

# Analysis of Scientific Cloud Computing requirements

**A. López García, E. Fernández-del-Castillo**  
*aloga@ifca.unican.es*  
*Spanish National Research Council (CSIC)*



**CSIC**

CONSEJO SUPERIOR DE INVESTIGACIONES CIENTÍFICAS



**IFCA**

Instituto de Física de Cantabria

Cloud and Scientific Computing

Cloud and Scientific Computing

Application use cases

Requirements for a Scientific Cloud

Conclusions

- IFCA — CSIC Cloud infrastructure is being user by a broad range of scientific users since 2011.
- Starting from a set of virtualized resources and a separated small cloud testbed, we have moved to a complete cloud integration
- We have worked closely with our users to adapt their workloads to a Cloud environment whenever we have seen that they will profit from it.
- From our experience, we have obtained a set of common requirements that are needed form a Cloud testbed to be suitable to accommodate scientific computing.

Cloud and Scientific Computing

Cloud and Scientific Computing

Application use cases

Requirements for a Scientific Cloud

Conclusions

## Cloud Computing

“a model for enabling ubiquitous, convenient, on-demand network access to a shared pool of configurable computing resources (e.g., networks, servers, storage, applications and services) that can be rapidly provisioned and released with minimal management effort or service provider interaction.”<sup>1</sup>

## Scientific Computing

“intersection of numeral mathematics, computer science and modelling”<sup>2</sup> to solve scientific problems.

---

<sup>1</sup> Mell and Grance. *The NIST Definition of Cloud Computing*. 2009.

<sup>2</sup> Karniadakis and Kirby. *Parallel Scientific Computing in C++ and MPI*. 2003.

Cloud and Scientific Computing

Cloud and Scientific Computing

Application use cases

Requirements for a Scientific Cloud

Conclusions

- Many of the Cloud benefits are present in other computational models:
  - Pay per use.
  - Access to distributed and heterogeneous computing resources
  - ...
- Cloud computing fills some gaps that are impossible or difficult to satisfy in current scientific datacenters.
  - Customized environments
  - On-demand access.
  - Elasticity.
  - Self service.
  - Non-conventional application models.

### Customized Environments.

- Freedom to choose or create its own execution environment.
- One of the biggest advantages against the Grid.

### On-demand access, elasticity, self service.

- Infinite capacity offered by commercial providers is not possible.
- On-demand access is a key for interactive workloads.
- Self-service interface can be either an advantage or a disadvantage (learning curve).

## Non-conventional application models.

- The cloud is not focused on the execution of atomic tasks (parallel or not).
- Complex and long-running execution environments.
- Simulation of dynamic software agents<sup>3</sup>, decision making process in urban management<sup>4</sup>, behavioural simulations using shared-nothing Map-reduce techniques<sup>5</sup>, etc.

---

<sup>3</sup> Sethia and Karlapalem. "A multi-agent simulation framework on small Hadoop cluster". 2011; Talia. "Cloud Computing and Software Agents: Towards Cloud Intelligent Services".

<sup>4</sup> Khan et al. "An architecture for integrated intelligence in urban management using cloud computing". 2012.

<sup>5</sup> Wang et al. "Behavioral simulations in mapreduce". 2010.

Cloud and Scientific Computing

Cloud and Scientific Computing

Application use cases

Requirements for a Scientific Cloud

Conclusions

## PROOF

- Used by the HEP community in the last phases of the analysis to produce final plots and numbers.
- Usually I/O bounded (analyze a several GB to a few TB) with relatively low CPU usage (mainly filtering).
- Requires a master node that serves as an entry point and distributes the workload to a variable set of workers.

### Benefits from the Cloud:

- An IaaS can support these interactive, on-demand and short lived sessions that are disposed after usage

## Particle Physics Phenomenology

- Multiple independent software packages used by the community.
- They are combined into complex workflows, where each step requires input from previous codes execution.
- Each scenario may require different versions of the software.
- Some of these packages need the usage of a proprietary software (Mathematica).
- Setting up this environment is a overwhelming task for the scientific.

