	63.03	102.18	177.28	276.96	502.60	1272.69	2555.89	3247.50	5516.63	10428.17	20449.20	43000.04	88526.72
	20 %	31 %	21 %	26 %	29 %	44 %	8 %	8 %	5 %	3 %	80 %	88 %	30 %	-6 %	-3 %	18 %	12 %	22 %	0 %	0 %	13 %
osu_allgatherv	1	2	4	8	16	32	64	128	256	512	1K	2K	4K	8K	16K	32K	64K	128K	256K	512K	1M
	25.66	27.61	25.41	27.68	31.34	40.98	78.11	133.69	207.27	505.08	529.16	558.27	689.62	1176.09	2461.61	3898.46	6573.11	11721.43	21588.64	43244.00	102924.40
	21.91	22.04	19.47	20.73	23.41	25.06	29.32	50.28	133.06	233.32	239.41	284.32	525.36	1254.16	2524.56	3167.98	5468.24	10026.13	20060.80	43166.45	87885.28
	14 %	20 %	23 %	25 %	25 %	38 %	62 %	62 %	35 %	53 %	54 %	49 %	23 %	-6 %	-2 %	18 %	16 %	14 %	7 %	0 %	14 %
osu_allreduce	4	8	16	32	64	128	256	512	1K	2K	4K	8K	16K	32K	64K	128K	256K	512K	1M	2M	4M
	36.82	32.90	32.66	33.15	35.42	36.90	41.44	51.83	69.81	107.22	48.50	51.63	86.36	120.66	184.02	296.64	521.36	1996.51	3221.70	6665.53	9984.70
	19.08	22.28	19.37	19.37	21.29	27.06	22.96	29.01	37.88	60.36	44.66	48.42	84.79	109.31	172.09	267.65	469.93	967.96	2649.74	6305.41	8972.08
	48 %	32 %	40 %	41 %	39 %	26 %	44 %	44 %	45 %	43 %	7 %	6 %	1 %	9 %	6 %	9 %	9 %	51 %	17 %	5 %	10 %

Designing (MPI+X) for Exascale

- Scalability for million to billion processors
 - Support for highly-efficient inter-node and intra-node communication (both two-sided and one-sided)
- Scalable Collective communication
 - Offloaded
 - Non-blocking
 - Topology-aware
- Balancing intra-node and inter-node communication for next generation multi-/many-core (128-1024 cores/node)
 - Multiple end-points per node
- Support for efficient multi-threading
- Integrated Support for GPGPUs and Accelerators
- Fault-tolerance/resiliency
- QoS support for communication and I/O
- Support for Hybrid MPI+PGAS programming
 - MPI + OpenMP, MPI + UPC, MPI + OpenSHMEM, CAF, MPI + UPC++...
- Virtualization
- Energy-Awareness

MVAPICH2 Software Family

Requirements	Library
MPI with Support for InfiniBand, Omni-Path, Ethernet/iWARP and, RoCE (v1/v2)	MVAPICH2
Optimized Support for Microsoft Azure Platform with InfiniBand	MVAPICH2-Azure
Advanced MPI features/support (UMR, ODP, DC, Core-Direct, SHArP, XPMEM), OSU INAM (InfiniBand Network Monitoring and Analysis),	MVAPICH2-X
Advanced MPI features (SRD and XPMEM) with support for Amazon Elastic Fabric Adapter (EFA)	MVAPICH2-X-AWS
Optimized MPI for clusters with NVIDIA GPUs and for GPU-enabled Deep Learning Applications	MVAPICH2-GDR
Energy-aware MPI with Support for InfiniBand, Omni-Path, Ethernet/iWARP and, RoCE (v1/v2)	MVAPICH2-EA
MPI Energy Monitoring Tool	OEMT
InfiniBand Network Analysis and Monitoring	OSU INAM
Microbenchmarks for Measuring MPI and PGAS Performance	ОМВ

Overview of MVAPICH2 Features

- Job start-up
- Transport Type Selection
- Process Mapping and Point-to-point Intra-node Protocols
- Collectives

Towards High Performance and Scalable Startup at Exascale



Job Startup Performance

• Near-constant MPI and OpenSHMEM initialization time at any process count

 10x and 30x improvement in startup time of MPI and OpenSHMEM respectively at 16,384 processes

 Memory consumption reduced for remote endpoint information by O(processes per node)

• 1GB Memory saved per node with 1M processes and 16 processes per node

(a) **On-demand Connection Management for OpenSHMEM and OpenSHMEM+MPI.** S. Chakraborty, H. Subramoni, J. Perkins, A. A. Awan, and D K Panda, 20th International Workshop on High-level Parallel Programming Models and Supportive Environments (HIPS '15)

(b) **PMI Extensions for Scalable MPI Startup.** S. Chakraborty, H. Subramoni, A. Moody, J. Perkins, M. Arnold, and D K Panda, Proceedings of the 21st European MPI Users' Group Meeting (EuroMPI/Asia '14)

C O Non-blocking PMI Extensions for Fast MPI Startup. S. Chakraborty, H. Subramoni, A. Moody, A. Venkatesh, J. Perkins, and D K Panda, 15th IEEE/ACM International Symposium on Cluster, Cloud and Grid Computing (CCGrid '15)

(e) SHMEMPMI – Shared Memory based PMI for Improved Performance and Scalability. S. Chakraborty, H. Subramoni, J. Perkins, and D K Panda, 16th IEEE/ACM International Symposium on Cluster, Cloud and Grid Computing (CCGrid '16)

Startup Performance on TACC Frontera

MPI_Init on Frontera



- MPI_Init takes 3.9 seconds on 57,344 processes on 1,024 nodes
- All numbers reported with 56 processes per node

New designs available since MVAPICH2-2.3.2

How to Get the Best Startup Performance with MVAPICH2?

