

# cloudScheduler V2: Distributed Cloud Computing

Randall Sobie

F. Berghaus, K. Casteels, C. Driemel, M. Ebert, F. F. Galindo C. Leavett-Brown,  
D. MacDonell, M. Paterson, R. Seuster, S. Tolkamp, J. Weldon

*University of Victoria*

*TRIUMF*

# cloudscheduler

Designed in 2009 to provision virtual machines (VMs) on clouds  
(dedicated/pledged and opportunistic, private and commercial clouds)

ATLAS and Belle II experiments

CANFar - Canadian Advanced Network For Astronomy Research

Typical HEP workloads 5000 cores (peaks up to 10,000 cores)

Clouds in North America and Europe

(sometimes in Australia)

Openstack, Open Nebula, Amazon EC2, Google GCE, Microsoft Azure

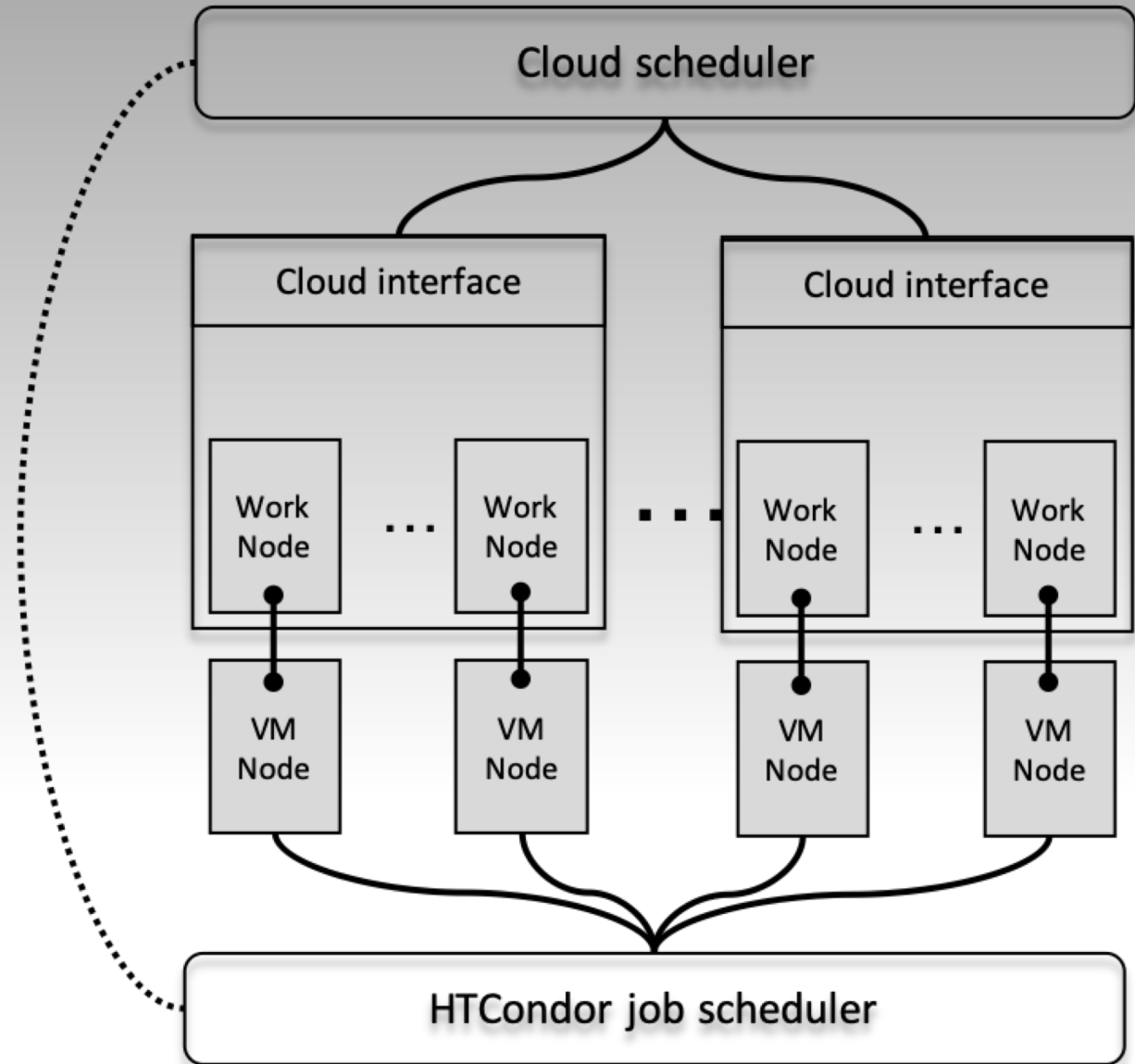
Dynafed data federator

Belle II Canadian Raw Data Centre (Tier-1)

## cloudscheduler Version 1

(HTCondor and CS run at Victoria)

1. User submits job to HTCondor
  - Idle job and no resources
2. CS queries HTCondor for job requirements
3. CS looks for cloud that meets job requirements
  - If a match, requests creation of VM instance
4. Boot and contextualize a CERN micro-VM (CernVM)
  - Use CVMFS for OS/software
  - VM joins HTCondor pool
5. User job has resources and job executes on VM
6. Job completes
  - Other jobs in queue can run on VM
  - Otherwise, CS retires VM instance
7. CS destroys VM when there are no more jobs