## Benefits from the Cloud:

- Scientific users can deploy a stable infrastructure with their exact requirements.
- Snapshotting a given image for later usage or recreating previous experiments.

## GIS Pattern recognition

- Usage of multiple software packages: Modis (satellite data analysis tools), GRASS and GDAL (geospatial libraries and tools), PROJ4 (cartographic data management) and R (statistical programming environment) to analyse the vegetation indexes (NDVI and EVI).
- Deploying the software on the Grid → long and exhausting process for the user.
  - Incompatibilities between system libraries and installed software.
  - Totally new computing environment for them!
- A database was needed to store the results.

## Benefits from the Cloud:

- Deploy the image and start computing in just a few hours.
- Deploy their own database system to store their data.
- They could start its environment (database and workers) and shut it off when they do not need it.

Cloud and Scientific Computing

Cloud and Scientific Computing

Application use cases

Requirements for a Scientific Cloud

Conclusions

We have identified several requirements that need to be tackled both by the Cloud middleware and by the resource providers.

- Application level requirements.
- Performance related issues (apart from CPU).
- Enhancement of scheduling policies.
- Absence of vendor lock-in.

### Use of scientific image catalogs:

- Building an image might be a time consuming task for non seasoned users.
- Most scientific users are not prone to create, manage and maintain their own images.
- A predefined set of images, with common scientific software packages should be available.
- Also, predefined images with licensed and institutional software.

### Contextualization:

- Process of installing, configuring and preparing software dynamically at boot time.
- This way predefined images can be smaller, paying the tax of a longer boot time.
- Useful for software with well defined dependencies that evolves fast

### High performance communications:

- Parallel applications need of low-latency high speed interconnects (Infiniband or 10GbE).
- Common in HPC environments, not so common in the Cloud.
- Cloud and virtualization software need to manage these devices (passthrough, SR-IOV).

### High performance data access:

- Data oriented workloads demand high speed access towards the data to be analysed.
- In a Cloud framework, data is normally decoupled from the instance (i.e. block device access a la EBS or object storage a la S3).

### Instance co-allocation.

- Some workloads require the parallel execution of tasks across several nodes.
- Large requests of independent nodes should be discriminated from tightly-coupled or parallel nodes.
- The former can be served on a first-come first-served basis; the later ones need some advanced scheduling and placement policies.

### Short start-up overhead.

- When a request is made some time is needed to prepare the instance (image transfer, image resize, data injection, etc.).
- Scheduling mechanism should take into account this overhead.
- Large requests can introduce a bigger penalty, for example because several images need to be transferred in parallel.
- Some smart image preparation mechanism is needed (pre-caching, shared catalogs, etc.)

Performance aware placement.

- Two virtual machines competing for the utilization of the same resource (for example two I/O bounded machines) should not share the same physical host.
- Access to specialized hardware (low-latency networks, GPGPUS, etc.) should be also taken into account.

Spot and preemptable instances.

- Long running tasks are common in scientific computing
- The instances hosting such executions can be transparently paused to let some higher priority tasks execute.
- Lower cost for the user.
- Better resource utilization.

Bare metal provisioning.

- The cloud middleware should foresee the provisioning of bare metal machines for the cases where a virtualized node is not an option.

### Interoperability, lock-ins

- Interoperability should be a must, since it will enable collaboration.
- The usage of open standards and recommendations (OCCL, CIMI, CDMI) at the API level is a key factor.
- The virtual machines also suffer from the "hypervisor lock-in"  
→ Open Virtualization Format (OVF).

Cloud and Scientific Computing

Cloud and Scientific Computing

Application use cases

Requirements for a Scientific Cloud

Conclusions

- Cloud computing is not the technology to rule them all.
- Cloud computing may be really useful for the long-tail scientific users.
- Cloud middleware needs to satisfy the scientific computing needs, not only industry needs.
- Scientific Computing datacenters should move to a mixed model: cloud computing in addition to more traditional computing power (batch system, Grid).

Thanks!