- MV2_HOMOGENEOUS_CLUSTER=1
- MV2_ON_DEMAND_UD_INFO_EXCHANGE=1

//Set for homogenous clusters

//Enable UD based address exchange

Using SLURM as launcher

• Use PMI2

- ./configure --with-pm=slurm --with-pmi=pmi2
- srun --mpi=pmi2 ./a.out

• Use PMI Extensions

- Patch for SLURM available at <u>http://mvapich.cse.ohio-state.edu/download/</u>
- Patches available for SLURM 15, 16, and 17
- PMI Extensions are automatically detected by MVAPICH2

Using mpirun_rsh as launcher

• MV2_MT_DEGREE

 degree of the hierarchical tree used by mpirun_rsh

• MV2_FASTSSH_THRESHOLD

#nodes beyond which hierarchical-ssh scheme is used

• MV2_NPROCS_THRESHOLD

#nodes beyond which file-based communication
 is used for hierarchical-ssh during start up

Transport Protocol Selection in MVAPICH2



Number of Processes

• Both UD and RC/XRC have benefits

- Hybrid for the best of both
- Enabled by configuring MVAPICH2 with the –enable-hybrid
- Available since MVAPICH2 1.7 as integrated interface

Parameter	Significance	Default	Notes
MV2_USE_UD_HYBRID	 Enable / Disable use of UD transport in Hybrid mode 	Enabled	• Always Enable
MV2_HYBRID_ENABLE_THRESHOLD_SIZE	 Job size in number of processes beyond which hybrid mode will be enabled 	1024	 Uses RC/XRC connection until job size < threshold
MV2_HYBRID_MAX_RC_CONN	 Maximum number of RC or XRC connections created per process Limits the amount of connection memory 	64	 Prevents HCA QP cache thrashing

- Refer to Running with Hybrid UD-RC/XRC section of MVAPICH2 user guide for more information
- <u>http://mvapich.cse.ohio-state.edu/static/media/mvapich/mvapich2-2.3a-userguide.html#x1-690006.11</u>

Process Mapping support in MVAPICH2



MVAPICH2 detects processor architecture at job-launch

Preset Process-binding Policies – Bunch

- "Core" level "Bunch" mapping (Default)
 - MV2_CPU_BINDING_POLICY=bunch



- "Socket/Numanode" level "Bunch" mapping
 - MV2_CPU_BINDING_LEVEL=socket MV2_CPU_BINDING_POLICY=bunch



Preset Process-binding Policies – Scatter

- "Core" level "Scatter" mapping
 - MV2_CPU_BINDING_POLICY=scatter



- "Socket/Numanode" level "Scatter" mapping
 - MV2_CPU_BINDING_LEVEL=socket MV2_CPU_BINDING_POLICY=scatter



Process and thread binding policies in hybrid MPI+Threads

- A new process binding policy "hybrid"
 - MV2_CPU_BINDING_POLICY = hybrid
- A new environment variable for co-locating Threads with MPI Processes
 - MV2_THREADS_PER_PROCESS = k
 - Automatically set to OMP_NUM_THREADS if OpenMP is being used
 - Provides a hint to the MPI runtime to spare resources for application threads.
- New variable for threads bindings with respect to parent process and architecture
 - MV2_HYBRID_BINDING_POLICY= {bunch|scatter|linear|compact|spread|numa}
 - Linear binds MPI ranks and OpenMP threads sequentially (one after the other)
 - Recommended to be used on non-hyper threaded systems with MPI+OpenMP
 - Compact binds MPI rank to physical-core and locates respective OpenMP threads on hardware threads
 - Recommended to be used on multi-/many-cores e.g., KNL, POWER8, and hyper-threaded Xeon, etc.

Binding Example in Hybrid (MPI+Threads)

- MPI Processes = 4, OpenMP Threads per Process = 4
- MV2_CPU_BINDING_POLICY = hybrid
- MV2_THREADS_PER_PROCESS = 4



- Detects hardware-threads support in architecture
- Assigns MPI ranks to physical cores and respective OpenMP Threads to HW threads

Binding Example in Hybrid (MPI+Threads) ---- Cont'd

- MPI Processes = 4, OpenMP Threads per Process = 4
- MV2_CPU_BINDING_POLICY = hybrid
- MV2_THREADS_PER_PROCESS = 4
- MV2_THREADS_BINDING_POLICY = linear



MPI Rank-0 with its 4-OpenMP threads gets bound on Core-0 through Core-3, and so on

Binding Example in Hybrid (MPI+Threads) ---- Cont'd

- MPI Processes = 16
- Example: AMD EPYC 7551 processor with 8 NUMA domains
- MV2_CPU_BINDING_POLICY = hybrid
- MV2_HYBRID_BINDING_POLICY = numa



User-Defined Process Mapping

- User has complete-control over process-mapping
- To run 4 processes on cores 0, 1, 4, 5:
 - \$ mpirun_rsh -np 4 -hostfile hosts MV2_CPU_MAPPING=0:1:4:5 ./a.out
- Use ',' or '-' to bind to a set of cores:
 - \$mpirun_rsh -np 64 -hostfile hosts MV2_CPU_MAPPING=0,2-4:1:5:6 ./a.out
- Is process binding working as expected?
 - MV2_SHOW_CPU_BINDING=1
 - Display CPU binding information
 - Launcher independent
 - Example
 - MV2_SHOW_CPU_BINDING=1 MV2_CPU_BINDING_POLICY=scatter

-----CPU AFFINITY-----

RANK:0 CPU_SET: 0

RANK:1 CPU_SET: 8

- Refer to Running with Efficient CPU (Core) Mapping section of MVAPICH2 user guide for more information
- http://mvapich.cse.ohio-state.edu/static/media/mvapich/mvapich2-2.3rc1-userguide.html#x1-600006.5

Collective Communication in MVAPICH2



Run-time flags:

All shared-memory based collectives : MV2_USE_SHMEM_COLL (Default: ON) Hardware Mcast-based collectives : MV2_USE_MCAST (Default : OFF) CMA and XPMEM-based collectives are in MVAPICH2-X

Hardware Multicast-aware MPI_Bcast on TACC Stampede



- MCAST-based designs improve latency of MPI_Bcast by up to **85%**
- Use MV2_USE_MCAST=1 to enable MCAST-based designs

MPI_Scatter - Benefits of using Hardware-Mcast



• Enabling MCAST-based designs for MPI_Scatter improves small message up to **75%**

	Parameter	Description	Default
	MV2_USE_MCAST = 1	Enables hardware Multicast features	Disabled
	enable-mcast	Configure flag to enable	Enabled
Network Based	Computing Laboratory	MUG'20	

Offloading with Scalable Hierarchical Aggregation Protocol (SHArP)

- Management and execution of MPI operations in the network by using SHArP
 - Manipulation of data while it is being transferred in the switch network
- SHArP provides an abstraction to realize the reduction operation
 - Defines Aggregation Nodes (AN), Aggregation Tree, and Aggregation Groups
 - AN logic is implemented as an InfiniBand Target Channel Adapter (TCA) integrated into the switch ASIC *
 - Uses RC for communication between ANs and between AN and hosts in the Aggregation Tree *

