

MPI usage at NERSC: Present and Future

Alice Koniges, Brandon Cook, Jack Deslippe, Thorston Kurth, Hongzhang Shan Lawrence Berkeley National Laboratory, Berkeley, CA USA



threading model (56%)

Abstract: We describe how MPI is used at the National Energy Research Scientific Computing Creter (NRSC), NEXRS is the production high-performance computing center for the US Department of Energy with more than 500 uses and 400 datation spreaks. We also compare the users are used to the strategiest and the strategiest and the strategiest and the strategiest and and the taking and the strategiest and exemptions to MPI are plausible and used to NEXRS (user, We also drockas percented strategiest and servicestors to MPI are plausible and used to NEXRS (user, Strategiest and extrategiest and exemptions) to MPI are plausible and used to NEXRS (user, Strategiest and extrategiest and exemptions). The MPI are plausible and used to NEXRS (user, Strategiest and extrategiest and exemptions) of MPI are plausible and used to NEXRS (user, Strategiest and extrategiest and the exemption of the plausible and used to NEXRS (user, Strategiest and extrategiest and the exemption of the strategiest and extrategiest and the extinct NEXRS (user, Strategiest) and the extinct NEXRS (user, Strategiest) and the extinct NEXRS (user, Strategiest) and NEXRS (User, Strategiest) How MPI is used at NERSC Portability issues wrt different MPI implementations In survey responses of 32 code groups, 10 gave the following comments on portability - New reduces the tell injuncementations seems to differ in two my inducement reastwine at a differ and the inducement of th Special thanks to Dr. Rolf Rabenseifner and colleagues at Stuttgart High Performance Computing Center for input in making the MPI NERSC User Survey. write at/). On NERSC's mad NERSC Hardware com regular mode and DMAEP mode). If the second sec ist supercomputer system (NERSC-8). It is named after American the first women to use a Nobel Drive in Discipleau or Medicine. amist Gerty Cori, ndary conditions Scatter/Gather AllToA Cori Phase 1 system is a Cray XC based on the intel Haswell multi-core processor. Cori Phase 2 is based the second generation of the Initial Xcon Phile Tamily on products, called the Knights Landing (KNL) Mary Integra Core (MIC) Architecture. Cori Phase I & Cori Phase II Integration begins - September 19, 2016 (Cr - 4 weaks). Comments/concerns/suggestions/requests wrt MPI There responses: Premarily had communication for integration and all-to-materiative IOne sided MPT Structure Communications, these ded communications one-sided and print-back between rank 0 and other Graph Intervensil / random access I masker PEC bit 10: generally this is minor amound of I deally all meaningful patterns are benchmarked Have you ever checked if your communication is serialized? urvey responses E code groups gave additional information cort have hopefit Lacd good gene uncore MM for Coly makes list of things hard. So many problems would have been soled much quicker with source. It would ago to document alconate gloridati protocial wholes an specific MM implementations (q.g. Con MM) for smill messages. That could equilan hypothol document alconate of a control of the sole of the specific MM implementations. Induced and the sole of the sole hypothol document control of the control of the sole of the specific MM is and the count of the sole o ions, threaded communications Phase III of a 200 compute nodes, 52,160 cores per nodes, 52,000 cores per nodes, 52,000 cores per node, 50,000 co $\overline{}$ Cray Ar Aggreg ment point-to-point comm Biocking MPI_Send / MPI_Recv Biocking MPI_Send / MPI_Recv Biocking MPI_Send / MPI_Recv Biocking MPI_Send Recv Edison has been the NERSC's primary supercomputer system for a number of years. It is named after U.S. Inventor and businessman Thomas Ava Edison. Do you overlap com and computation? Edison, a Cray XC30, has a peak performance of 2.57 petalfopulse, 133.824 compute cores for running scat-applications, 357 Terabytes of memory, and 7.59 Petala online disk storage with a peak I/O bandwidth of 168 -inahutes (IGB) per second. ::: the Edisorta Application codes used by MPI survey responders NIMROD NCAR LES GINGER FEL code WRF
 NWChem
 NWChem, CTF
 MLC DOD simulations colors

 Paralla Research Knuth
 epi58P
 PPS-38 (1646 element based M40 color)

 ALPSCOR
 Liam
 WWMPPCNAF - no applications of the second MPI (A.19) (Cat MPI (A -11 (35.5%) -11 (35.5%) **NERSC Research Areas and Users** WRF, UMWM, HYCOM, ESMF CSU global atmospheric model -10 (32.3%) Lattice OCD 241 Materials Sc Fusion Energy Climate 223 -8 (25.8%) **Responses from Survey of MPI Alternatives** 2015 NERSC Usage By Di Engineering 24 What error handling do you u do MPI Init ration? hat language(s) is(are) your application written in? Are you working on a proxy app for a major application? application? (Please name full application.) LDUNS)
LDUNS
LDUNS Default Predefined error handlers Chenistry 244 with MPI_Init with MPI_Init_thread computation is greatly simplified
 NEK5000 is the full app. At least, proper loca matrix integration is missing. ro MPI+Y2 +X? MPI-OpenMP MPI-OpenMP MPI-OpenAPC MPI-OpenAPC -0 (5%) Other -0 (5%) -0 (2) Other response: UPC, python, bash UPC, UPC+ Python Javas Cols. Python, PERL UPC R, haskell, ocard, python -18 (56.3%) 20 (62.5) used you to consider approaches other than MPI? new approaches other than MPI?
