

OpenStack Cloud Tuning for High Performance Computing

¹TARINI PRASAD PATNAIK, *Gandhi Institute of Excellent Technocrats, Bhubaneswar, India*

²SIBASIS RATH, *Indus College of Engineering, Bhubaneswar, Odisha, India*

Abstract—High-Performance computing (HPC) is scarcely attempted in clouds because of slow and inefficient Inter-VM communication on the same server as well as huge latency between remote units. This was changed by introduction of *ivshmem*, a PCI device-based shared memory between VMs on the same server, but unfortunately, this mechanism became broken with Linux update few years ago. We have restored this shared memory system and made, for the first time, full cloud integration using latest versions of OpenStack, Linux, QEMU, libvirt and MPICH. Also, the analyses of different factors influencing both TCP/IP and *ivshmem* communication is presented along with tuning techniques that could significantly increase performance. Finally, we have created *ivshmem* communication channel that can replace standard Neutron TCP/IP network, resulting three to six times performance improvement.

Keywords-

OpenStack; cloud computing; high performance computing; cloud tuning

I. INTRODUCTION

Nowadays, Cloud Computing is the dominant general-purpose computing paradigm, and OpenStack [1] is the most popular open-source cloud operating system for private clouds. Unfortunately, High-Performance-Computing (HPC) in a cloud was not possible in the past, because of the huge overhead for inter-VM communication on the same server and between servers as well, as it was shown for example by [2], [3]. However, with the advent of a remake of *ivshmem* [4], a shared memory between VMs on the same server became possible again, thus making HPC in principle feasible for an OpenStack cloud.

In this paper, a set of OpenStack tuning measures are discussed that augment the possibilities users have for HPC in OpenStack, provided that MPICH [5] is engaged. For that purpose, we investigated the case that each MPICH process is allocated to one VM and measured the inter-VM communication on the same server by using the calls `MPI_PUT` for data-exchange and `MPI_WIN_LOCK` for data-synchronization. Both calls were wrapped by us around the original MPICH calls of the same name in order to be able to run the *ivshmem* remake in OpenStack.

We will show in the following various performance tuning methods for HPC in OpenStack via *ivshmem*, as well as for the classical inter-VM communication via TCP/IP that is based on Neutron's [6] Open vSwitch [7] architecture. By our measures, inter-VM bandwidth and latency could be improved by a factor of up to six, from worst case to

best case, as our performance measurements have shown.

The rest of the paper is organized as follows: in chapter 2, the state of the art in inter-VM communication in OpenStack is given. Chapter 3 presents a tuning for the classical TCP/IP communication that is based on level-3 caching and a custom *virtio* [8] network bridge. Chapter 4 explains tuning measures for *ivshmem* based on proper NUMA allocation and vCPU pinning. In this chapter, also the best TCP/IP method is compared to the best *ivshmem* method demonstrating a superior improvement factor for the latter. The paper ends with a conclusion and a reference list.

II. STATE OF THE ART IN INTER-

VM COMMUNICATION IN OPENSTACK

As any cloud, OpenStack is a distributed system, even if the cloud is physically located in the same rack or in the same computing center where TCP/IP would not be needed, because L2 switching would be sufficient. By studying

[7], we were able to draw a block diagram of the resulting software overhead (Fig. 1) for the case that two MPICH processes are executed by two VMs on the same server. According to Fig. 1, any data frame has to go two times through the following stages: TCP/IP stack, device driver, virtual network interface, *qbr* Linux Bridge, Open vSwitch (OVS) Integration Bridge, OVS VLAN Bridge, physical Ethernet Interface, physical Ethernet Switch. Although OVS is part of OpenStack's Neutron network service and certain port configurations could reduce number of intermediate interfaces, it still produces significant overhead. Together with the VM communication overhead, the consequence is an unacceptable low HPC performance.

