

Point to Point Communication and Collective MPI Performance Test

Yitao Shen
Rensselaer Polytechnic Institute
Troy, NY, 12180
larst@affiliation.org

Jinqiang Jiang
Rensselaer Polytechnic Institute
Troy, NY, 12180
larst@affiliation.org

Chau-Lin Huang
Rensselaer Polytechnic Institute
Troy, NY, 12180
huangc11@rpi.edu

Lorson Blair
Rensselaer Polytechnic Institute
Troy, NY, 12180
larst@affiliation.org

Christopher D. Carothers
Rensselaer Polytechnic Institute
Troy, NY, 12180
chrisc@cs.rpi.edu

ABSTRACT

Communication among multiple nodes plays a critical role in the performance of High Performance Computing. MPI [?] offers a great number of libraries to maximize and test the communication performance in the parallel computing networks. The message passing has been observed to spend additional time in transferring information from one node to another, and several parallel models have been devised to properly describe the phenomenon, which in terms serve as criteria for future benchmarking. In this work, we build a collective MPI performance test to evaluate the performance among different models. To accomplish this, we have a parallel version of Game of Life program optimized with MPI communication scheme and CUDA for GPU parallelization. As far as the authors are concerned, this work could be the first intend to verify the networking performance of a GPU-aided program in terms of different parallel models. In terms of hardware, the benchmarking takes place on AiMOS [?], an eight petaflop supercomputer using a heterogenous system architecture built with IBM POWER9 CPUs and NVIDIA GPUs. On the software side, the program relies on IBM Spectrum MPI and NVidia CUDA math library [?].

KEYWORDS

Parallel Computing, High Performance Computing, CUDA, MPI

ACM Reference Format:

Yitao Shen, Jinqiang Jiang, Chau-Lin Huang, Lorson Blair, and Christopher D. Carothers. 2020. Point to Point Communication and Collective MPI Performance Test. In *Woodstock '18: ACM Symposium on Neural Gaze Detection, June 03–05, 2018, Woodstock, NY*. ACM, New York, NY, USA, 7 pages. <https://doi.org/10.1145/1122445.1122456>

Unpublished working draft. Not for distribution.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted by ACM, provided that the copies are not made for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

Woodstock '18, June 03–05, 2018, Woodstock, NY

© 2020 Association for Computing Machinery.

ACM ISBN 978-1-4503-XXXX-X/18/06...\$15.00

<https://doi.org/10.1145/1122445.1122456>

2020-05-05 19:17. Page 1 of 1–7.

1 INTRODUCTION

2 RELATED WORKS

There are several works discussing the benchmarking of MPI performance on parallel systems using different libraries. Pješivac-Grbović et al. [?] give a thorough overview of the parallel models, and compare the performances on inter-cluster MPI collective operations on two systems. The authors in the mentioned article also demonstrate that the gap between the message sendings depends on the number of unique destination nodes.

There are also works whose authors discuss the communication performance on the parallel version of Conway's Game of Life. L. Ma et al. [?] demonstrate the performance boosting by implementing the algorithm into its parallel counterpart. The article's author implemented the parallel program with OpenMP library, which takes advantage of the multithreading in the CPU rather than seeking time efficiency improvement from multi-node perspective.

In our work we referenced an open-source framework for implementing network benchmarks presented by T. Hoefler et al. [?] This framework separate communication patterns from communication moddules which allows independently added benchmark types and network protocols. The authors also presented a LogGPS pattern [?] which supports measurement of LogP and LogGP models parameters such as latency, overhead and gap per bytes over MPI.

3 GAME OF LIFE AS AN ALGORITHM

The Game of Life as invented by the British mathematician John Horton Conway in 1970 [?] is a cellular automaton. The algorithm is a zero-player game and as the game evolves throughout undetermined number of iterations, the outcome is determined by the given initial configuration. The game is taken place on a two-dimensional orthogonal grid of square cells. The cell status is atomic, that is it can only be found as alive or dead. Each cell's status is determined by eight adjacent neighboring cells. At each step in time, any cell with fewer than two live neighbors dies due to under-population. On the contrary, if the cell lives with two or three live neighbors survives to the next step. If there are more than three immediately adjacent cells, the cell perishes due to over-population. Lastly, the cell resucitates with exactly three live neighbors and can be seen as the result of cellular reproduction.

The operations involved in the Game of Life include instantiation of the world configuration, the update of the cells in the world, and swapping the newly updated world with the previous one. The operations not only is possible to implement serially, but also with parallel speed-up mechanism to take advantage of the efficient memory manipulation offered by NVIDIA CUDA math library, or through the de-facto networking of various nodes through MPI libraries when the dimension of the world increases toward dimension of high orders of magnitude. This latter deployment scenario is the focus of our study.