newson approaches other than MPI?
below approaches
below app Aldo Sadi' Na Pichica I Ingin Assess rines -13 (48.19 Ther responses: MP1 + TBB, C(++) But we prefer MP1 UPC, UPC++, GA 0 2 4 1 11 threads -only, possibly with shared memory optimizations (MP -11 -9 (33.3%) -6 (22.2%) lah.tre If you constant MP fullh anther communication model, e.g., Mr + Generation, Mr + Second M del, e.g., MPI + Co-Array Fortran, MPI + UPC, etc., which m —10 (37%) —8 (29.6%) 2015 NERSC Usage by Ins 2015 Acade mic Usage at NERSC Other responses: 0 - complex DAG-like workflow - compler transformations or - Binary code without access - Much less syntax to accem-10 12 processing vs. MPI MLLONS) Surgic University Talmet Univ National Univ U. Peer Lowis State U. Methods Scheel Mit U. Methods Scheel Mit U. Trees Anelle U. Deblever U. Deb 25 25 25 21 22 20 MET William & Mary U. Arkona UC Sin Dispo UC Sin delay U. Basia U- C U. Washington UC Sivies UCA U. Coloreste Der Prinzetze Unit Gallech U. Dicapo U. Biertacky Coloresta Divi - ----OTELES IN Original Constants UN If yes, which window creation method(s) MPI_Win_create MPI_Win_allocate MPI_Win_create_dynamic Funded to implement Load balancing MP1 is too low level. ::)/approach(s) do vou use? --8 (30.8%) 9 (34.6%) Charn Kok If yes, do you use overlapping w If you use overlapping windows, please specify the reason -0 (0%) -1 (3.8%) Snapshot of User Comparisons of MPI to other Use Casper papers I use Casper for RMA asynchronous progress which requires overlapped windows to offload user RMA communication to bedretravent a co-------3 (11.5%) -2 (7.7%) Programming Models at NERSC Clear responses: UP
 --8 (30.8%)
 --8 (30.8%)
 --8 (30.8%)
 --8 (30.8%)
 --8 (30.8%)
 --8 (30.8%)
 --8 (30.8%)
 --8 (30.8%)
 --8 (30.8%)
 --8 (30.8%)
 --8 (30.8%)
 --8 (30.8%)
 --8 (30.8%)
 --8 (30.8%)
 --8 (30.8%)
 --8 (30.8%)
 --8 (30.8%)
 --8 (30.8%)
 --8 (30.8%)
 --8 (30.8%)
 --8 (30.8%)
 --8 (30.8%)
 --8 (30.8%)
 --8 (30.8%)
 --8 (30.8%)
 --8 (30.8%)
 --8 (30.8%)
 --8 (30.8%)
 --8 (30.8%)
 --8 (30.8%)
 --8 (30.8%)
 --8 (30.8%)
 --8 (30.8%)
 --8 (30.8%)
 --8 (30.8%)
 --8 (30.8%)
 --8 (30.8%)
 --8 (30.8%)
 --8 (30.8%)
 --8 (30.8%)
 --8 (30.8%)
 --8 (30.8%)
 --8 (30.8%)
 --8 (30.8%)
 --8 (30.8%)
 --8 (30.8%)
 --8 (30.8%)
 --8 (30.8%)
 --8 (30.8%)
 --8 (30.8%)
 --8 (30.8%)
 --8 (30.8%)
 --8 (30.8%)
 --8 (30.8%)
 --8 (30.8%)
 --8 (30.8%)
 --8 (30.8%)
 --8 (30.8%)
 --8 (30.8%)
 --8 (30.8%)
 --8 (30.8%)
 --8 (30.8%)
 --8 (30.8%)
 --8 (30.8%)
 --8 (30.8%)
 --8 (30.8%)
 --8 (30.8%)
 --8 (30.8%)
 --8 (30.8%)
 --8 (30.8%)
 --8 (30.8%)
 --8 (30.8%)
 --8 (30.8%)
 --8 (30.8%)
 --8 (30.8%)
 --8 (30.8%)
 --8 (30.8%)
 --8 (30.8%)
 --8 (30.8%)
 --8 (30.8%)
 --8 (30.8%)
 --8 (30.8%)
 --8 (30.8%)
 --8 (30.8%)
 --8 (30.8%)
 --8 (30.8%)
 --8 (30.8%)
 --8 (30.8%)
 --8 (30.8%)
 --8 (30.8%)
 --8 (30.8%)
 --8 (30.8%)
 --8 (30.8%)
 --8 (30.8%)
 --8 (30.8%)
 --8 (30.8%)
 --8 (30.8%)
 --8 (30.8%)
 --8 (30.8%)
 --8 (30.8%)
 --8 (30.8%)
 --8 (30.8%)
 --8 (30.8%)
 --8 (30.8%)
 --8 (30.8%)
 --8 (30.8%)
 --8 (30.8%)
 --8 (30.8%)
 --8 (30.8%)
 --8 (30.8%)
 --8 (30.8%)
 --8 (30.8%)
 --8 (30.8%)
 --8 (30.8%)
 --8 (30.8%)
 --8 (30.8%)
 --8 (30.8%)
 --8 (30.8%)
 --8 (30.8%)
 --8 (30.8%)
 --8 (30.8%)
 --8 (30.8%)
 --8 (30.8%)
 --8 (30.8%)
 --8 (30.8%)
 --8 (30.8%)
