### **Clouds for research computing**

Randall Sobie Institute of Particle Physics University of Victoria

Collaboration UVIC, NRC (Ottawa), NRC-HIA (Victoria)

1



## Outline

- Research computing
- Cloud computing
- Clouds for research computing
- Distributed research cloud for particle physics and astronomy
- Outlook

Why is the universe not made of equal amounts matter and antimatter?

We build instruments large detectors to record the collisions of matter and antimatter

collisions

"Events"











#### **Desktop computing**

12+ processor cores 10+ TB storage 10G network cards



#### **Applications**

Interactive computing Visualization Remote sensor control

#### **Cluster computing**

300+ cores/rack 1000+ TB storage Interconnectivity



#### Application

High Throughput HTC Small scale parallel jobs (physical sciences)



#### **Super Computers**

TOP500 100K cores

High-speed interconnectivity



#### **Applications**

Climate modeling Cosmology

## **Research Computing**

#### **Computing Grids**

Distributed compute clusters 100K+ cores 1000TB ++ storage



#### **Applications**

Very High Throughput Particle physics Earth sciences Health sciences





#### laaS

Delivers computing infrastructure as a service

Science clouds Commercial clouds (EC2)

HTC platform

### **Cloud computing**



#### PaaS

**Delivers computing platform or software stack as a service** Commercial clouds (EC2)

eg Instrumentation control





#### SaaS

**Delivers applications or software stack as a service** Commercial clouds (EC2)

eg Mathematica or Matlab



## Why use clouds?

The BaBar project stopped recording data in 2008.

Complex code developed over 15 years Limited to specific operating systems and libraries Diminishing resources and few people.

Virtualization is the only solution for preserving the software environment



### Data preservation

Topical issue for many research fields



Particle physics community has largely ignored this issue often because new facilities make the old data obsolete.

Projects are now decades long and unlikely to be repeated for many years.

We need to ensure the data is accessible for the long term.

Challenge We need to preserve the software environment



### Requirements

Sophisticated user communities Non-GUI users Batch computing Complex software packages and demanding system requirements Specific OS system Specific application libraries and compilers





Medium-scale data sets (100s TBs) Data accessed (on-demand) from remote repositories

## Conceptual design of a distributed cloud



#### **Design goals:**

Leverage existing work in grid computing (authentication, data management, networks)

Use existing research computing facilities and get access to new resources

Boot user-customized VMs in a familiar batch computing environment

Simplify systems configuration by removing the application dependence

Use the network to move data to the clouds

#### Sky Computing or Grid of Clouds

### **Components**



Application encapsulation Image replication eg Xen, KVM

> Dynamic resources eg Condor, SGE

### VM management : Repoman



## **CERN-VM Filesystem (CVMFS)**



## **Authentication**

We use X509 certificates for authentication (except Amazon EC2)

Used in particle physics (LHC/CERN) and also by Westgrid (Compute Canada)

Certificates issued by Grid Canada

**X509** is an **ITU Telecommunication** standard for a public key infrastructure (PKI) for single sign-on

We use it to X509 certificates for user job management, access to Repoman and access to the data storage



### The Interactive System

User saves the modified environment as a new image



### laaS cloud resources





CS looks at the job queue and sends a request to the next available cloud to boot the User-VM



User view of the system is the same as a standard batch environment

Job script contains a link to the user's VM required for the job

### **Cloud Scheduler**



Cloud scheduler looks at the job queue

Makes a request to boot a user-VM on a cloud

The cloud retrieves the user-VM from the repository

The user-VM attaches itself to the Condor pool and Condor sends the user-job to the user-VM.

The user-VM stays active if there are more jobs that require it.

### **Simulation Production**



Randall Sobie IPP/University of Victoria

## Astronomy applications

CANFAR Project Canadian Advanced Network for Astronomical Research UIVC, UBC, NRC-HIA CANARIE-funded project



Distributed cloud used to process survey data

In production for 8 months using different laaS cloud resources

Compute Canada cloud site at UVIC

Enabling system for user analysis as well as production jobs

See Tuesday presentation

### **Current status**

- Run up to 500 simultaneous jobs over 8 clouds
  - > 100,000 successful jobs
  - Should scale for low-IO applications
- Testing user-analysis (chaotic, high IO)
  - Early test :10TB analysis in 2 days
- Caching of user-VMs (only transfer 1 16G VM)
  - Testing squid-cache
- Network challenges
  - Needed to work with CANARIE to resolve some issues
  - CANARIE will soon connect with Amazon EC2

## Summary

- Our distributed cloud focuses on applications in physical sciences with large high-throughput (HTC) workloads and a knowledgeable user community
- Fault-tolerant system using multiple-laaS (commercial or science) cloud resources
- Based on open-source components with two new in-house elements
  - Cloud scheduler and Repoman (VM repository)
- Production use by astronomers and BaBar
- >100,000 VMs booted

Support provided by CANARIE, NSERC, NRC, Amazon, Google, FutureGrid (NSF)

Randall Sobie IPP/University of Victoria

# Components & References

- Open Source code developed by University of Victoria:
  - Cloud Scheduler >=0.11.1, https://github.com/hep-gc/cloud-scheduler
  - Repoman, https://github.com/hep-gc/repoman
- Other Open Source components used:
  - Scientific Linux 5.x (Xen, KVM), http://www.scientificlinux.org
  - Nimbus >=2.5, http://www.nimbusproject.org
  - Condor >=7.4, http://www.cs.wisc.edu/condor
  - MyProxy, http://grid.ncsa.illinois.edu/myproxy
  - Xrootd, http://xrootd.slac.stanford.edu
  - Lustre >=1.8.3, http://wiki.lustre.org/index.php/Main\_Page
  - Squid 2.7.STABLE8, http://www.squid-cache.org
  - Munin 1.4.5 (epel repository), http://munin-monitoring.org