4 GPU AND CUDA TOOLKIT

According to [?], AiMOS uses NVidia Tesla V100 GPUs in conjunction with the compute nodes. The NVIDIA Tesla V100 accelerator contains the Volta GV100 GPU. Volta is the codename for NVidia's GPU microarchitecture release on December 7, 2017. Volta was NVidia's first chip to feature Tensor Cores, designed specially to yield higher deep learning performance than regular CUDA cores. [?]. The architecture is implemented with TSMC's 12 nm FinFET process.

Tesla V100 delivers 7.8 TFLOPS of double precision floating point (FP64) performance, 15.7 TFLOPS of single precision (FP32) performance, and 125 Tensor TFLOPS based on GPU Boost clock.

In AiMOS cluster, there are 16 nodes each containing four NVidia Tesla V100 GPUs with 16 GiB of memory each. In addition, there are 252 nodes each possessing six of the same accelerators with 32 GiB of memory each [?].

The CUDA Toolkit version used in AiMOS is CUDA 9.1 and 10.0 [?].

5 MPI

AiMOS adopts MPI Spectrum as its message passing interface library. The MPI Spectrum is developed by IBM as a high-performance, production quality implementation of MPI to accelerate end-to-end communication. MPI Spectrum is based on the open-source MPI library, but is integrated with improved RDMA networking add supports NVIDIA GPUs based on IBM PAMI backend. Another feature is that MPI Spectrum enhances collective library running blocking and non-blocking algorithms [?].

6 MODELS TO DESCRIBE THE PARALLEL SYSTEMS

To verify the communication performance of the GPU-accelerated MPI-aided Game of Life program, different configurations of nodes are set. The collective MPI performance is expressed in terms of the four common networking performance models: Hockney [?], LogP [?], LogGP [?], and PLogP [?]. The parallel models can be seen as a sequence of proposals toward establishing the proper description for both point-to-point and collective communication time consumption under any parallel computing system.

The Hockney model is considered the simplest parallel model of communication. The model assumes that the time taken to send a message is

$$T = \alpha + \beta m$$

where α is the size of the message, and β is the inverse of the bandwidth, while m represents the message size. The model is suitable to describe point-to-point communication systems.

The LogP model intends to offer a simple yet more detail view to facilitate the finding of bottlenecks in possible communication latency. The model is described with four parameters: the latency L , overhead o , gap between the sending of messages g , and the number of processors or nodes involved in the communication P . The model assumes that only small amount of messages is transferred simultaneously. The time needed to transfer messages between nodes takes

$$T = L + 2o$$

where L is the latency, and o as the overhead.

Since LogP does not monitor transmission of long messages, LogGP further extend such aspect in its description. A new parameter Gap per byte (G) is taken into account, which is defined as the time per byte for a long message [?]. Unlike the LogP model which restricts to constant small size messages, LogGP allows sending larger messages. Typically, time taken to transfer a message is:

$$T = L + 2o + (m - 1)G$$

where G is the gap per byte and m is the size of the message.

In the work of T. Kielmann et al. [?], parametrized LogP is introduced as a slight extension of LogP and LogGP models that is capable of monitoring the minimal gap between two messages without saturating the network for each message size. In addition to the parameters contained in LogP, additional parameters are included: the sender and receiver overheads, message transfer time and data copying time to and from the network interfaces are included in the latency. Moreover, the gap parameter is defined as the minimum time interval between consecutive message transmission or reception. The overall time spent in the message transferring can be expressed as:

$$T = L + g(m)$$

where $g(m)$ is the gap per message. The worth pointing out that the sender $o_s(m)$ and receiver $o_r(m)$ overheads depend on the message size.

7 PERFORMAMNCE METRICS

To estimate the execution time according to Hockney's model, we used the formula $t(s) = l + s / b$. In the formula, s stands for the message size, l stands for the latency of the network, b stands for the bandwidth of the network.

$$T = \alpha + \beta$$

$$\alpha = l = \text{Latency}$$

$$\beta = \frac{s}{b} = \frac{\text{MessageSize}}{\text{Bandwidth}}$$

Then it came with the question "How can we determine the latency and the bandwidth?" In this benchmark model, it assumes that some process A sends a message to process B, while process B sends message back. The advantage is that it will not require any synchronized clocks between these two processes. While the disadvantage is that it will presume the communication performance or the costs between two points is totally symmetric.

To determine latency, it requires to execute the benchmark for iteration time equal to zero.

