

# The Transition of High Performance Computing Applications from Supercomputers and Clusters to the Cloud

Maria Efthymiadou  
MSc. Computer Science  
Vrije Universiteit  
Amsterdam Netherlands  
Email: maria.g.efthymiadou@gmail.com

Alex Oberhauser  
MSc. Computer Science  
Vrije Universiteit  
Amsterdam Netherlands  
Email: a.oberhauser@student.vu.nl

**Abstract**—In this paper the transition of HPC from supercomputers and clusters to the cloud is outlined. After a short introduction the term HPC is defined. On the base of HPC application requirements the paper describes why the first generation cloud is not suitable for HPC and how the major players in the cloud space try to fix that by going back to the cluster architecture. The paper concludes with an outlook to future development.

## I. INTRODUCTION

This paper is the result of a literature study that answers the question if public cloud providers are suitable for High Performance Computing (HPC) applications. In order to answer this non-trivial question we define the term HPC in section II. In this section we also define some properties that are important for HPC applications. Based on this properties we explain why or why not public clouds are suitable to host HPC applications. In section III we explain the shortcomings and drawbacks of clouds without dedicated HPC support. For the scope of this paper we call this type of cloud *first generation cloud*. The *second generation cloud*, as explained in section IV, solves this problems by going back to the cluster architecture. Finally the paper concludes with section V with a short outlook to future developments.

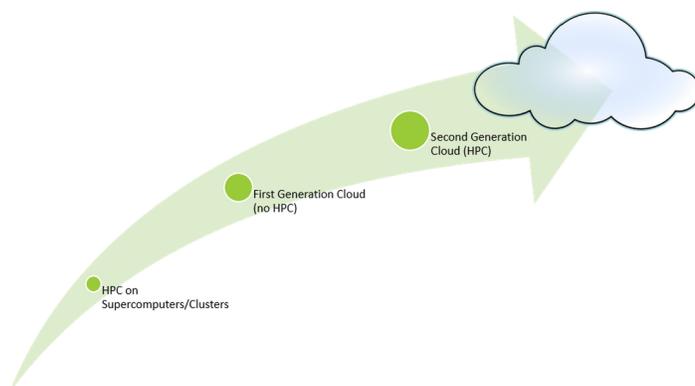


Fig. 1. Transition of HPC to the Cloud

In Figure 1 the transition of HPC on supercomputers to a more cloud friendly approach is outlined. The document struc-

ture follows exactly this transitions and hence the evolution of HPC.

## II. HIGH PERFORMANCE COMPUTING AND SUPERCOMPUTER

The word Supercomputer is closely related to a specific time period. For example, every person nowadays has a Supercomputer compared to the possibilities of the first ones invented back in 1960's. Which needs drove scientists and manufactures to construct such powerful machines 50 years ago? The answer is simply the increase of data and the need of computational power. These are the same reasons for which nowadays supercomputers and the cloud are used, to enable handling High Performance Computing applications. At this point, it is necessary to define in a structure way the term High Performance Computing:

"High-performance computing (HPC) is the use of parallel processing for running advanced application programs efficiently, reliably and quickly. HPC uses supercomputers and computer clusters to solve advanced computation problems." [6]

Derived from the above definition HPC applications have the following properties:

**Fast Computation Power** In order to process and analyse the mass of data HPC applications need high computation power. For specific use cases, specific hardware is used, e.g. GPUs, FPGAs or ASICs.

**High Throughput** The computation power is spread over multiple physical (or virtual) machines. This makes it necessary to transfer the data over the network. Additionally, it is important that the local I/O operations are fast enough to handle the workload.

**Load Balancing** In order to achieve scalability a load balancer is needed. This load balancer distributes the workload over multiple instances.

High Performance Computing either based on supercomputers or computer cluster, depending on the needed performance it can be costly and there is a big overhead to maintain. This has led the cloud providers to explore and

deliver support for HPC applications which offer the benefits of cloud computing for these demanding applications.