 -2.00 -1 (3.8%) NERSC users are trying a variety of programming models and often comparing them to MPI. While these comparisons are obviously implementation and application dependent, we give here some current results. (See links in MPI Alternatives Survey for more results. -2 (11 5%) If yes to one-sided communication with normal windows, which RMA method(s)? -1 (3.8%) If yes to one-sided communication with windows, which RMA memory model? Bi-Directional bandwidth test on Cori Phase 1 using two processes, one per node MPI_Put MPI_Get ore than Other responses: • MPL_Fetch_and_op • MPL_Fetch_and_op • FOp and CAS If yes, which models? LIPC++: A PGAS Extension for C++ MPI and multiprocessing
 MPI and multiprocessing
 most of tham
 MPI and Coarray Fortran
 MPI and SHME
 benchmarks in multiple models
 MPI vs OCR MPI MPI MPI, OpenMP MPI, OpenMP ::: (MB/s) If yes to one-sided communication with normal windows, which synchronization methods? MPI_Win_post/star/complete/wait MPI_Win_post/star/complete/wait MPI_Win_post/star/communication.eg. MPI_Rout UPC++ palat_pints-atuaterts-para, size asym_nige(on, deal, even) Asym_ater(event, tender) fastar() golat_pints-p.pin_size() Bandwidth If yes, are there results available, e.g., paper, talk, link? http://ganet.bl.gov/performanca/ http://dx.doi.org/10.1070/787-3-319-41221-1_17, https://github.com/Pa https://dx.doi.org/10.1077/04242405655110 See http://gutec.bl.gov/publications/ bttp://gutec.bl.gov/publications/ Have you compared your approach to MPI? : See http://wp.cbl.gov/publications/ http://www.competer.org/scal/publications/ http://www.competer.org/scal/publications/ http://dx.doi.org/10.1109/IPDPSW.2014.83, http://www.gpi-alia.com/gpi2/banchr http://dx.doi.org/10.1109/IPDPSW.2014.83, http://www.gpi-alia.com/gpi2/banchr Zhang, Yil, et al. "UPC++: a PGAG Extension for C++ Parallel and Distributed Processing Symposium, 2014 Do you use MPI shared memory windows (allocated with MPI_Win_allocate_shared)? Message Size (Bytes) BoxLib using cray-mpich7.4.0 on 64 nodes normal language (C/ normal language (C/Fr If yes, on which architecture(s) was it optimized? rtran) assignments for "remote store" on the shared me MPI_Put and/or MPI_Get on the shared memory wi it optimized? • many - Intel CPU cluster, Cray, eti • x86_64 • x86 • x86 / InfiniBand Cray
 Cray
 Cray
 Cray XC (Edison/Cori)
 See individual papers for details BoxLib: An AMR Software Framework for building parallel adaptive mesh refinement (AMR) codes abving multi-acaleniumis-physics parall differential equations Available at https://github.com/BoxLib-Codes Managed by Winigun Zhang, et al. C++ and Fortan versions MPI + OpenMPI (+ UPC++) n method do you use? micn synchronization method do you use? MPL_Barrien jin MPI and sync all in CAF win sync or us, as appropriate MPL_Win_lock_all/unloca_all ush_jooral_all MPL_Barrien_bas MPL_SenseReev of empty message MPL_Wart MPL_WartAL MPL_Win_sync for shared memory windows multiple methods w do you think MPI could be improved to better address your needs? Much simpler programming model Active Message support. Hether FMA support fuult tearnor, ungenernation quality Vierage High Performance ParalleX (HPX)
* Liphowsight multi-threading
- Dukidis work into annalier tasks
- Increases concurrency
* Massing-driven computation
- Move work todat
- Keeps work todat