For this experiment, the message size changes from 1 to 1024 by the power of two. And the rank number per node changes from 1 to 6 for up to two nodes. After the testing in different conditions, it is found that the latency of the network is almost static. More specifically, to calculate the total execution time, it simply subtracts the end timestamp generated after the for loop of iteration by the start timestamp, both of which are generated via the function `MPI_Time()`.

8 PERFORMANCE RESULTS

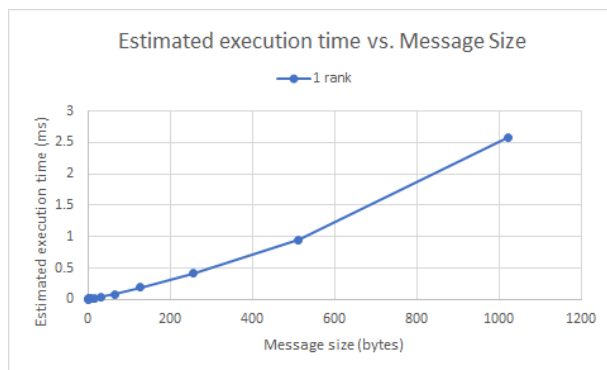
For this experiment, the message size changes from 1 to 1024 by the power of two. And the rank number per node changes from 1 to 6 for up to two nodes. After the testing in different conditions, it is found that the latency of the network is almost static which is fixed around 0.004 ms.

9 CONCLUSION AND FUTURE WORKS

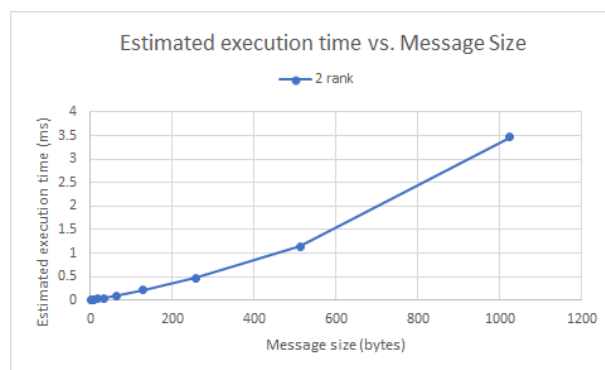
In this work we perform benchmarking on four different parallel models to evaluate the performance of the CUDA-aided Game of Life program. The GPU-accelerated program is integrated with MPI's communication mechanism to effectively manage the memory access as well as transferring among different nodes.

For future works, there are two possible extensions on the program to study how different mechanism may further improve or deteriorate the performance of the current version of Game of Life program. The first possible direction is to incorporate multithreading to the current version to observe how sharing resources on the same node may affect the message communication, and another possible study involves the execution of the same GPU-accelerated program on different topologies such as the Fat Tree [?], Dragonfly [?] or Slimfly [?] networks to investigate the difference in performance described with the mentioned parallel models.