## cloudscheduler Version 1

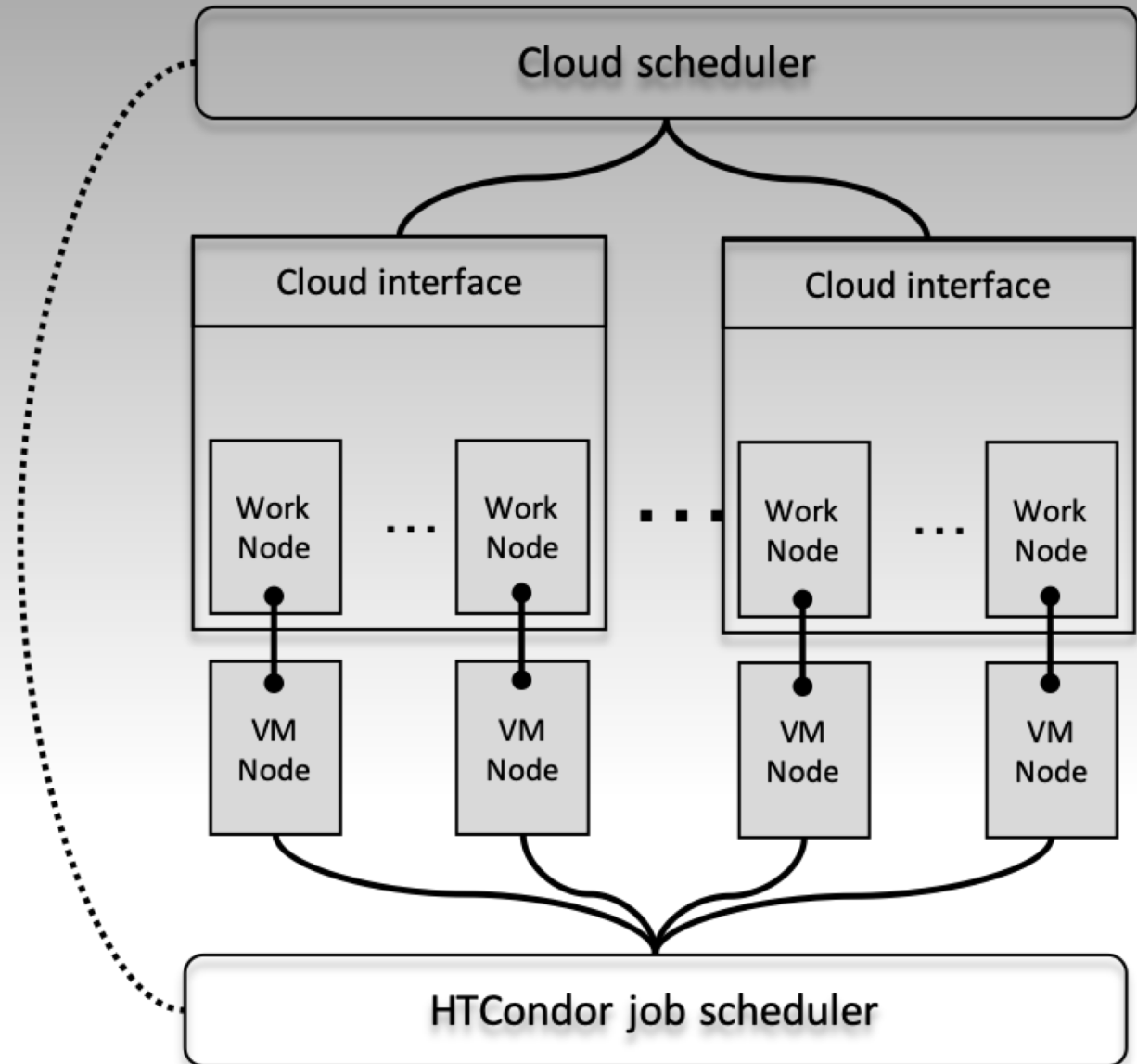
(HTCondor and CS run at Victoria)

1. User submits job to HTCondor
  - Idle job and no resources
2. CS queries HTCondor for job requirements
3. CS looks for cloud that meets job requirements
  - If a match, requests creation of VM instance
4. Boot and contextualize a CERN micro-VM (CernVM)
  - Use CVMFS for OS/software
  - VM joins HTCondor pool
5. User job has resources and job executes on VM
6. Job completes
  - Other jobs in queue can run on VM
  - Otherwise, CS retires VM instance
7. CS destroys VM when there are no more jobs