III. TUNING THE CLASSICAL TCP/IP INTER-

VM COMMUNICATION

For TCP/IP tuning and the subsequent chapters, we used OpenStack Juno, Ubuntu 16.04 as guest OS, CentOS 7.1 as host OS, QEMU 2.9.50, libvirt 2.0, MPICH 3.2 and *virtio* 1.1.1 as software environment. Additionally, we sent data from one MPI process to the other, while varying its size from 2^2 to 2^{20} bytes. The transmission was accomplished by one-sided `MPI_PUT` via the standard MPICH *Nemesis-sock* channel. However, because of the fact that there is normally no shared memory between different VMs, even on the same server, MPICH emulates this functionality by TCP/IP communication, sending data packets back and forth that carry shared variables.

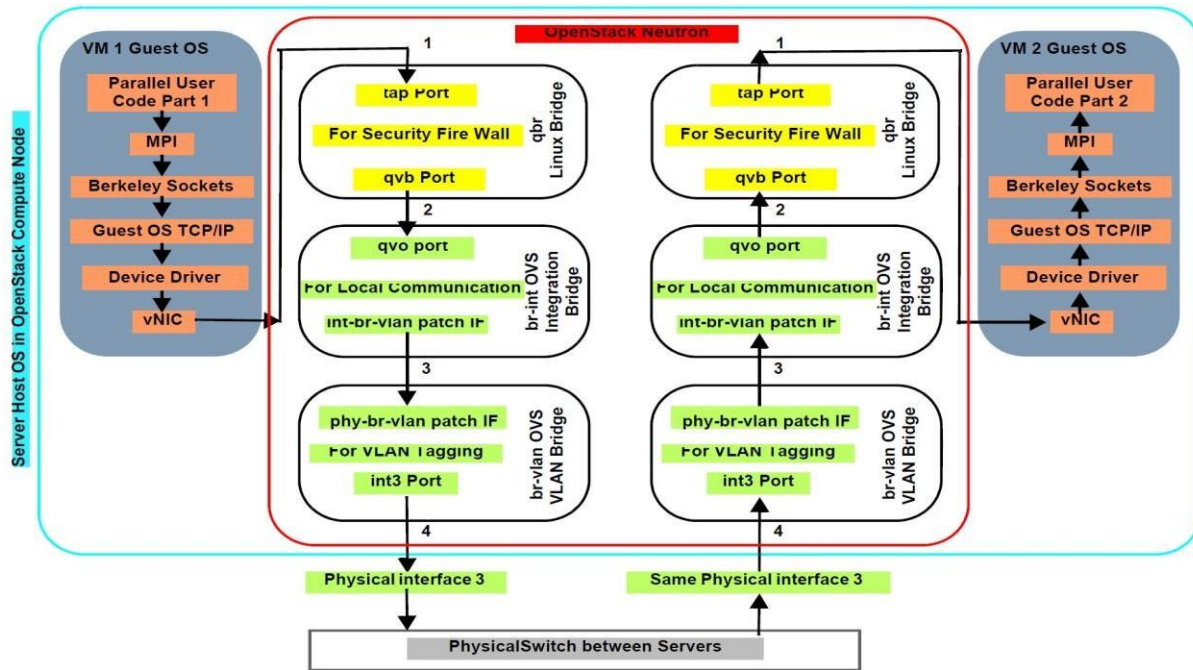


Figure 1. Software overhead in OpenStack for inter-VM communication. 1,2,3 = Ethernet Data Frames, 4 = VLAN-tagged Ethernet Frames, OVS = OpenvSwitch.

A. Level-3 Caching

Now, the influence on OVS of applied level-3 caching is discussed in comparison to level-2 caching only. As a reference, we used the elapsed time for transferring data for the case of no cloud (VMs only) for both, OVS with level-3 caching and without. The results are shown in the orange, grey, and yellow curves of Fig. 2 and Fig. 3.

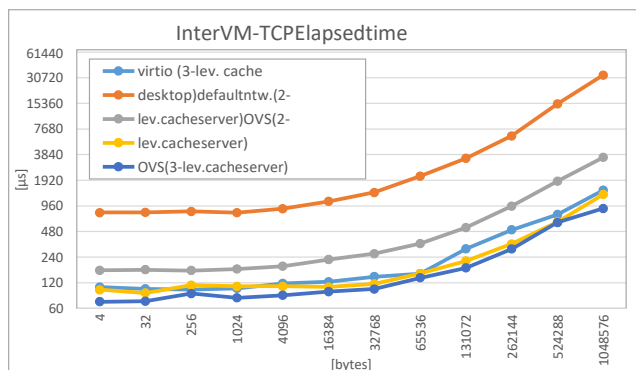


Figure 2. Elapsed times for TCP/IP-based inter-VM communication.