- Keeps work todat 256 512 1026 Number of Ranks 2018 use hints via the info obje If yes, do you aved some of the problems, but we have a paper in preparation to highlight the remainder. (efficiency), simpler API, and remote function execution MPI 3 has resolv. Lighter weight (ef Fault Tolerance Task-based N-Body using LibGeoDecomp associated with iated with commission of the cess_style, chu QuedRMEM on solid communation has an advantage one ton solid messaging for gold communation patterns where MP ten color messaging contracts proceeding unsequence messages, etc. can generate logit overhause. MP RMA could be used in place of OperGMEM if the performance of MP RMA wave optimised to where of SMEM. The second se Declarative criteria for
 Event driven
 Eliminates global barr yes, can you specify which MPL_File readwrte_at_alt() and readwrte_at() MPL_FILE_WRITE_ALL MPL_FILE_wRITE_ALL MPL_FILE_wat() and MPL_FILe_wrte_all write_at_all_write_read at MPL_FILE_WRITE_AT_ALL MPL_FILE_WRITE_AT_ALL Do you use MPI-IO? :sared name space - Global address space - Simplifies random gathers Does the application use dependent libraries? * Hyes, which libraries does it depend upon? • numpy.sdpy_adropy • ZenakQ and SKAR (https://ghtub.com/hyperiols) • COFUB Multiple
 Multiple
 BLAS, ScaLaPack
 HDF5 400 600 800 100 ____ Ř If yes, which data r "internal" "external32" Other If yes, which pos MiniGhost Weak Scaling Individual filep... Shared filepointers --th MPI_File_set_view) Explicit Offsets idual filepointers Does the application use Python? Can your appro * - 0,000 -MPI OpenMP P S M (40 variables - 20 time Other responses: • TBB • Starm • C++, Chapel, Legi 2000 1-Consolitional Frencation (consolition) Consolitional Proceeding (consolition) Consolitional Proceeding (consolition) Second Annual Consolition (consolition) Second Annual Consolition) Do you use an I/O package like HDF5? If yes, do you use, HDF5, NetCDF, or other 1900 :: 11 of 16 responses use HDF5
 7 of 16 responses use NetCDF
 No response specify other § 1000 HPX-5 is the High Performance ParalleX runtime library from Indiana University. (http://tpx.mov.met.ka.edu/ UNIX-3.edu/UNIX-Indiana.edu/Dentity Summary 500 Comments on Alternatives to MPI inding to our MPI Survey. We assume ou C's 5000+ users with 32 code groups resp ch tools, if any, do yo hich tools, if any, do you use to Allinea DDT Totalview, print statements Lots of print statements to dobug! write(",") for debugging :) TAU, own tools Intel Vtune, Allinea Forge ITAC ddt or debug the MPI portion of your application? gdb Capper for optimizing asynchronous progress of MPI RMA Sometimes VTune, Intel Trace Analyzer and Collector Compare results with setal code. CDR for debugging. Recordly stande performing instruction analysis using intel SDE In addition to UPC++ and HPX, NERSC users are using UPC, Charm++, Kokiso, Chapel, Coarary Fortran, OCR – based, Legion, Parallel Python, Hadsop/MapReduce, SHMEM, OpenMP, GASNet, Global Arrays, andGPI-2 / GASPI The majority are using MPI only or MPI + OpenMP with some users using MPI + ptreads, CUDA, OpenCL, TBB, UPC, UPC++, GASNe About 1/3 of the responders use one-sided communication (normal windows) with MPI_Win_create, MPI_Win_allocate, and MPI_Win_create Dated Lightin, Parametri rystoti, Instance many second and the standard second Accord to a the texplorability and international databases of the second combination with a wider range of other parallel approaches as seen in our b by using MPI alternatives and MPI + X can be used to improve future version odels as well as the one-sided implementation of MPI are preliminary and imp Tau internal timers + MPI_P* asynchronous execution (74%) overlap of communication with computation (70 uniform interface, not combined with separate