## Time for a rewrite after 10 years

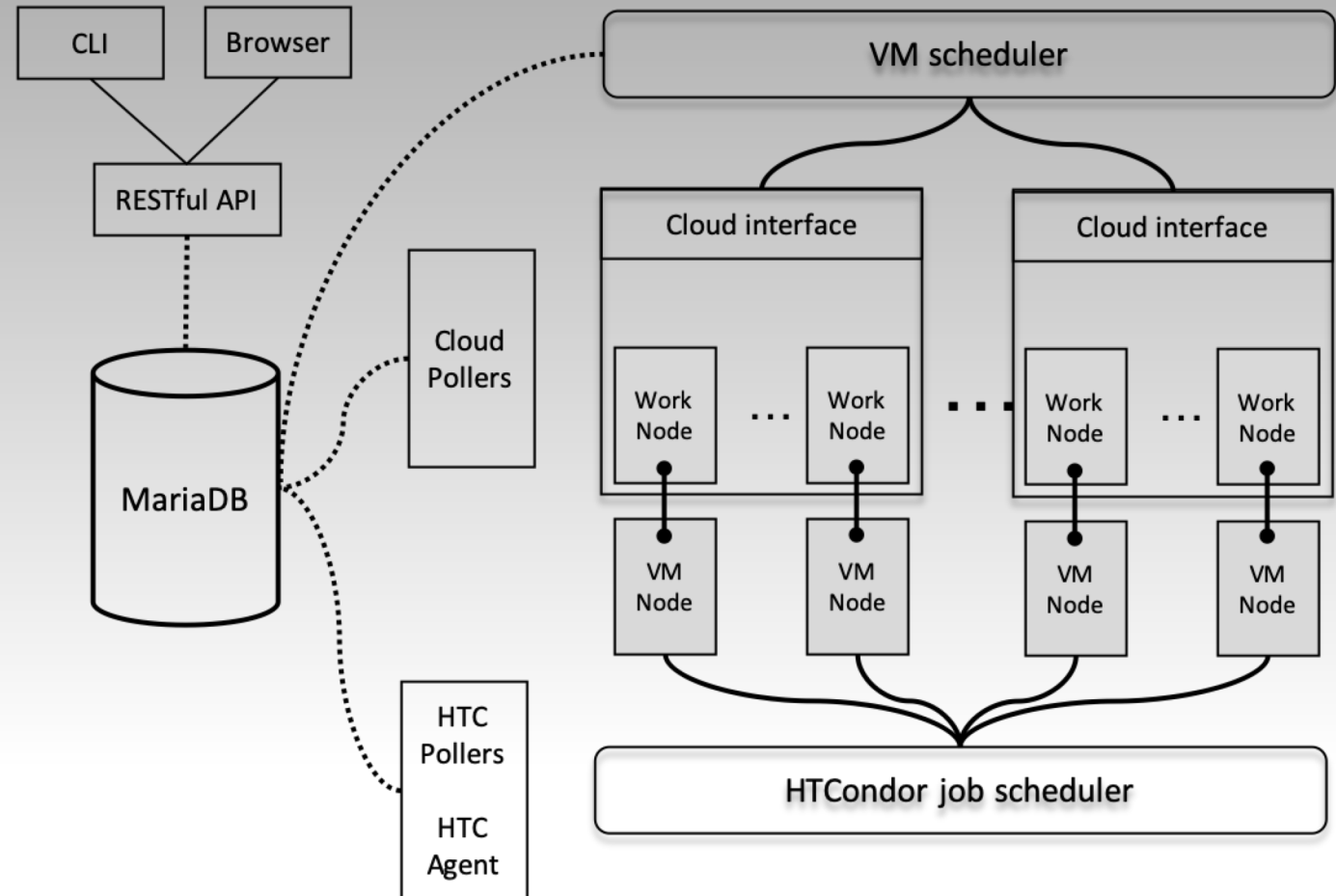
Python 3, redesigned architecture, GUI, ..

New functionality, opportunistic use of clouds, ..