The two orange curves depict the elapsed time and bandwidth for inter-VM data exchange without cloud. In this case, each data packet must go to the first IP router of the server where it is reflected back thus forming a loop. As one can see, if OpenStack is engaged with OVS (grey curves), performance is much better than without cloud. The reason for that is that each packet must only travel to the first switch of the server, as shown in Fig. 1, and not to the first IP router. Either at the switch or at the router, packets are reflected

back. Engaging level-3 caching further improves performance (yellow curves). Similarly, replacing the server by a desktop PC (light blue curves) showed also good performance, but disadvantage was limited scalability and cloud incompatibility. Finally, the only measure that made significant difference was replacement of Neutron's OVS by virtio-net.

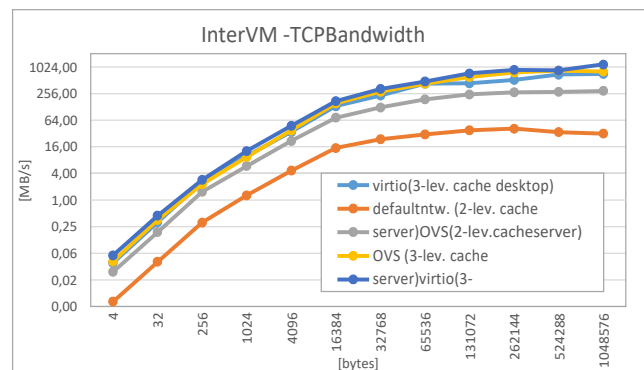


Figure 3. Bandwidth for TCP/IP-based inter-VM communication.

B. Virtio-Net

The dark blue curves in Fig. 2 and Fig. 3 show the effect of virtio instead of OVS. This replacement is possible because the KVM hypervisor of OpenStack has two interfaces: the first is used by the various QEMUs [9] to cooperate with. The second API is virtio that provides for "paravirtualization". This is a more efficient IO virtualization method than the so-called full software emulation made by OpenStack by means of KVM/QEMU. Besides being an API,

virtio is also a library of paravirtualized device drivers for the guest OS. Virtio-net is the paravirtualized device-driver for a virtual Ethernet card (vNIC). Internet-based inter-VM communication can profit from virtio-net if a Berkeley send call results in calling virtio-net. This driver is aware that it is executed in a VM, and it is therefore actively cooperating with its QEMU. With virtio-net, KVM does not need to intercept guest-OS device-driver accesses to emulated

PC devices, because they are not performed. Instead, the parameters for virtio calls are directly forwarded to QEMU. Technically, the Berkeley send call is not mapped onto a vNIC in guest OS, but via QEMU onto a virtio-net send queue ("virtqueue") in host OS. Virtqueues are much simpler and thus faster than vNICs, because they are only buffers. The rest in virtio-net happens as described in section A. Please note also that the paravirtualized guest-OS device-driver is called front-end driver, while the modified host-OS Ethernet-driver is termed back-end driver. The original host-OS back-end driver cannot be used in virtio-net because its input is a Linux data-structure called sk_buf, while the front-end driver outputs so-called virtqueue entries. Both, the front-end and the modified back-end driver, are contained in the virtio library. The resulting block diagram is shown in Fig. 4, where only the first half of the communication is displayed because it is symmetric.

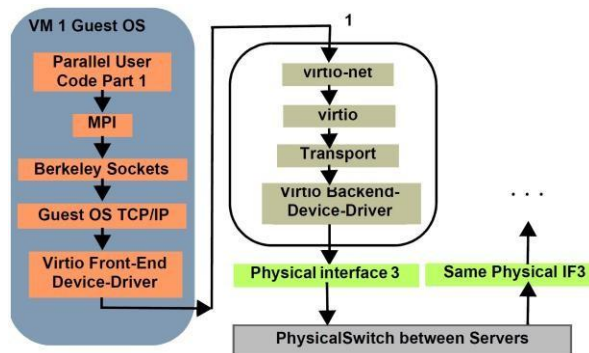


Figure 4. Virtio-net inter-VM communication architecture.