### III. FIRST GENERATION CLOUD: CLOUD WITHOUT DEDICATED HPC SUPPORT

In the scope of the paper we use the term *first generation cloud* service to indicate that the cloud services which were developed before HPC as a Service were added to the portfolio of the cloud providers. This first generation of public cloud services that were developed by Amazon, Microsoft, Google and other providers were limited for High Performance Computing since the target group of this type of public cloud services were mainly business customers and not the scientific community. As mentioned in [4] the main problem with HPC in cloud is partially the computation power, but mainly the network latency. High network latency results in increased execution time for applications that use the Message Passing Interface, MPI. For other applications, especially when the network latency is not so important the classical public cloud service could be an alternative [5]. Dependent on the underlying application and different non-functional attributes, like the trade-off between costs and performance, the drawbacks could be negligible. The total execution time in the cloud is slower than private clusters and supercomputers, but through the on-demand model it could be a cheaper and more customizable alternative. The term customizable, means that different software versions could be installed by the person that executes or develops the HPC application.

One of the most important and attractive features of cloud is the on demand scaling of the resources with related on demand billing. Cloud provider have experts that setup and optimize the environment which will host the HPC application. Moreover, the client receives infrastructure when needed and in different price and performance categories. [1] [6]

An important, but most times overlooked, decision criteria for the choice between a private cluster, private cloud, hybrid cloud, grid or public cloud is privacy. Under some conditions the data to be processed could not be stored in current public clouds, through privacy concerns. One such example of privacy aware data are confidential medical data.

The architectural decisions made during the development of cloud data centres, that allow big scalability on-demand, are the biggest obstacles for HPC. One such problem is the geographically distribution of leased virtual machines. For cloud users it is not possible to control the geographical location of the underlying physical host. This results in most cases to increased network latencies. Another obstacle was the virtualization overhead. One problem to mention was the slow access to PCI devices via the virtualization stack. This issue was solved via PCI passthrough by all major virtualization solutions, although it is not clear if these changes were incorporated into the current public cloud provider solutions. Additionally, the trade-off between security and performance in form of an intermediate firewall, that is controlled by the provider (e.g. the Amazon Firewall), increases the network latency, but also security. Also the use of off-the-shelf hardware could result in the deterioration of performance. Finally the fact that many applications require a high throughput for MPI messages can decrease the performance significantly [2].

a) *Advantages of first generation cloud over supercomputers/cluster:*

- ✓ On demand payment without overhead of managing the infrastructure.
- ✓ Customizable down to the software versions and operating system.
- ✓ Easy extension with new instances.

b) *Disadvantages of first generation cloud over supercomputers/cluster:*

- ✗ Not suitable for high throughput applications (network) through high and/or fluctuating network latency.
- ✗ Only suitable for smaller HPC applications.
- ✗ Virtualization overhead
- ✗ No support for special use case, like GPU clusters.

These problems and challenges have resulted in a general purpose cloud that is suitable for businesses, but not necessary for High Performance Computing. One solution is to re-design current HPC applications to fit the cloud model, if possible. Another solution, that comes directly from the cloud providers, are cloud service on top of clusters and supercomputers.

### IV. SECOND GENERATION CLOUD: CLOUD WITH DEDICATED HPC SUPPORT

The second generation of HPC cloud efforts are moving towards the solutions to the problems mentioned in the previous section. The two cloud leaders (Amazon and Microsoft) acknowledge the necessity of returning to the original infrastructure used for HPC, which are the clusters and the supercomputers, and they provide more suitable hardware for that. The cluster instances they offer have higher performance, compared to the ones from the first generation, because they have greater CPU, memory and for specific applications also GPU instances. These instances improve specific HPC application types, through increased computation speed. It is important to mention that this solution is not suitable for all kinds of computations. Besides that there is no possibility for every instance to have a GPU, which is not necessarily negative considering that there is a point where further addition of GPUs decreases the overall performance. [3]

The biggest cloud providers are going back to the clusters because they realised that the HPC service they were providing did not perform well. More specifically, concerning network connection speed Amazon offers 10 Gigabit/s and Microsoft Intraband. Microsoft offers A8 (Intel Xeon E5-2670 with 8 virtual cores @ 2.6 GHz and Memory 56 GB) and A9 (Intel Xeon E5-2670 with 16 virtual cores @ 2.6 GHz and Memory 112 GB) instances for MPI applications. Moreover the company claims, according to there official website, that they provide high throughput network based on remote direct memory access (RDMA) technology. On the other hand Amazon offers EC2 C3 Instance cluster - Amazon EC2 Cluster with number of cores equal to 26, 496 and 105, 984 GB memory. Cloud providers are starting to support in the second generation the HPC use case.