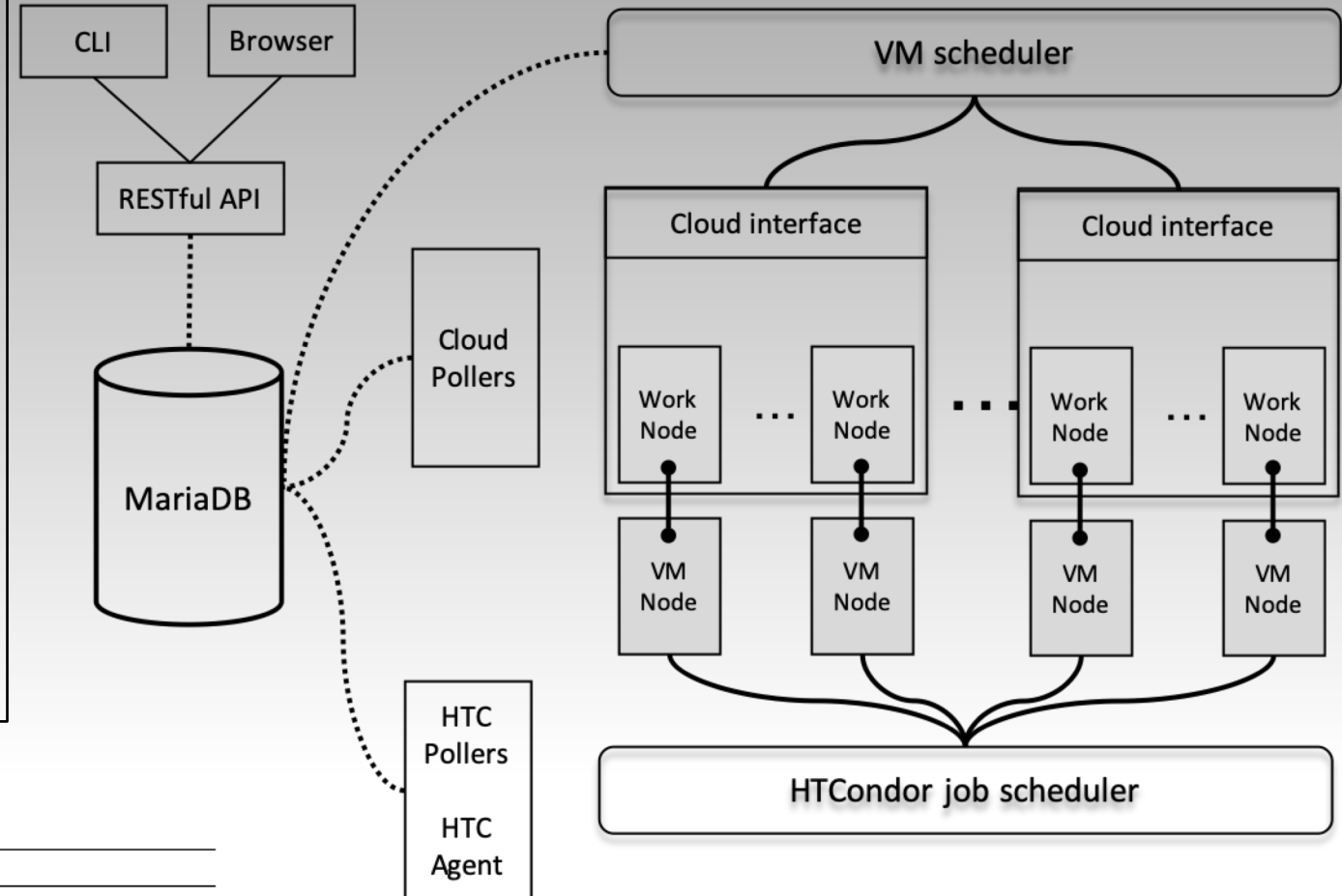
## Cloudscheduler Version 2

- In-memory data of CSV1 are replaced with a MariaDB database
  - Keep track of system state in MariaDB
  - More responsive with MariaDB
  - Resilient to outages and easier to maintain
- Independent scheduling, polling, and user interface processes that track the state of the clouds and the HTCondor pool
  - Improved VM and job scheduling
- RESTful web user interface for
  - Improved administration and management
  - Monitoring of the clouds and jobs
- Expanded functionality
  - Multiple projects (experiments)
  - Multiple HTCondor instances
  - Opportunistic sharing in a cloud



## Cloudscheduler Version 2 Workflow (Similar to CSV1)

- Determined by configuration of the system, clouds, job information and the state of VMs
- Mandatory contextualization
  - metadata passed to “cloud-init”
- Additional contextualization
  - Project configuration
  - Benchmark VM
  - Activate ES reporting by the VM
  - Apel accounting



VM state	Description
<i>starting</i>	VM is booting/contextualizing
<i>unregistered</i>	VM is running and has not registered in HTCondor pool
<i>idle</i>	VM is running, registered in HTCondor pool and not running jobs
<i>running</i>	VM is running, registered in HTCondor pool and running jobs
<i>retiring</i>	VM is running, retired in HTCondor pool and will complete running jobs
<i>manual</i>	VM is flagged as being manually used and will be ignored by the <i>VM Scheduler</i>
<i>error</i>	VM in error state according to the cloud information

## Cloudscheduler Version 2 GUI/CLI

(CLI and GUI have nearly identical functionality)

- CSV2 has users and groups (group = expt)
  - Access rights configurable
- Clouds configuration
  - Quotas, instance types
  - Opportunistic sharing
- VM images and SSH keys
- Distribute VM images to Openstack clouds
  - (formerly “glint” service)

The screenshot displays the Cloudscheduler Version 2 GUI. The top navigation bar includes tabs for 'atlas', 'Status', 'Clouds' (highlighted in green), 'Aliases', 'Defaults', 'Images', 'Keys', 'Users', and 'Groups'. On the left, a sidebar lists cloud configurations: 'arbutus' (selected), 'Settings', 'Metadata', 'Exclusions', 'arbutus-k8s', 'arbutus-nf', 'cc-east', 'chameleon', 'lrz', 'otter', and a '+' icon. The main content area shows the configuration for the 'arbutus' cloud. It includes a list of settings on the left (Enabled, Priority, Cloud type, URL, Region, Project, Username, Password, CA certificate, User domain name, Project domain name) and a list of settings on the right (Security Group, VM Keyname, VM Network, VM Image, VM Flavor, VM Keep Alive, Spot Price, Cores Softmax, Cores, RAM). The 'Enabled' checkbox is checked. The 'Cloud type' is set to 'openstack'. The 'URL' is 'https://arbutus.cloud.co'. The 'Region' is 'RegionOne'. The 'Project' is 'ATLAS-Services'. The 'Username' is 'atlascs'. The 'Password' is 'Update password'. The 'CA certificate' is '/etc/ssl/certs/CABundle.'. The 'User domain name' is 'CCDB'. The 'Project domain name' is 'Default'. The 'Security Group' is 'condorWorker'. The 'VM Keyname' is empty. The 'VM Network' is 'VLAN3327'. The 'VM Image' is empty. The 'VM Flavor' is 'c8-30gb-186'. The 'VM Keep Alive' is '-1'. The 'Spot Price' is '-1.0'. The 'Cores Softmax' is '1900'. The 'Cores' is a slider from -1 to 5000. The 'RAM' is a slider from -1 to 13107200 KB. An 'Update Cloud' button is at the bottom.

