

# Snowmass Letter of Interest - Cloud Computing - CompF4

Burt Holzman <[burt@fnal.gov](mailto:burt@fnal.gov)> (Fermilab), Andrew Norman (Fermilab)

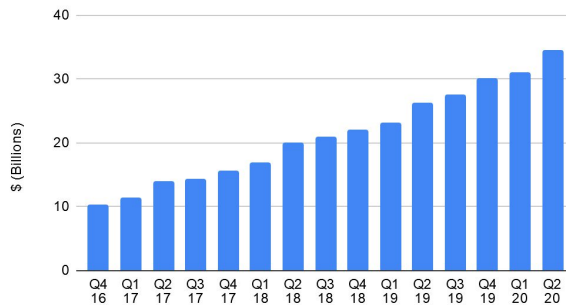
## Overview

The world currently spends more than \$30B per quarter on the consumption of Cloud Computing services [1]. This is 17 times the size of the entire FY20 budget for the Office of Science at the Department of Energy. These resources have been successfully used for scientific computing in HEP and elsewhere [2-4] under a pay-as-you-go model where users are billed monthly based on the resources they have consumed.

There are a wide range of Cloud services, but we categorize them into “capability” and “capacity”. *Capability* services represent a unique set of features that we have not

provisioned on-premises for a variety of reasons (cost-effectiveness, power consumption, proprietary solutions, etc.) *Capacity* services are services that allow us to scale out commodity services; historically we have focused on high-throughput (batch) computing.

Global Cloud Spending

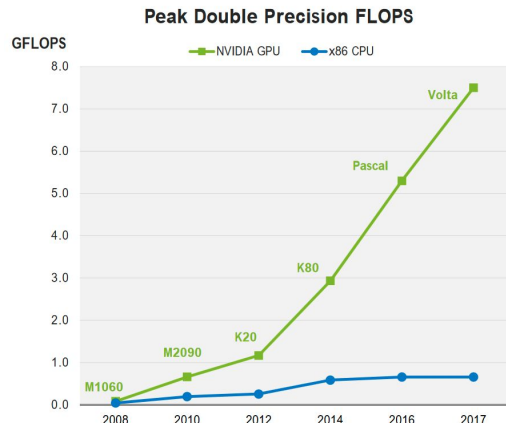


## R&D - Capability

In addition to seemingly boundless capacity, cloud computing also offers capabilities that may not be feasible, available, or cost-effective to purchase and deploy on-premises. The big three companies in this space (Amazon AWS, Microsoft, Google) provide heterogeneous computing architectures - specifically, architectures that are not CPUs (GPUs, FPGAs, ASICs). The cloud can provide a research and development environment in order to evaluate these technologies before they make the transition from capability to capacity. Some technology is vendor-specific and only available off-premises, but may be a very good fit for classes of problems. For example, both Google and Amazon offer optimized ASIC solutions for machine learning problems (Google TPU, AWS Inferentia) while Microsoft offers a custom FPGA configuration (Project Brainwave).

Heterogeneous hardware evolves rapidly (as seen in the plot of GPU trends from [5]). This implies that hardware purchased for our data centers approaches obsolescence more quickly than traditional computing. Additionally, due to market forces, these technologies are first available only in the cloud, with deployment at on-premises and shared facilities coming significantly later. These resources are not utilized at 100%; even the world’s largest

supercomputer, Oak Ridge's Summit, is at most 87% utilized [6]. Depending on usage, it may



be more cost-effective to subscribe to the cloud's pay-as-you-go model rather than amortize the entire cost up-front. Additionally, there are services of interest to the community that are only available on cloud; for example, poster sessions for the Neutrino 2020 conference and the Fermilab Annual Users' Meeting were conducted in a push-to-deploy cloud-hosted virtual reality environment [7].

For this capability-based R&D, we facilitate research on commercial cloud by providing direct access to the providers' own interfaces rather than proxying through a bespoke service. Privileges are limited in order to

manage risk and costs; for example, only administrators deploy new GPU instances, but the end-user can log in and power it on and off at will; the instances incur compute costs by the minute only when powered on. Additional monitoring is provided - state changes are recorded and broadcast to all users. An administrator-configurable timer can be used to power off instances to help minimize costs.

## HEPCloud - Capacity

The current HEPCloud program allows for dedicated on-premises batch computing resources to be significantly augmented by transient cloud resources during periods of peak demand. This has the potential to realize a significant long term cost savings compared to sizing on-premises computing to the brief peaks in the collider and neutrino program's analysis needs.

Since its inception, both the global computing landscape and needs of the HEP community have shifted dramatically. The scientific community needs to access diverse computing hardware and heterogeneous computing platforms, both as testbeds for new analysis and computing techniques, as prototypes and onramps to the DoE sponsored leadership computing facilities, and as capacity-based resources in order to meet the computing demands of the future. HEPCloud will grow to provide on-demand access to heterogeneous computing architectures, without the researchers needing to become domain experts in the mechanics and infrastructure that are required to interface HEP analysis and workflows with large-scale commercial cloud platforms. It will build upon the foundation from the R&D effort to transition these capability-based resources to capacity. This allows for an elastic expansion of unique, costly, and short lived specialized hardware to support AI and ML application training and other infrequent resource-intensive tasks.

We propose the evolution of HEPCloud to provide a portal and interface that allow researchers to concentrate on their applications and not infrastructure, while ensuring that the complexity,

security and cost optimization associated with these advanced resources can be provided and managed in a centralized fashion.

## References

- [1] Canalys.com. 2020. [online] Available at: <[https://www.canalys.com/static/press\\_release/2020/Canalys-cloudq220.pdf](https://www.canalys.com/static/press_release/2020/Canalys-cloudq220.pdf)> [Accessed 27 August 2020]
- [2] Foster, I. and Gannon, D., 2017. *Cloud Computing For Science And Engineering*.
- [3] Holzman, B., Bauerdick, L.A.T., Bockelman, B. et al. HEPCloud, a New Paradigm for HEP Facilities: CMS Amazon Web Services Investigation. *Comput Softw Big Sci* **1**, 1 (2017). <https://doi.org/10.1007/s41781-017-0001-9>
- [4] Sfiligoi, I., Schultz, D., Riedel, B., Wuerthwein, F., Barnet, S. and Brik, V., 2020. Demonstrating a Pre-Exascale, Cost-Effective Multi-Cloud Environment for Scientific Computing. *Practice and Experience in Advanced Research Computing*,.
- [5] Hester, K., 2017. *HPC | Volta*. [online] SEG Annual Meeting. Available at: <[http://hpcsociety.memberlodge.com/resources/Documents/SEG2017/Ken%20Hester%20SEG\\_HPC\\_NVIDIA.pdf](http://hpcsociety.memberlodge.com/resources/Documents/SEG2017/Ken%20Hester%20SEG_HPC_NVIDIA.pdf)> [Accessed 27 August 2020].
- [6] Oak Ridge Leadership Computing Facility, High Performance Computing Facility Operational Assessment 2019 (2020). ORNL/SPR-2020/1499
- [7] Hubs.mozilla.com. 2020. *Hubs By Mozilla*. [online] Available at: <<https://hubs.mozilla.com/>> [Accessed 27 August 2020].