More details in the tutorial "SHARPv2: In-Network Scalable Streaming Hierarchical Aggregation and Reduction Protocol" by Devendar Bureddy (NVIDIA/Mellanox)



* Bloch et al. Scalable Hierarchical Aggregation Protocol (SHArP): A Hardware Architecture for Efficient Data Reduction

Benefits of SHARP Allreduce at Application Level



Avg DDOT Allreduce time of HPCG

More details in the talk "Impact of SHARP and Adaptive Routing on Applications on Frontera" by John Cazes (TACC) on Wednesday (08/26/2020) from 12:00 PM - 12:30 PM EDT

SHARP support available since MVAPICH2 2.3a

Parameter	Description	Default		
MV2_ENABLE_SHARP=1	Enables SHARP-based collectives	Disabled		
enable-sharp	Configure flag to enable SHARP	Disabled		

- Refer to Running Collectives with Hardware based SHARP support section of MVAPICH2 user guide for more information
- http://mvapich.cse.ohio-state.edu/static/media/mvapich/mvapich2-2.3-userguide.html#x1-990006.26

MVAPICH2 Software Family

Requirements	Library
MPI with Support for InfiniBand, Omni-Path, Ethernet/iWARP and, RoCE (v1/v2)	MVAPICH2
Optimized Support for Microsoft Azure Platform with InfiniBand	MVAPICH2-Azure
Advanced MPI features/support (UMR, ODP, DC, Core-Direct, SHArP, XPMEM), OSU INAM (InfiniBand Network Monitoring and Analysis),	MVAPICH2-X
Advanced MPI features (SRD and XPMEM) with support for Amazon Elastic Fabric Adapter (EFA)	MVAPICH2-X-AWS
Optimized MPI for clusters with NVIDIA GPUs and for GPU-enabled Deep Learning Applications	MVAPICH2-GDR
Energy-aware MPI with Support for InfiniBand, Omni-Path, Ethernet/iWARP and, RoCE (v1/v2)	MVAPICH2-EA
MPI Energy Monitoring Tool	OEMT
InfiniBand Network Analysis and Monitoring	OSU INAM
Microbenchmarks for Measuring MPI and PGAS Performance	ОМВ

MVAPICH2-X for MPI and Hybrid MPI + PGAS Applications

High Performance Parallel Programming Models							
MPI PGAS Hybrid MPI + X							
	Message Passing Interface	(UPC, OpenSHMEM, CAF, UPC++)	(MPI + PGAS + OpenMP/Cilk)				

High Performance and Scalable Unified Communication Runtime											
Diverse APIs and Mechanisms											
Optimized Point- to-point Primitives	Optimized Point- to-point Primitives Remote Memory Access Active Messages Active Messages (Blocking and Non-Blocking) Non-Blocking) Scalable Job Startup Fault Tolerance Introspection & Analysis with OSU INAM										
Support fo (Infin	r Modern Netwo niBand, iWARP, RoCE	rking Technologie , Omni-Path)	:S	Supp	oort for Modern I (Intel-X	Multi-/Many-cor eon, OpenPower)	e Architectures				

- Current Model Separate Runtimes for OpenSHMEM/UPC/UPC++/CAF and MPI
 - Possible deadlock if both runtimes are not progressed
 - Consumes more network resource
- Unified communication runtime for MPI, UPC, UPC++, OpenSHMEM, CAF
 - Available with since 2012 (starting with MVAPICH2-X 1.9)
 - <u>http://mvapich.cse.ohio-state.edu</u>

MVAPICH2-X 2.3

- Released on 06/08/2020
- Major Features and Enhancements
 - MPI Features
 - Based on MVAPICH2 2.3.4
 - OFA-IB-CH3, OFA-IB-RoCE, PSM-CH3, and PSM2-CH3 interfaces
 - Enhanced point-to-point and collective tunings for AMD EPYC, Catalyst@EPCC, Mayer@Sandia, Auzre@Microsoft, AWS, and Frontera@TACC
 - MPI (Advanced) Features
 - Optimized support for large message MPI_Allreduce and MPI_Reduce
 - OFA-IB-CH3 and OFA-IB-RoCE interfaces
 - Improved performance for communication using DC transport
 - OFA-IB-CH3 interface
 - Enhanced support for AWS EFA adapter and SRD transport protocol
 - OFA-IB-CH3 interface
 - Enhanced point-to-point and collective tuning for AWS EFA adapter and SRD transport protocol
 - OFA-IB-CH3 interface
 - Add multiple MPI_T PVARs and CVARs for point-to-point and collective operations

- Tuning for MPI collective operations for Intel Broadwell, Intel CascadeLake, Azure HB (AMD EPYC), and Azure HC (Intel Skylake) systems
 - OFA-IB-CH3, OFA-IB-RoCE, PSM-CH3, and PSM2-CH3 interface
- Support for OSU InfiniBand Network Analysis and Management (OSU INAM) Tool v0.9.6
- Unified Runtime Features
 - Based on MVAPICH2 2.3.4 (OFA-IB-CH3 interface). All the runtime features enabled by default in OFA-IB-CH3 and OFA-IB-RoCE interface of MVAPICH2 2.3.4 are available in MVAPICH2-X 2.3

MVAPICH2-X Feature Table

Features for InfiniBand (OFA-IB-CH3) and RoCE (OFA-RoCE-CH3)	Basic	Basic-XPMEM	Intermediate	Advanced
Architecture Specific Point-to-point and Collective Optimizations for x86, OpenPOWER, and ARM	\checkmark	~	~	
Optimized Support for PGAS models (UPC, UPC++, OpenSHMEM, CAF) and Hybrid MPI+PGAS models	\checkmark	~	~	~
CMA-Aware Collectives	\checkmark	\checkmark	\checkmark	\checkmark
Optimized Asynchronous Progress*	\checkmark	\checkmark	 	\checkmark
InfiniBand Hardware Multicast-based MPI_Bcast*+	\checkmark	\checkmark	\checkmark	\checkmark
OSU InfiniBand Network Analysis and Monitoring (INAM)*+				\checkmark
XPMEM-based Point-to-Point and Collectives		\checkmark	\checkmark	\checkmark
Direct Connected (DC) Transport Protocol*+			\checkmark	\checkmark
User mode Memory Registration (UMR)*+				\checkmark
On Demand Paging (ODP)*+				\checkmark
Core-direct based Collective Offload*+				\checkmark
SHARP-based Collective Offload*+				\checkmark

- * indicates disabled by default at runtime. Must use appropriate environment variable in MVAPICH2-X user guide to enable it.
- + indicates features only tested with InfiniBand network