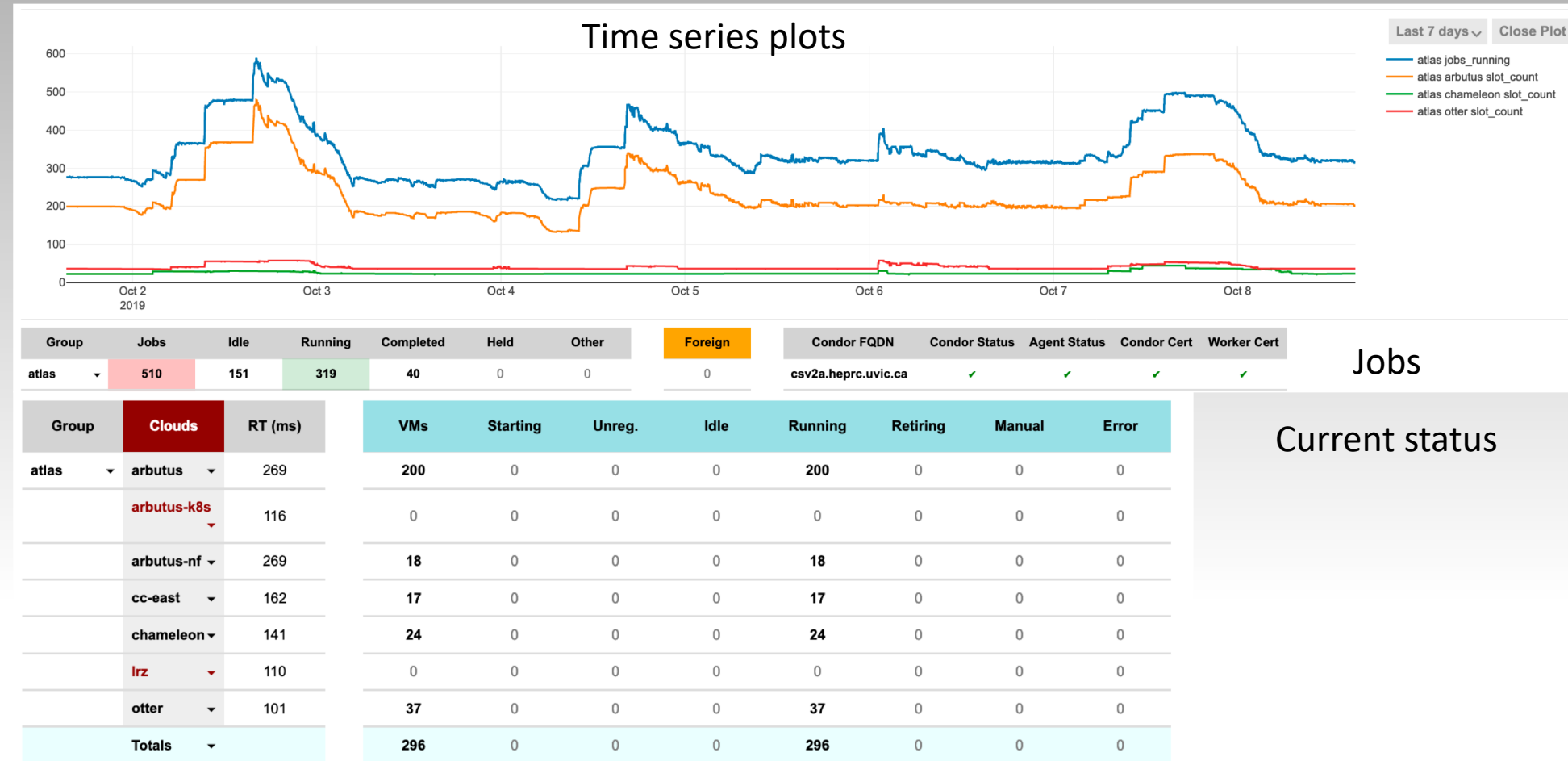
Setting	Value
Enabled	<input checked="" type="checkbox"/>
Priority	0
Cloud type	openstack
URL	https://arbutus.cloud.co
Region	RegionOne
Project	ATLAS-Services
Username	atlascs
Password	Update password
CA certificate	/etc/ssl/certs/CABundle.
User domain name	CCDB
Project domain name	Default
Security Group	condorWorker
VM Keyname	
VM Network	VLAN3327
VM Image	
VM Flavor	c8-30gb-186
VM Keep Alive	-1
Spot Price	-1.0
Cores Softmax	1900
Cores	-1 / 5000
RAM	-1 / 13107200 KB