The increased network speed and the selection of the same geographical location of instances in a cluster improves the network latency problem of the first generation cloud. In combination with virtualizations that support PCI passthrough and the support for specific use cases, such as GPU clusters, this improves the performance further. Although the on-demand model and the sharing of a supercomputer with other customers holds the performance not yet as high as classical supercomputers and clusters do, the pay-as-you-go model and the increased customizable makes the HPC as a Service attractive for researchers, especially for smaller to mid-size computations.

## V. CONCLUSION

The first generation cloud or cloud without dedicated HPC support was not targeted towards a scientific use case. The main target market were and are enterprises. The second generation cloud extends the first generation with dedicated HPC support. Currently we see two different strategies how the public cloud market leader support HPC. Amazon provides direct support for the technology with there *Amazon HPC* service. This service is mainly IaaS and does not come with additional services that are tailored towards the HPC scenario. Microsoft takes another approach and supports real HPC as a Service with software tool support (PaaS). Additionally, they support a big part of the community behind HPC with there *Microsoft Azure for Research* campaign. Both solve the shortcomings of the first generation public cloud by using the architecture of clusters and supercomputers and providing services on top. This combines the power of supercomputers with the convenience and advantages of the cloud. In the Table I different non-functional metrics are compared.

	Cluster	First Gen	Second Gen
Network Latency	Low	High	Low
Computation	Fastest	Slow	Fast
Virtualization	None	Overhead	Optimized
Costs	Highest	Low	High
Customizable	No	High	High

TABLE I. COMPARISON OF SUPERCOMPUTER/CLUSTER WITH FIRST AND SECOND GENERATION CLOUD FOR HPC APPLICATIONS.

During the early stage of the cloud hype we have seen one solution fits all approaches. In the future we will see more specialized solutions, but with the same services and interfaces on top of the underlying technology. This includes not only HPC as a Service, but also mobile clouds and any other cloud types or hybrid solutions that will appear.

Despite the fact that the cloud providers have adapted to the changing market for HPC applications not all problems are solved. For example the performance fluctuation, due to resource sharing with other customers, could be a problem. There will be always room for improvement also from the customer side. The architecture of the application itself needs to be designed and adapted in such a way that it fully exploits the cloud architecture. For example, parallel computation when suitable, better load balancing and finally in case the bandwidth is low the client needs to have the data close, so the search of a more suitable solution is important.

Finally, we conclude that the current public cloud providers are suitable for the most HPC applications and provide special-

ized services for this purposes. Exceptions are computations that need a whole super computer and not only parts of it, for example the mapping of the human brain or global climate research.

## REFERENCES

- [1] P Rizwan Ahmed. The state of high performance computing in the cloud.
- [2] Jaliya Ekanayake and Geoffrey Fox. High performance parallel computing with clouds and cloud technologies. In *Cloud Computing*, pages 20–38. Springer, 2010.
- [3] Roberto R Expósito, Guillermo L Taboada, Sabela Ramos, Juan Touriño, and Ramón Doallo. General-purpose computation on gpus for high performance cloud computing. *Concurrency and Computation: Practice and Experience*, 25(12):1628–1642, 2013.
- [4] Keith R Jackson, Lavanya Ramakrishnan, Krishna Muriki, Shane Canon, Shreyas Cholia, John Shalf, Harvey J Wasserman, and Nicholas J Wright. Performance analysis of high performance computing applications on the amazon web services cloud. In *Cloud Computing Technology and Science (CloudCom), 2010 IEEE Second International Conference on*, pages 159–168. IEEE, 2010.
- [5] Christian Vecchiola, Suraj Pandey, and Rajkumar Buyya. High-performance cloud computing: A view of scientific applications. In *Pervasive Systems, Algorithms, and Networks (ISPAN), 2009 10th International Symposium on*, pages 4–16. IEEE, 2009.
- [6] Ye Xiaotao, Lv Aili, and Zhao Lin. Research of high performance computing with cloud. 2010.