Overview of MVAPICH2-X Features

- Direct Connect (DC) Transport
 - Available from MVAPICH2-X 2.3rc1 onwards
- CMA-based Collectives
 - Available from MVAPICH2-X 2.3rc1 onwards
- Asynchronous Progress
 - Available from MVAPICH2-X 2.3rc1 onwards
- XPMEM-based Reduction Collectives
 - Available from MVAPICH2-X 2.3rc1 onw ards
- XPMEM-based Non-reduction Collectives
 - Available from MVAPICH2-X 2.3rc2 onw ards

Impact of DC Transport Protocol on Neuron

Neuron with YuEtAl2012



Overhead of RC protocol for connection establishment and communication

- Up to 76% benefits over MVAPICH2 for Neuron using Direct Connected transport protocol at scale
 - VERSION 7.6.2 master (f5a1284) 2018-08-15
- Numbers taken on bbpv2.epfl.ch
 - Knights Landing nodes with 64 ppn
 - ./x86_64/special -mpi -c stop_time=2000 -c is_split=1 parinit.hoc
 - Used "runtime" reported by execution to measure performance
- Environment variables used
 - MV2_USE_DC=1
 - MV2_NUM_DC_TGT=64
 - MV2_SMALL_MSG_DC_POOL=96
 - MV2_LARGE_MSG_DC_POOL=96
 - MV2_USE_RDMA_CM=0

Available from MVAPICH2-X 2.3rc2 onwards

More details in talk "Building Brain Circuits: Experiences with shuffling terabytes of data over MPI" by Matthias Wolf, Blue Brain Project, EPFL, Switzerland on Wednesday (08/26/2020) from 1:30 PM - 2:00 PM EDT.

Optimized CMA-based Collectives for Large Messages



Performance of MPI_Gather on KNL nodes (64PPN)

- Significant improvement over existing implementation for Scatter/Gather with 1MB messages (up to 4x on KNL, 2x on Broadwell, 14x on OpenPower)
- New two-level algorithms for better scalability
- Improved performance for other collectives (Bcast, Allgather, and Alltoall)

S. Chakraborty, H. Subramoni, and D. K. Panda, Contention Aware Kernel-Assisted MPI Collectives for Multi/Many-core Systems, IEEE Cluster '17, BEST Paper Finalist Available since MVAPICH2-X 2.3b

Benefits of the New Asynchronous Progress Design: Broadwell + InfiniBand



<u>P3DFFT</u>

High Performance Linpack (HPL)



Up to 33% performance improvement in P3DFFT application with <u>448 processes</u> Up to 29% performance improvement in HPL application with <u>896 processes</u>

A. Ruhela, H. Subramoni, S. Chakraborty, M. Bayatpour, P. Kousha, and D.K. Panda, "Efficient design for MPI Asynchronous Progress without Dedicated Resources", Parallel Computing 2019

Shared Address Space (XPMEM)-based Collectives Design



- "Shared Address Space"-based true zero-copy Reduction collective designs in MVAPICH2
- Offloaded computation/communication to peers ranks in reduction collective operation
- Up to **4X** improvement for 4MB Reduce and up to **1.8X** improvement for 4M AllReduce

J. Hashmi, S. Chakraborty, M. Bayatpour, H. Subramoni, and D. Panda, Designing Efficient Shared Address Space Reduction Available since MVAPICH2-X 2.3rc1 Collectives for Multi-/Many-cores, International Parallel & Distributed Processing Symposium (IPDPS '18), May 2018.

Performance of Non-Reduction Collectives with XPMEM



- 28 MPI Processes on single dual-socket Broadwell E5-2680v4, 2x14 core processor
- Used osu_bcast from OSU Microbenchmarks v5.5
Application Level Benefits of XPMEM-based Designs

CNTK AlexNet Training (B.S=default, iteration=50, ppn=28)

MiniAMR (dual-socket, ppn=16)



- Intel XeonCPU E5-2687W v3 @ 3.10GHz (10-core, 2-socket)
- Up to 20% benefits over IMPI for CNTK DNN training using AllReduce
- Up to 27% benefits over IMPI and up to 15% improvement over MVAPICH2 for MiniAMR application kernel

MVAPICH2 Software Family

Requirements	Library
MPI with Support for InfiniBand, Omni-Path, Ethernet/iWARP and, RoCE (v1/v2)	MVAPICH2
Optimized Support for Microsoft Azure Platform with InfiniBand	MVAPICH2-Azure
Advanced MPI features/support (UMR, ODP, DC, Core-Direct, SHArP, XPMEM), OSU INAM (InfiniBand Network Monitoring and Analysis),	MVAPICH2-X
Advanced MPI features (SRD and XPMEM) with support for Amazon Elastic Fabric Adapter (EFA)	MVAPICH2-X-AWS
Optimized MPI for clusters with NVIDIA GPUs and for GPU-enabled Deep Learning Applications	MVAPICH2-GDR
Energy-aware MPI with Support for InfiniBand, Omni-Path, Ethernet/iWARP and, RoCE (v1/v2)	MVAPICH2-EA
MPI Energy Monitoring Tool	OEMT
InfiniBand Network Analysis and Monitoring	OSU INAM
Microbenchmarks for Measuring MPI and PGAS Performance	ОМВ

MPI + CUDA - Naive

• Data movement in applications with standard MPI and CUDA interfaces

At Sender:

cudaMemcpy(s_hostbuf, s_devbuf, . . .); MPI_Send(s_hostbuf, size, . . .);

At Receiver:

MPI_Recv(r_hostbuf, size, . . .); cudaMemcpy(r_devbuf, r_hostbuf, . . .);

High Productivity and Low Performance



MPI + CUDA - Advanced

• Pipelining at user level with non-blocking MPI and CUDA interfaces

At Sender:

```
for (j = 0; j < pipeline\_len; j++)
   cudaMemcpyAsync(s_hostbuf + j * blk, s_devbuf + j *
     blksz, ...);
for (j = 0; j < pipeline\_len; j++) {
     while (result != cudaSucess) {
         result = cudaStreamQuery(...);
         if(j > 0) MPI Test(...);
      MPI_Isend(s_hostbuf + j * block_sz, blksz . . .);
```

MPI_Waitall();

<<Similar at receiver>>

Low Productivity and High Performance



GPU-Aware (CUDA-Aware) MPI Library: MVAPICH2-GPU

- Standard MPI interfaces used for unified data movement
- Takes advantage of Unified Virtual Addressing (>= CUDA 4.0)
- Overlaps data movement from GPU with RDMA transfers