Update Cloud

## Monitoring

Jobs  
VMs

Also  
HTCondor slots  
Core usage  
Memory  
Central services



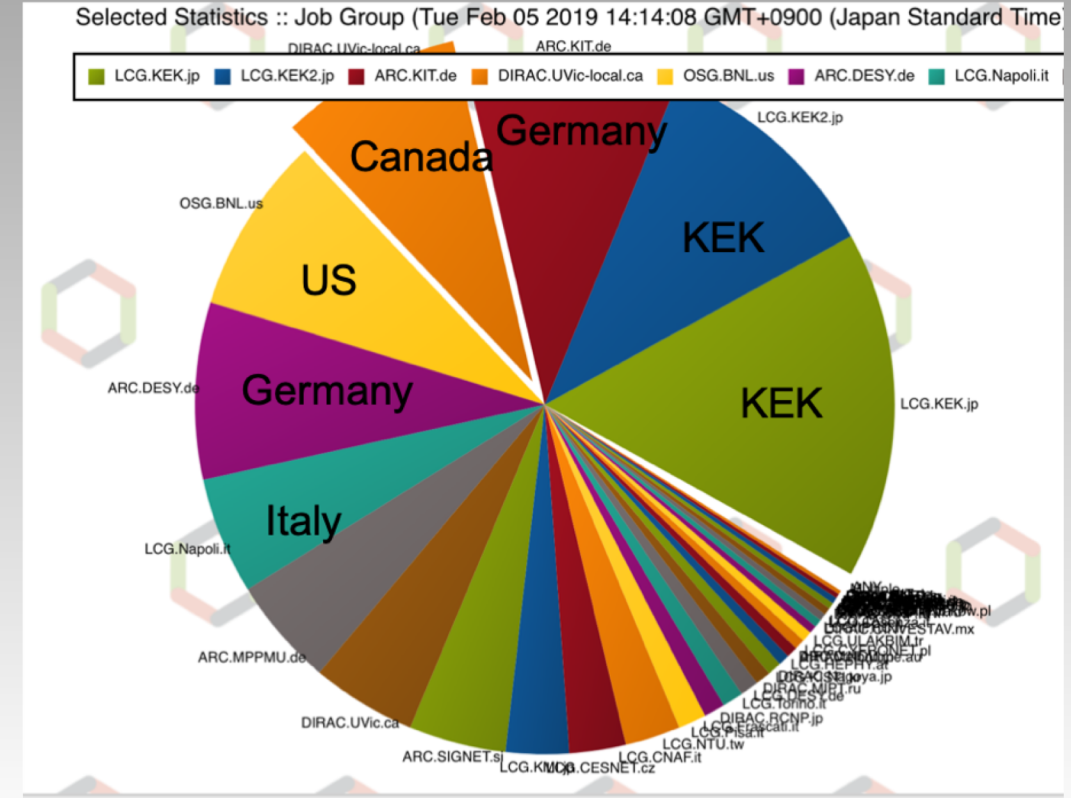
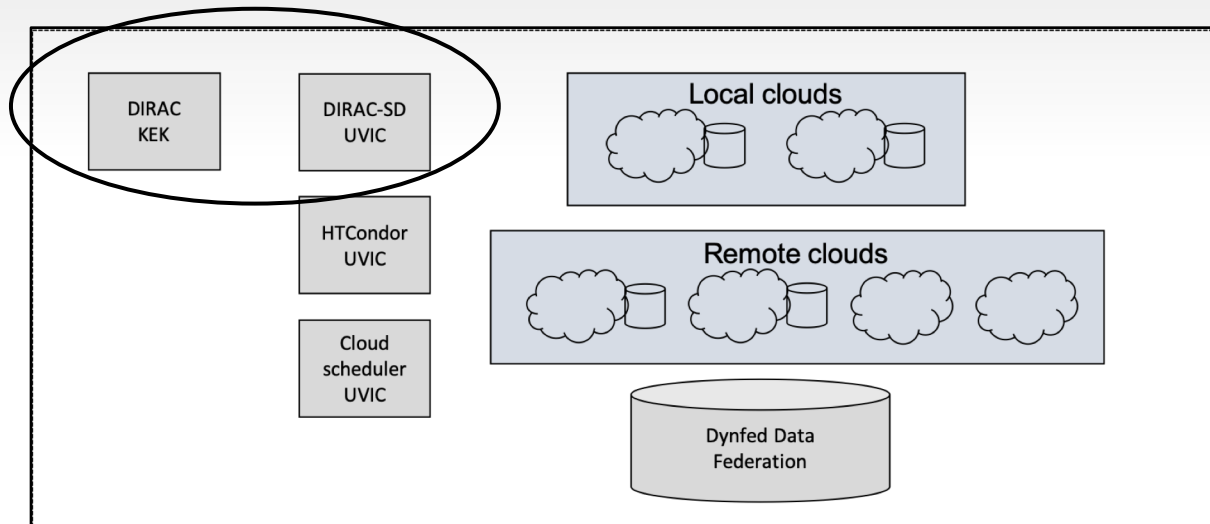
Jobs

Current status



# Belle II

- Uses the DIRAC workload management system
  - 3 DIRAC “Site Directors” (SD) in Victoria
  - SD submits pilot jobs to HTCondor based on queues in KEK
- CS system shared with ATLAS (separate HTCondor systems)
- Storage element (grid storage), Ceph (object storage)
- Dynafed federation for managing remote cloud storage