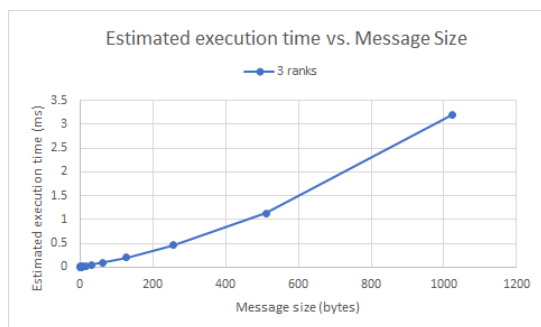
10 ACKNOWLEDGMENTS



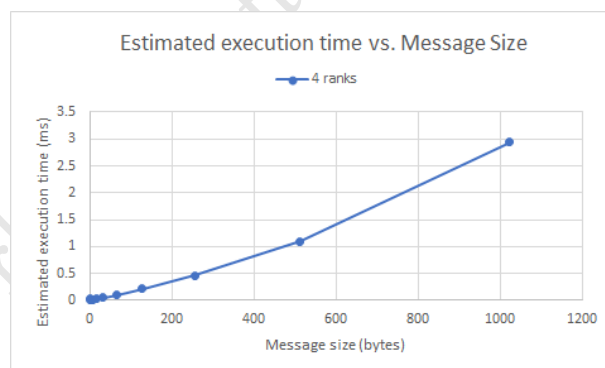
(a) Rank 1



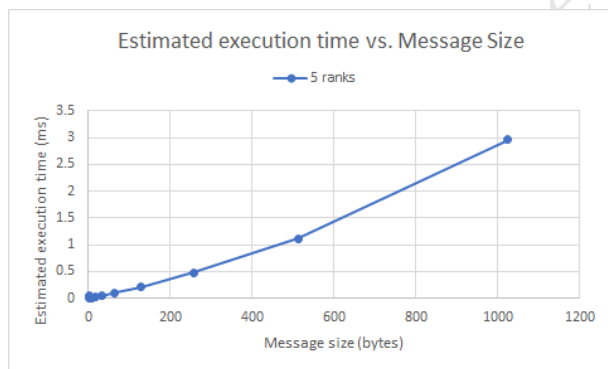
(b) Rank 2



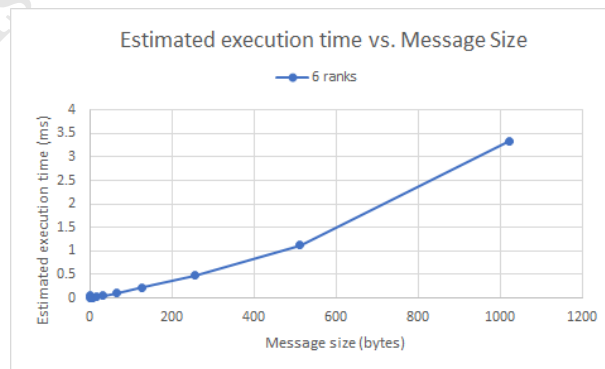
(c) Rank 3



(d) Rank 4



(e) Rank 5



(f) Rank 6

Figure 1: Estimated execution time vs message size along different ranks.

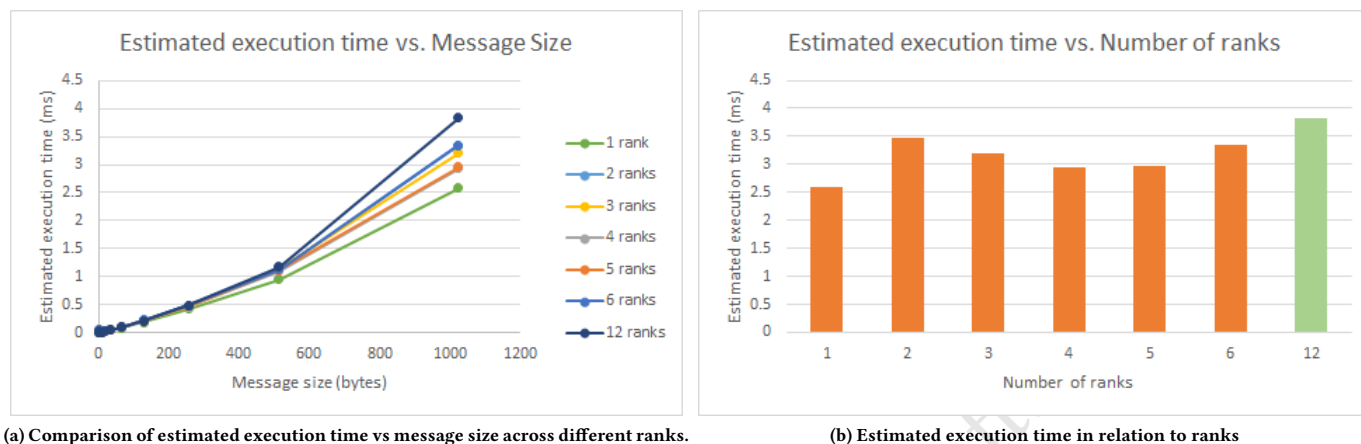


Figure 2: The relationship between message sizes and number of ranks to estimated execution time.

581
582
583
584
585
586
587
588
589
590
591
592
593
594
595
596
597
598
599
600
601
602
603
604
605
606
607
608
609
610
611
612
613
614
615
616
617
618
619
620
621
622
623
624
625
626
627
628
629
630
631
632
633
634
635
636
637
638

639
640
641
642
643
644
645
646
647
648
649
650
651
652
653
654
655
656
657
658
659
660
661
662
663
664
665
666
667
668
669
670
671
672
673
674
675
676
677
678
679
680
681
682
683
684
685
686
687
688
689
690
691
692
693
694
695
696

Unpublished working draft.
Not for distribution.

697	755
698	756
699	757
700	758
701	759
702	760
703	761
704	762
705	763
706	764
707	765
708	766
709	767
710	768
711	769
712	770
713	771
714	772
715	773
716	774
717	775
718	776
719	777
720	778
721	779
722	780
723	781
724	782
725	783
726	784
727	785
728	786
729	787
730	788
731	789
732	790
733	791
734	792
735	793
736	794
737	795
738	796
739	797
740	798
741	799
742	800
743	801
744	802
745	803
746	804
747	805
748	806
749	807
750	808
751	809
752	810
753	811
754	812