CUDA-Aware MPI: MVAPICH2-GDR 1.8-2.3.4 Releases

- Support for MPI communication from NVIDIA GPU device memory
- High performance RDMA-based inter-node point-to-point communication (GPU-GPU, GPU-Host and Host-GPU)
- High performance intra-node point-to-point communication for multi-GPU adapters/node (GPU-GPU, GPU-Host and Host-GPU)
- Taking advantage of CUDA IPC (available since CUDA 4.1) in intra-node communication for multiple GPU adapters/node
- Optimized and tuned collectives for GPU device buffers
- MPI datatype support for point-to-point and collective communication from GPU device buffers
- Unified memory

MVAPICH2-GDR: Pre-requisites for OpenPOWER & x86 Systems

- MVAPICH2-GDR 2.3.4 requires the following software to be installed on your system:
 - 1. Mellanox OFED 3.2 and later
 - 2. NVIDIA Driver 367.48 or later
 - 3. NVIDIA CUDA Toolkit 7.5 and later
 - 4. NVIDIA Peer Memory (nv peer mem) module to enable GPUDirect RDMA (GDR) support
- Strongly Recommended for Best Performance
 - 5. GDRCOPY Library by NVIDIA: <u>https://github.com/NVIDIA/gdrcopy</u>
- Comprehensive Instructions can be seen from the MVAPICH2-GDR User Guide:
 - <u>http://mvapich.cse.ohio-state.edu/userguide/gdr/</u>

MVAPICH2-GDR: Download and Setup on OpenPOWER & x86 Systems

- Simple Installation steps for both systems
- Pick the right MVAPICH2-GDR RPM from Downloads page:
 - <u>http://mvapich.cse.ohio-state.edu/downloads/</u>
 - e.g. <u>http://mvapich.cse.ohio-state.edu/download/mvapich/gdr/2.3/mofed4.5/mvapich2-gdr-mcast.cuda10.0.mofed4.5.gnu4.8.5-2.3-1.el7.x86_64.rpm</u> (== <mv2-gdr-rpm-name>.rpm)
- \$ wget http://mvapich.cse.ohio-state.edu/download/mvapich/gdr/2.3/<mv2-gdr-rpm-name>.rpm <u>Root Users:</u>

\$ rpm -Uvh --nodeps <mv2-gdr-rpm-name>.rpm

Non-Root Users:

- \$ rpm2cpio <mv2-gdr-rpm-name>.rpm | cpio id
- Contact MVAPICH help list with any questions related to the package <u>mvapich-help@cse.ohio-state.edu</u>

ROCE and Optimized Collectives Support

- RoCE V1 and V2 support
- RDMA_CM connection support
- CUDA-Aware Collective Tuning
 - Point-point Tuning (available since MVAPICH2-GDR 2.0)
 - Tuned thresholds for the different communication patterns and features
 - Depending on the system configuration (CPU, HCA and GPU models)
 - Tuning Framework for GPU based collectives
 - Select the best algorithm depending on message size, system size and system configuration
 - Support for Bcast and Gather operations for different GDR-enabled systems
- Available since MVAPICH2-GDR 2.2RC1 release

MVAPICH2-GDR 2.3.4

- Released on 06/04/2020
- Major Features and Enhancements
 - Based on MVAPICH2 2.3.4
 - Enhanced MPI_Allreduce performance on DGX-2 systems
 - Enhanced MPI_Allreduce performance on POWER9 systems
 - Reduced the CUDA interception overhead for non-CUDA symbols
 - Enhanced performance for point-to-point and collective operations on Frontera's RTX nodes
 - Add new runtime variable 'MV2_SUPPORT_DL' to replace 'MV2_SUPPORT_TENSOR_FLOW'
 - Added compilation and runtime methods for checking CUDA support
 - Enhanced GDR output for runtime variable MV2_SHOW_ENV_INFO
 - Tested with Horovod and common DL Frameworks (TensorFlow, PyTorch, and MXNet)
 - Tested with PyTorch Distributed
 - Support for CUDA 10.1
 - Support for PGI 20.x
 - Enhanced GPU communication support in MPI_THREAD_MULTIPLE mode
 - Enhanced performance of datatype support for GPU-resident data
 - Zero-copy transfer when P2P access is available between GPUs through NVLink/PCIe
 - Enhanced GPU-based point-to-point and collective tuning
 - OpenPOWER systems such as ORNL Summit and LLNL Sierra ABCI system @AIST, Owens and Pitzer systems @Ohio Supercomputer Center
 - Scaled Allreduce to 24,576 Volta GPUs on Summit

Tuning GDRCOPY Designs in MVAPICH2-GDR

Parameter	Significance	Default	Notes
MV2_USE_GDRCOPY	• Enable / Disable GDRCOPY- based designs	1 (Enabled)	• Always enable
MV2_GDRCOPY_LIMIT	 Controls messages size until which GDRCOPY is used 	8 KByte	 Tune for your system GPU type, host architecture. Impacts the eager performance
MV2_GPUDIRECT_GDR COPY_LIB	 Path to the GDRCOPY library 	Unset	• Always set
MV2_USE_GPUDIRECT_ D2H_GDRCOPY_LIMIT	 Controls messages size until which GDRCOPY is used at sender 	16Bytes	 Tune for your systems CPU and GPU type

- Refer to Tuning and Usage Parameters section of MVAPICH2-GDR user guide for more information
- <u>http://mvapich.cse.ohio-state.edu/userguide/gdr/#_tuning_and_usage_parameters</u>

Tuning Loopback Designs in MVAPICH2-GDR

Parameter	Significance	Default	Notes
MV2_USE_GPUDIRECT_ LOOPBACK	 Enable / Disable LOOPBACK-based designs 	1 (Enabled)	 Always enable
MV2_GPUDIRECT_LOO PBACK_LIMIT	 Controls messages size until which LOOPBACK is used 	8 KByte	 Tune for your system GPU type, host architecture and HCA. Impacts the eager performance Sensitive to the P2P issue

• Refer to Tuning and Usage Parameters section of MVAPICH2-GDR user guide for more information

• <u>http://mvapich.cse.ohio-state.edu/userguide/gdr/#_tuning_and_usage_parameters</u>