The first disadvantage of virtio is that TCP/IP is still engaged. The second is that the virtqueue entries have to be handled by two QEMUs each, one for the source, the other for the target VM. Finally, for every data frame sent by the real Ethernet card, each QEMU has to make a system call to KVM, which is a time-consuming procedure, because it requires a full process context switch, with all MMU page table entries reloaded. Because of that, we strive for a better inter-VM communication method.

IV. IVSHMEM

Ivshmem is a virtual PCI device in a guest OS which is emulated by KVM/QEMU. It establishes a Linux/POSIX shared memory (SHM) between the VM and its host OS. Ivshmem thus enables zero-copy VM-to-Host communication and vice versa, which is very efficient with respect to bandwidth and latency, because no internal data buffer exists. Ivshmem can also be used for inter-VM

communication with the host OS SHM as an intermediate step. Ivshmem is implemented by mapping its virtual PCI device memory to the host OS SHM. This is possible, because the memory is emulated by QEMU as a data structure inside of itself, and because multiple QEMUs can communicate with each other in host OS.

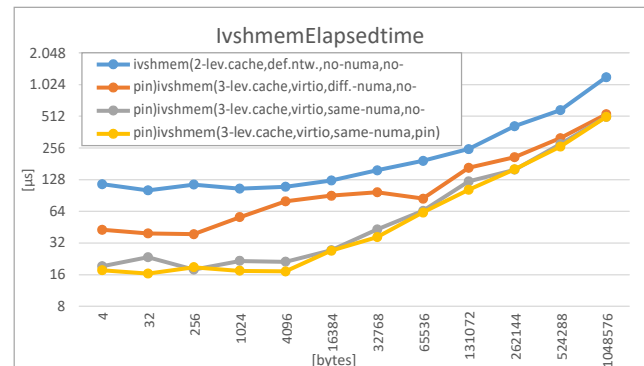


Figure 5. Elapsed times for ivshmem-based inter-VM communication.

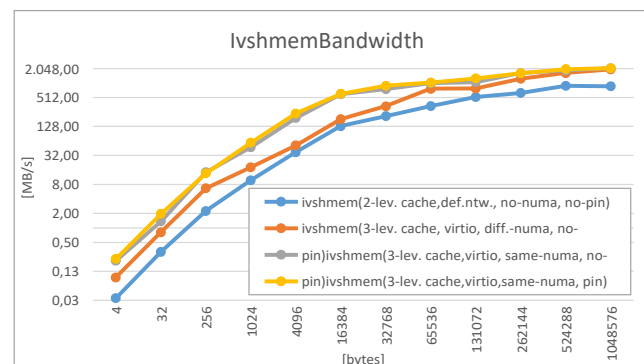


Figure 6. Bandwidth for ivshmem-based inter-VM communication.

Ivshmem was supported for a few years by libvirt and QEMU as a virtual PCI device that allowed for shared memory (SHM) between guests and host. A synchronization mechanism called Shared-Memory-Server (SMS) was created, which became part of QEMU as its ivshmem-server. This server could send interrupts to VMs located on the same compute node, thus allowing unblocking of an application in a target VM, which was waiting for data from an application in a source VM. Mutual blocking and unblocking was used for SHM access synchronization without consuming CPU cycles (no spinlock), which is indispensable for HPC. However, this valuable feature could only be implemented if both applications were implemented as the user part of Linux user-IO device-drivers (uio). Additionally, it was the task of the user to write, by means of code examples of both, the kernel and user part of the uio driver and to integrate his application into this driver. This was not applicable for HPC with MPI, because MPI processes are typically not device drivers. The reason for ivshmem being limited to uio device drivers was that the one and only possible spinlock-free synchronization between VMs in Linux is possible by means of a blocking read in the user part of uio. As soon as

an interrupt arrives at a VM, this user part is called as an interrupt service routine and the blocking read at the target VM is unblocked. Unfortunately, with newer Linux versions, unblocking did not work anymore, with the consequence that `ivshmem` was nearly useless because of missing SHM-access synchronization. Our contribution was to patch codes for both, the kernel and the user part of `uio` device drivers for `ivshmem` and to find proper version and configuration matches for Linux, QEMU and libvirt so that everything works again. We have also created a new SHM communication and synchronization channel inside of MPICH. In particular, we achieved `ivshmem` to run for the first time in a cloud with OpenStack, resulting in a three to six times performance improvement compared to TCP/IP (Fig. 7 and Fig. 8).