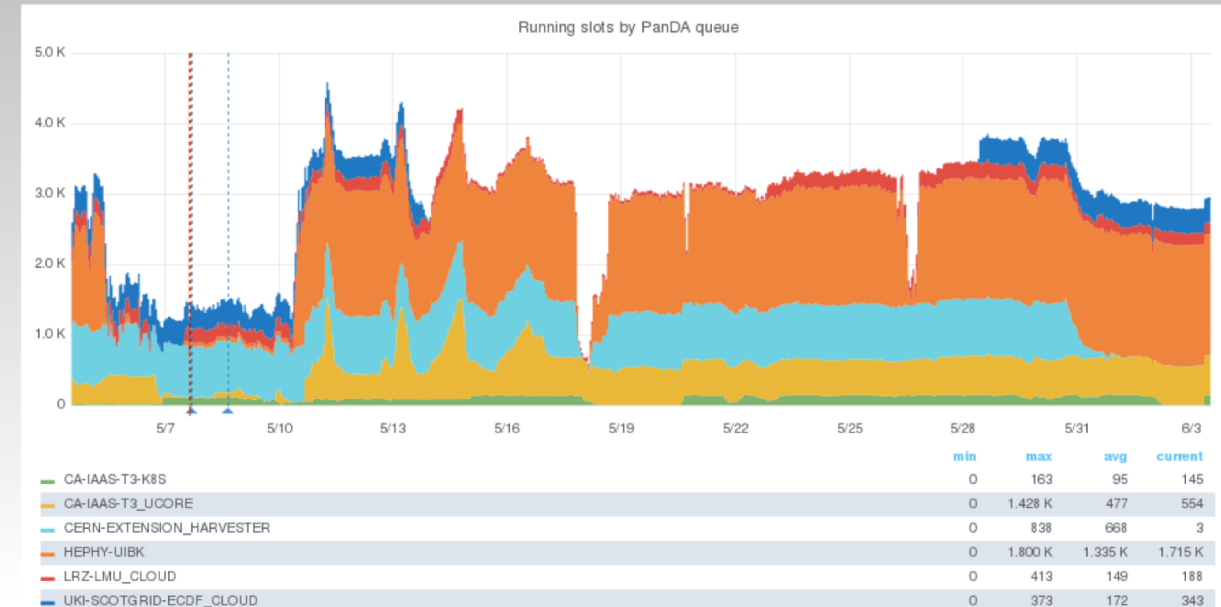


Distributed cloud provides 10% of B2 computing

Canadian Raw Data Centre will store 15% of the 2<sup>nd</sup> copy of the raw data in 2021 (using the “cloud” infrastructure)

# ATLAS

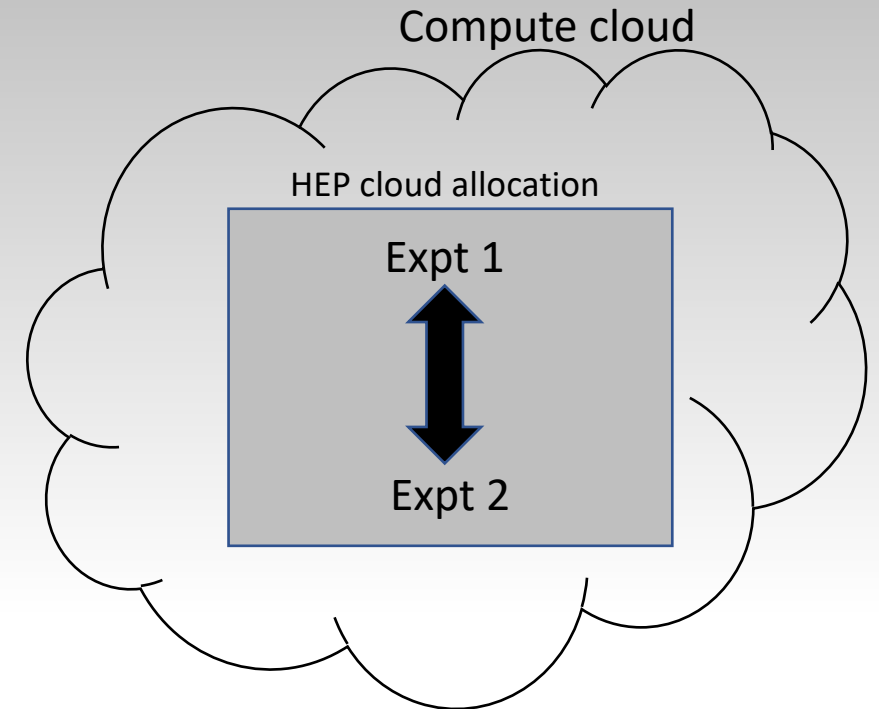
- Two systems:
  - CSV2 in Victoria for Belle II/ATLAS using clouds in North America
  - CSV1 in CERN for ATLAS using clouds in Europe
- Separate PanDa queues for each European cloud
  - Each site gets credit for their resources
- Support clouds in Edinburgh (ECDF), Munich (LRZ), Innsbruck and CERN (around 3000 cores)
- Plan to migrate all clouds to CSV2 in Victoria
  - Belle II HTCondor instance
  - Multiple ATLAS HTCondor instances



Distributed cloud provides 1% of ATLAS resources  
(comparable to a Canadian Tier-2 site)

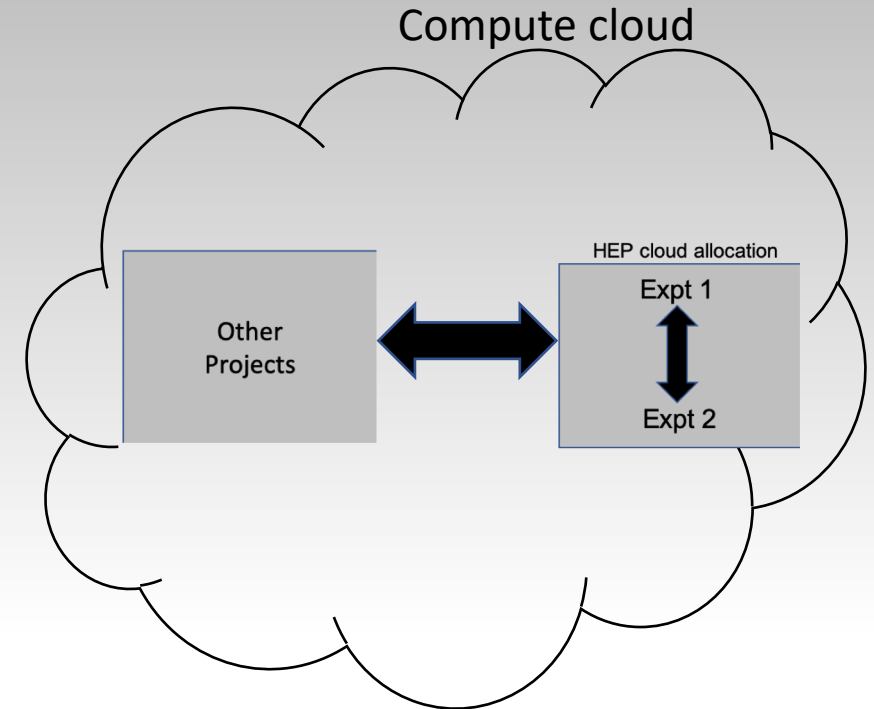
## Opportunistic use of idle HEP resources