Tuning GPUDirect RDMA (GDR) Designs in MVAPICH2-GDR

Parameter	Significance	Default	Notes
MV2_USE_GPUDIRECT	 Enable / Disable GDR-based designs 	1 (Enabled)	• Always enable
MV2_GPUDIRECT_LIMIT	 Controls messages size until which GPUDirect RDMA is used 	8 KByte	 Tune for your system GPU type, host architecture and CUDA version: impact pipelining overheads and P2P bandwidth bottlenecks
MV2_USE_GPUDIRECT_ RECEIVE_LIMIT	 Controls messages size until which 1 hop design is used (GDR Write at the receiver) 	256KBytes	 Tune for your system GPU type, HCA type and configuration

- Refer to Tuning and Usage Parameters section of MVAPICH2-GDR user guide for more information
- <u>http://mvapich.cse.ohio-state.edu/userguide/gdr/#_tuning_and_usage_parameters</u>

Application-Level Evaluation (HOOMD-blue)

64K Particles

3500 2500 3000 Average Time Steps per second MV2 MV2+GDR Average Time Steps per second (TPS) 2500 2000 **2X 2X** 2000 1500 (TPS) 1500 1000 1000 500 500 0 0 8 16 32 4 8 16 32 4 Number of Processes Number of Processes

256K Particles

- Platform: Wilkes (Intel Ivy Bridge + NVIDIA Tesla K20c + Mellanox Connect-IB)
- HoomDBlue Version 1.0.5
 - GDRCOPY enabled: MV2_USE_CUDA=1 MV2_IBA_HCA=mlx5_0 MV2_IBA_EAGER_THRESHOLD=32768 MV2_VBUF_TOTAL_SIZE=32768 MV2_USE_GPUDIRECT_LOOPBACK_LIMIT=32768 MV2_USE_GPUDIRECT_GDRCOPY=1 MV2_USE_GPUDIRECT_GDRCOPY_LIMIT=16384

MPI Datatype support in MVAPICH2

- Datatypes support in MPI
 - Operate on customized datatypes to improve productivity
 - Enable MPI library to optimize non-contiguous data

At Sender:

```
MPI_Type_vector (n_blocks, n_elements, stride, old_type, &new_type);
MPI_Type_commit(&new_type);
```

... MPI Send(s buf, size, new type, dest, tag, MPI COMM WORLD);

- Inside MVAPICH2
 - Use datatype specific CUDA Kernels to pack data in chunks
 - Efficiently move data between nodes using RDMA
 - In progress currently optimizes *vector* and *hindexed* datatypes
 - Transparent to the user

H. Wang, S. Potluri, D. Bureddy, C. Rosales and D. K. Panda, GPU-aware MPI on RDMA-Enabled Clusters: Design, Implementation and Evaluation, IEEE Transactions on Parallel and Distributed Systems, Accepted for Publication.

MPI Datatype Processing (Computation Optimization)

- Comprehensive support
 - Targeted kernels for regular datatypes vector, subarray, indexed_block
 - Generic kernels for all other irregular datatypes
- Separate non-blocking stream for kernels launched by MPI library
 - Avoids stream conflicts with application kernels
- Flexible set of parameters for users to tune kernels
 - Vector
 - MV2_CUDA_KERNEL_VECTOR_TIDBLK_SIZE
 - MV2_CUDA_KERNEL_VECTOR_YSIZE
 - Subarray
 - MV2_CUDA_KERNEL_SUBARR_TIDBLK_SIZE
 - MV2_CUDA_KERNEL_SUBARR_XDIM
 - MV2_CUDA_KERNEL_SUBARR_YDIM
 - MV2_CUDA_KERNEL_SUBARR_ZDIM
 - Indexed_block
 - MV2_CUDA_KERNEL_IDXBLK_XDIM

Performance of Stencil3D (3D subarray)

Stencil3D communication kernel on 2 GPUs with various X, Y, Z dimensions using MPI_Isend/Irecv

- DT: Direct Transfer, TR: Targeted Kernel
- Optimized design gains up to 15%, 15% and 22% compared to TR, and more than 86% compared to DT on X, Y and Z respectively





MPI Datatype Processing (Communication Optimization)

Common Scenario

```
MPI_Isend (A,.. Datatype,...)
MPI_Isend (B,.. Datatype,...)
MPI_Isend (C,.. Datatype,...)
MPI_Isend (D,.. Datatype,...)
...
```

MPI_Waitall (...);

*A, B...contain non-contiguous MPI Datatype

Waste of computing resources on CPU and GPU



Application-Level Evaluation (Cosmo) and Weather Forecasting in Switzerland





- 2X improvement on 32 GPUs nodes
- 30% improvement on 96 GPU nodes (8 GPUs/node)

<u>Cosmo model: http://www2.cosmo-model.org/content</u> /tasks/operational/meteoSwiss/

On-going collaboration with CSCS and MeteoSwiss (Switzerland) in co-designing MV2-GDR and Cosmo Application

C. Chu, K. Hamidouche, A. Venkatesh, D. Banerjee, H. Subramoni, and D. K. Panda, Exploiting Maximal Overlap for Non-Contiguous Data Movement Processing on Modern GPU-enabled Systems, IPDPS'16

MVAPICH2-GDR: Enhanced Derived Datatype

- Kernel-based and GDRCOPY-based one-shot packing for inter-socket and inter-node communication
- Zero-copy (packing-free) for GPUs with peer-to-peer direct access over PCIe/NVLink



Deep Learning: New Challenges for Runtimes

- Scale-up: Intra-node Communication
 - Many improvements like:
 - NVIDIA cuDNN, cuBLAS, NCCL, etc.
 - CUDA 9 Co-operative Groups
- Scale-out: Inter-node Communication
 - DL Frameworks most are optimized for singlenode only
 - Distributed (Parallel) Training is an emerging trend
 - OSU-Caffe MPI-based
 - Microsoft CNTK MPI/NCCL2
 - Google TensorFlow gRPC-based/MPI/NCCL2
 - Facebook Caffe2 Hybrid (NCCL2/Gloo/MPI)
 - PyTorch



MVAPICH2 (MPI)-driven Infrastructure for ML/DL Training



MVAPICH2-GDR vs. NCCL2 – Allreduce Operation (DGX-2)

- **Optimized designs in MVAPICH2-GDR offer better/comparable performance for most cases**
- MPI_Allreduce (MVAPICH2-GDR) vs. ncclAllreduce (NCCL2) on 1 DGX-2 node (16 Volta GPUs)



Platform: Nvidia DGX-2 system (16 Nvidia Volta GPUs connected with NVSwitch), CUDA 10.1

C.-H. Chu, P. Kousha, A. Awan, K. S. Khorassani, H. Subramoni and D. K. Panda, "NV-Group: Link-Efficient Reductions for Distributed Deep Learning on Modern Dense GPU Systems, "ICS-2020, June-July 2020.