The isolated `ivshmem` performance results are depicted in Fig. 5 and Fig. 6. The blue curves show that implementing default emulated TCP/IP network, as required for initial MPI synchronization, has still very poor performance, but far better than non-SHM solution, shown in Fig. 2 and Fig. 3. Performance improves further if the same tuning measures are engaged as before, i.e. level-3 caching and `virtio` instead of `OVS`. The effect is shown in the orange curves. Additional gains are possible by proper NUMA allocation and CPU pinning.

A. Proper NUMA Allocation, no-CPU Pinning

In the gray curves of Fig. 5 and Fig. 6, it becomes visible that manually allocating the communicating VMs into the same non-uniform memory-access (NUMA) region yields in further significant improvements. In such a uniform memory access region, all VMs have the same access latency to the physical shared memory in host OS. On the other hand, allocating the MPI VMs to different NUMA regions leads to frequent cache misses, and data access times are not equal anymore as well.

B. vCPUPinning

If cache misses occur together with a rescheduling of the vCPU, that executes a VM, from one physical processing unit (i.e. core) to another, the result is a non-monotonously increasing performance for increased message size.

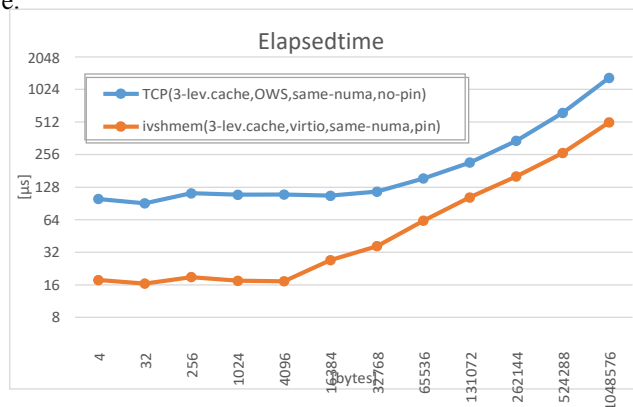


Figure 7. Elapsed times for TCP and `ivshmem`-based inter-VM communication

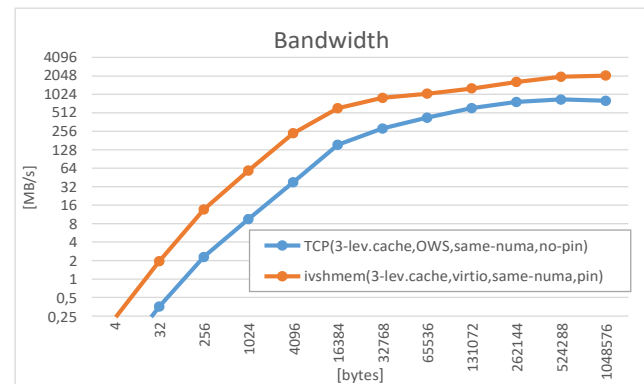


Figure 8. Bandwidth for TCP and `ivshmem`-based inter-VM communication.

This is noticeable in the gray curve of Fig. 5 for sizes from 4B to 1KB. Therefore, each VCPU from a VM was manually bound by us to a specific physical processing unit which is known as CPU pinning.

The utmost performance boost is achievable if `ivshmem` is used together with level-3 caching, `virtio`, proper NUMA allocation and CPU pinning. This is demonstrated by the yellow curves in Fig. 5 and Fig. 6. In this case, bandwidth reaches 2GB/s for a message size of 1MB, which is twice as much as for the best TCP/IP case. For smaller message sizes, the difference is even bigger. The results were achieved without using SR-IOV [10] as hardware accelerator for communication.