- We have clouds for ATLAS, Belle II or HEP
- There are times when the workload on one project is low
- If the cloud is open to opportunistic use, then we shift the idle resources to the other experiment
- Once jobs appear, the opportunistic VMs are retired, the jobs are allowed to complete and the VMs destroyed
  - Typical turnover in 6-12 hours
- CSV2 can be configured for opportunistic use within a cloud
  - Can use a configurable fraction of the other resources



## Opportunistic use of idle non-HEP resources

- The goal is to use the idle resources of other projects
- We have no information on the resources allocated to other users
  - Exploring what information is required
  - How to have the cloud-operators pass it to CSV2
  - Automated system rather than manual
- Consider how to manage these resources:
  - Gentle transition (letting jobs finish)
  - Immediate termination
  - Single event generation



# Accounting

## Experiments:

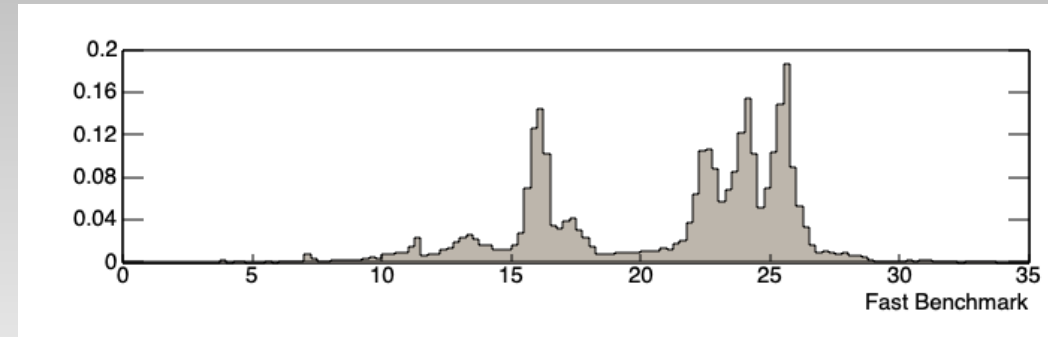
- Belle II keeps track within DIRAC using the fast-benchmark (DB12)
- ATLAS resources are tracked with PanDA

## Distributed cloud

- Measure the CPU performance of each VM at boot (DB12 with HEPiX normalization)
- Record the CPU/wallclock time and write to ES

## Apel Accounting

- Now storing the variables required for Apel accounting
- Written the code to extract and upload the information to MariaDB
- Developing script to upload to Apel database
- Should be operational Nov 2019
- Credit assigned to group, cloud or country



***We do not know the underlying hardware***

HEPSpec06 takes hours to run (license issues)

Run the fast-benchmark (DB12) at boot on each VM

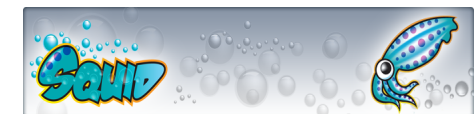
Not an ideal benchmark – waiting for new one

## Distributed cloud (using cloudscheduler) is a platform of external and in-house components

Preference is to use external components and develop those that do not exist

### External services

CernVM  
CVMFS  
DB12 and HEPspec06 Benchmarks  
Squids  
DynaFed data federator  
HTCondor job scheduler  
Elasticsearch database  
Apel accounting  
Ceph object storage  
Openstack/EC2/GCE/Azure clouds  
MariaDB  
Docker/Singularity containers  
HTTP/WebDAV  
InfluxDB



# Summary

The distributed compute cloud using cloudscheduler and HTCondor has been running production jobs for ATLAS and Belle II for many years  
(also used for Canadian astronomy workloads)

We have utilized dedicated and opportunistic, research and commercial clouds  
All usage to be reported to Apel accounting for each site/region

The system has undergone an extensive revision (cloudscheduler V2, MariaDB, ..)  
Significant improvement in reliability, responsiveness and management  
Opportunistic sharing of HEP resources (and non-HEP resources soon)

The Canadian Belle II Raw Data Centre will use this technology

Related talks at this conference:

Using Containers for managing infrastructure in ATLAS

Using Dynafed as a distributed storage element

Sim@P1 (ATLAS), a cloud system utilizing the HLT cluster at CERN

## Dynafed

The Dynafed data federator as grid site storage element, F.Berghaus, CHEP2019

Using a dynamic data federation for running Belle-II simulation applications in a distributed cloud environment, CHEP2018

<https://doi.org/10.1051/epjconf/201921404026>

Integrating a dynamic data federation into the ATLAS distributed data management system, CHEP 2018,

<https://doi.org/10.1051/epjconf/201921407009>

## cloudscheduler

High-throughput cloud computing with the cloudscheduler VM provisioning service. Submitted to Software and Computing in Big Science 2019

<http://heprcdocs.phys.uvic.ca/papers/cs.pdf>

Quasi-online accounting and monitoring system for distributed clouds, CHEP 2018,

<https://doi.org/10.1051/epjconf/201921407035>

## sim@P1

Title: ATLAS Sim@P1 upgrades during long shutdown two, F. Berghaus CHEP 2019

## Containers

Using Kubernetes as an ATLAS computing site, CHEP 2019