MVAPICH2-GDR: MPI_Allreduce at Scale (ORNL Summit)

- Optimized designs in MVAPICH2-GDR offer better performance for most cases
- MPI_Allreduce (MVAPICH2-GDR) vs. ncclAllreduce (NCCL2) up to 1,536 GPUs



C.-H. Chu, P. Kousha, A. Awan, K. S. Khorassani, H. Subramoni and D. K. Panda, "NV-Group: Link-Efficient Reductions for Distributed Deep Learning on Modern Dense GPU Systems, "ICS-2020, June-July 2020.

Distributed Training with TensorFlow and MVAPICH2-GDR

• ResNet-50 Training using TensorFlow benchmark on 1 DGX-2 node (16 Volta GPUs)



Platform: Nvidia DGX-2 system (16 Nvidia Volta GPUs connected with NVSwitch), CUDA 9.2

MUG'20

Distributed Training with TensorFlow and MVAPICH2-GDR

ResNet-50 Training using TensorFlow benchmark on 400 SUMMIT -- 1536 Volta ImageNet-1k has 1.2 million images 350 (Thousands) GPUs! 300 MVAPICH2-GDR reaching ~0.35 million 250 images per second for ImageNet-1k! 1,281,167 (1.2 mil.) images Image per second 200 150 Time/epoch = 3.6 seconds 100 50 Total Time (90 epochs) 0 $= 3.6 \times 90 = 332$ seconds = 1 2 4 6 12 24 48 96 192 384 768 1536 Number of GPUs 5.5 minutes! NCCL-2.4 MVAPICH2-GDR-2.3.2

*We observed errors for NCCL2 beyond 96 GPUs

Platform: The Summit Supercomputer (#1 on Top500.org) – 6 NVIDIA Volta GPUs per node connected with NVLink, CUDA 9.2

Distributed TensorFlow on TACC Frontera (2048 CPU nodes)

- Scaled TensorFlow to 2048 nodes on Frontera using MVAPICH2 and IntelMPI
- MVAPICH2 and IntelMPI give similar performance for DNN training
- Report a peak of 260,000 images/sec on 2048 nodes
- On 2048 nodes, ResNet-50 can be trained in 7 minutes!





Applications-Level Tuning: Compilation of Best Practices

- MPI runtime has many parameters
- Tuning a set of parameters can help you to extract higher performance
- Compiled a list of such contributions through the MVAPICH Website
 - <u>http://mvapich.cse.ohio-state.edu/best_practices/</u>
- Initial list of applications
 - Amber
 - HoomDBlue
 - HPCG
 - Lulesh
 - MILC
 - Neuron
 - SMG2000
 - Cloverleaf
 - SPEC (LAMMPS, POP2, TERA_TF, WRF2)
- Soliciting additional contributions, send your results to mvapich-help at cse.ohio-state.edu.
- We will link these results with credits to you.

Amber: Impact of Tuning Eager Threshold



Data Submitted by: Dong Ju Choi @ UCSD

- Tuning the Eager threshold has a significant impact on application performance by avoiding the synchronization of rendezvous protocol and thus yielding better communication computation overlap
- 19% improvement in overall execution time at256 processes
- Library Version: MVAPICH2 2.2
- MVAPICH Flags used
 - MV2_IBA_EAGER_THRESHOLD=131072
 - MV2_VBUF_TOTAL_SIZE=131072
- Input files used
 - Small: MDIN
 - Large: <u>PMTOP</u>

MiniAMR: Impact of Tuning Eager Threshold



MiniAMR

- Tuning the Eager threshold has a significant impact on application performance by avoiding the synchronization of rendezvous protocol and thus yielding better communication computation overlap
- 8% percent reduction in total communication time
- Library Version: MVAPICH2 2.2
- MVAPICH Flags used
 - MV2_IBA_EAGER_THRESHOLD=32768
 - MV2_VBUF_TOTAL_SIZE=32768

Data Submitted by Karen Tomko @ OSC and Dong Ju Choi @ UCSD

SMG2000: Impact of Tuning Transport Protocol



- UD-based transport protocol selection benefits the SMG2000 application
- 22% and 6% on 1,024 and 4,096 cores, respectively
- Library Version: MVAPICH2 2.1
- MVAPICH Flags used
 - MV2_USE_ONLY_UD=1
- System Details
 - Stampede@ TACC
 - Sandybridge architecture with dual 8-cores nodes and ConnectX-3 FDR network

Neuron: Impact of Tuning Transport Protocol



- UD-based transport protocol selection benefits the SMG2000 application
- 15% and 27% improvement is seen for 768 and 1,024 processes respectively
- Library Version: MVAPICH2 2.2
- MVAPICH Flags used
 - MV2_USE_ONLY_UD=1
- Input File
 - <u>YuEtAl2012</u>
- System Details
 - Comet@SDSC
 - Haswell nodes with dual 12-cores socket per node and Mellanox FDR (56 Gbps) network.

HPCG: Impact of Collective Tuning for MPI+OpenMP Programming Model



- Partial subscription nature of hybrid MPI+OpenMP programming requires a new level of collective tuning
 - For PPN=2 (Processes Per Node), the tuned version of MPI_Reduce shows 51% improvement on 2,048 cores
- 24% improvement on 512 cores
 - 8 OpenMP threads per MPI processes
- Library Version: MVAPICH2 2.1
- MVAPICH Flags used
 - The tuning parameters for hybrid MPI+OpenMP programming models is on by default from MVAPICH2-2.1 onward
- System Details
 - Stampede@ TACC
 - Sandybridge architecture with dual 8-cores nodes and ConnectX-3 FDR network

Data Submitted by Jerome Vienne and Carlos Rosales-Fernandez @ TACC

HPCG

HOOMD-blue: Impact of GPUDirect RDMA Based Tuning



- HOOMD-blue is a Molecular Dynamics simulation using a custom force field.
- GPUDirect specific features selection and tuning significantly benefit the HOOMD-blue application. We observe a factor of 2X improvement on 32 GPU nodes, with both 64K and 256K particles
- Library Version: MVAPICH2-GDR 2.2
- MVAPICH-GDR Flags used
 - MV2_USE_CUDA=1
 - MV2_USE_GPUDIRECT=1
 - MV2_GPUDIRECT_GDRCOPY=1
- System Details
 - Wilkes@Cambridge
 - 128 Ivybridge nodes, each node is a dual 6cores socket with Mellanox FDR