Finally, direct comparison of standard TCP/IP cloud (`OVS`) performance with our `ivshmem` integration is depicted in Fig. 7 and Fig. 8, for Elapsed time and Bandwidth respectively. We have deliberately used same-NUMA region for measurements, since NUMA allocation is done randomly by OpenStack scheduler and with each instance creation could engage different memory region. For the smaller size messages, when synchronization is dominating communication, performance difference is factor of six in favor of `ivshmem`. With further increase of message block size, data flow becomes main contributor to overall communication time and the difference drops to factor of three, for the biggest messages. This is remarkable result, considering that NUMA tuning was not applied, due to random nature of OpenStack scheduler.

V. CONCLUSION

High-Performance-Computing in a cloud was not possible in the past, because of the huge overhead generated during inter-VM communication on the same server and consequently between remote units. This was changed with the advent of a remake of `ivshmem`, which is a physical shared memory between VMs on the same host server. As one of the main results, we made it possible for the first time to use the `ivshmem` remake in an OpenStack-based private cloud. Additionally, we introduced several tuning methods, such as proper NUMA allocation and CPU pinning and therefore, adapted `ivshmem` to perform even better. Also, we gave effective measures for TCP/IP-based emulation of

physical shared memory, which are level-3 caching and virtio instead of Neutron's Open vSwitch. All methods were evaluated and compared with each other by measurements, showing that the best ivshmem scenario is at least a few times as fast as the best TCP/IP case. All measurements were made with our wrapper versions of MPICH's MPI_PUT for data-exchange and MPI_WIN_LOCK for shared memory synchronization.

REFERENCES

- [1] Open source software for creating private and public clouds, <https://www.openstack.org/>
- [2] R. Ledyayev, H. Richter, High Performance Computing in a Cloud Using OpenStack, The Fifth International Conference on Cloud Computing, GRIDs, and Virtualization, CLOUD COMPUTING 2014, <http://www.iaria.org/conferences2014/CLOUDCOMPUTING14.html>, Venice, Italy, May 25-29, 2014.
- [3] H. Richter, About the Suitability of Clouds in High-Performance Computing, Proc. Sixth International Conference on Computer Science and Information Technology (CCSIT 2016), Journal Computer Science and Information Technology (CC&IT), Volume 6, Number 1, January 2016, pp. 23-33, Volume Editors: Jan Zizka, Dhinakaran Nagamalai, ISBN: 978-1-921987-45-8, DOI: 10.5121/csit.2016.60103, <http://airccj.org/CSCP/vol6/csit64803.pdf>, AIRCC Publishing Cooperation, Zurich, Switzerland, January 02-03, 2016.
- [4] P. Ivanovic, H. Richter, Performance Analysis of ivshmem for High-Performance Computing in Virtual Machines, Proc. 2nd International Conference on Virtualization Application and Technology (ICVAT 2017), Shenzhen, China, Nov. 17-19, 2017.
- [5] A. Amer, P. Balaji, W. Bland, W. Gropp, R. Latham, H. Lu, MPICH User's Guide, Version 3.2, Mathematics and Computer Science Division-Argonne National Laboratory, Nov. 11, 2015
- [6] Introduction to OpenStack Networking (neutron), <https://docs.openstack.org/liberty/networking-guide/intro-os-networking.html>
- [7] Open vSwitch in OpenStack, <https://docs.openstack.org/liberty/networking-guide/scenario-classic-ovs.html>
- [8] Virtual I/O Device (VIRTIO) Version 1.0. Committee Specification Draft 01/Public Review Draft 01, <http://docs.oasis-open.org/virtio/virtio/v1.0/csprd01/virtio-v1.0-csprd01.pdf>, Dec. 03, 2013.
- [9] S. Weil, QEMU version 2.10.93 User Documentation, <https://qemu.weilnetz.de/doc/qemu-doc.html>
- [10] P. Kutch, B. Johnson, SR-IOV for NVF Solutions - Practical Considerations and Thoughts, rev. 1.0, Networking Division, <http://www.intel.com/content/dam/www/public/us/en/documents/technology-briefs/sr-iov-nfv-tech-brief.pdf>, Feb. 23, 2017