Application Scalability on Skylake and KNL with Omni-Path



Courtesy: Mahidhar Tatineni @SDSC, Dong Ju (DJ) Choi@SDSC, and Samuel Khuvis@OSC ---- Testbed: TACC Stampede2 using MVAPICH2-2.3b

Runtime parameters: MV2_SMPI_LENGTH_QUEUE=524288 PSM2_MQ_RNDV_SHM_THRESH=128K PSM2_MQ_RNDV_HFI_THRESH=128K

SPEC MPI 2007 Benchmarks: Broadwell + InfiniBand



MVAPICH2-X outperforms Intel MPI by up to 31%

Configuration: 448 processes on 16 Intel E5-2680v4 (Broadwell) nodes having 28 PPN and interconnected

with 100Gbps Mellanox MT4115 EDR ConnectX-4 HCA

Execution Time in (s)
MVAPICH2 – Plans for Exascale

- Performance and Memory scalability toward 1-10M cores
- Hybrid programming (MPI + OpenSHMEM, MPI + UPC, MPI + CAF ...)
 - MPI + Task*
- Enhanced Optimization for GPU Support and Accelerators
- Taking advantage of advanced features of Mellanox InfiniBand
 - Tag Matching*
 - Adapter Memory*
 - Bluefield based offload*
- Enhanced communication schemes for upcoming architectures
 - ROCm*
 - Intel Optane*
 - BlueField*
 - CAPI*
- Extended topology-aware collectives
- Extended Energy-aware designs and Virtualization Support
- Extended Support for MPI Tools Interface (as in MPI 3.0)
- Extended FT support
- Support for * features will be available in future MVAPICH2 Releases

For More Details on other MVAPICH2 Libraries/Features

• MPI_T Support

 More details in the talk "Performance Engineering using MVAPICH and TAU" by Sameer Shende (Paratools/UO) on Tuesday (08/25/2020) from 3:30 PM - 4:00 PM EDT

• MVAPICH2-Azure

 More details in the talk "MVAPICH2 on Microsoft Azure HPC" by Jithin Jose (Microsoft, Azure) on Tuesday (08/25/2020) from 1:30 PM - 2:00 PM EDT

• MVAPICH2-AWS

 More details in the talk "Scaling Message Passing on Amazon Web Services with Elastic Fabric Adapter" by Raghunath Rajachandrasekar (AWS) on Tuesday (08/25/2020) from 1:00 PM - 1:30 PM EDT

• Support for SHARP

 More details in the talk "Impact of SHARP and Adaptive Routing on Applications on Frontera" by John Cazes (TACC) on Wednesday (08/26/2020) from 12:00 PM - 12:30 PM EDT

Funding Acknowledgments

Funding Support by



Acknowledgments to all the Heroes (Past/Current Students and Staffs)

Current Students (Graduate)			Current Research Scientists	Current Post-docs
 Q. Anthony (Ph.D.) M. Bayatpour (Ph.D.) CH. Chu (Ph.D.) A. Jain (Ph.D.) M. Kedia (M.S.) 	 K. S. Khorassani (Ph.D.) P. Kousha (Ph.D.) N. S. Kumar (M.S.) B. Ramesh (Ph.D.) K. K. Suresh (Ph.D.) 	 N. Sarkauskas (Ph.D.) S. Srivastava (M.S.) S. Xu (Ph.D.) Q. Zhou (Ph.D.) 	 A. Shafi H. Subramoni <i>Current Senior Research Associate</i> J. Hashmi <i>Current Software Engineer</i> A. Reifsteck 	 M. S. Ghazimeersaeed K. Manian <i>Current Research Specialist</i> J. Smith
Past Students				
 A. Awan (Ph.D.) A. Augustine (M.S.) P. Balaji (Ph.D.) R. Biswas (M.S.) S. Bhagvat (M.S.) A. Bhat (M.S.) D. Buntinas (Ph.D.) 	 T. Gangadharappa (M.S.) K. Gopalakrishnan (M.S.) J. Hashmi (Ph.D.) W. Huang (Ph.D.) W. Jiang (M.S.) J. Jose (Ph.D.) S. Kini (M.S.) M. Koon (Ph.D.) 	 P. Lai (M.S.) J. Liu (Ph.D.) M. Luo (Ph.D.) A. Mamidala (Ph.D.) G. Marsh (M.S.) V. Meshram (M.S.) A. Moody (M.S.) S. Narayula (Ph.D.) 	 R. Rajachandrasekar (Ph.D.) D. Shankar (Ph.D.) G. Santhanaraman (Ph.D.) N. Sarkauskas (B.S.) A. Singh (Ph.D.) J. Sridhar (M.S.) S. Sur (Ph.D.) H. Subramoni (Ph.D.) 	 <i>Past Research Scientists</i> K. Hamidouche S. Sur X. Lu <i>Past Programmers</i> D. Bureddy J. Perkins
 L. Chai (Ph.D.) B. Chandrasekharan (M.S S. Chakraborthy (Ph.D.) N. Dandapanthula (M.S.) V. Dhanraj (M.S.) Past Post-Docs	.) – K. Kulkarni (M.S.) – R. Kumar (M.S.) – S. Krishnamoorthy (M.S.) – K. Kandalla (Ph.D.) – M. Li (Ph.D.)	 R. Noronha (Ph.D.) R. Noronha (Ph.D.) X. Ouyang (Ph.D.) S. Pai (M.S.) S. Potluri (Ph.D.) K. Raj (M.S.) 	 K. Vaidyanathan (Ph.D.) A. Vishnu (Ph.D.) J. Wu (Ph.D.) W. Yu (Ph.D.) J. Zhang (Ph.D.) 	Past Research Specialist – M. Arnold
 D. Banerjee X. Besseron HW. Jin 	– J. Lin – M. Luo – E. Mancini	 S. Marcarelli A. Ruhela I. Vienne 	– H. Wang	

J. Vienne

E. Mancini

Thank You!

panda@cse.ohio-state.edu



Network-Based Computing Laboratory <u>http://nowlab.cse.ohio-state.edu/</u>



The High-Performance MPI/PGAS Project http://mvapich.cse.ohio-state.edu/ Follow us on Twitter: @mvapich



High-Performance Big Data

The High-Performance Big Data Project <u>http://hibd.cse.ohio-state.edu/</u>



The High-Performance Deep Learning Project <u>http://hidl.cse.ohio-state.edu/</u>