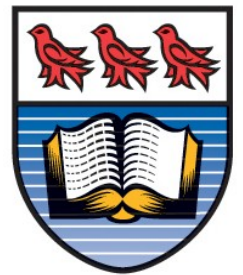# Dynamic Resource Provisioning for Batch Computing in the Cloud

## Frank Berghaus

(University of Victoria)

On behalf of the ATLAS Cloud Computing Group
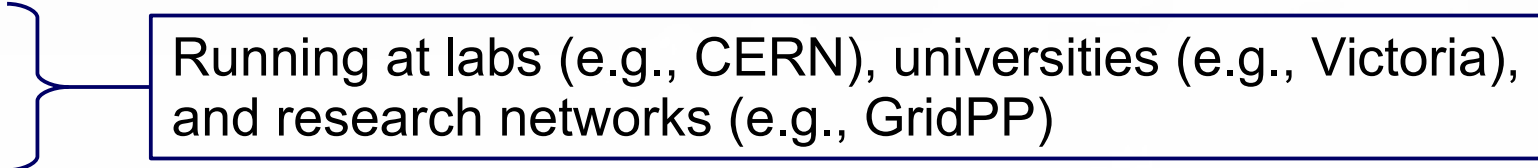
Belle II

ATLAS EXPERIMENT

University of Victoria

# Overview

- Overview of the Cloud Scheduler System

- Worker Node Virtual Machines

- Batch Configuration for Dynamic Job Requirements:

  - Single Core, Multi-Core, High Memory, etc.

- Current Deployment

- Squid Discovery with Shoal

# Infrastructure-as-a-Service (IaaS) Clouds

- IaaS Cloud: A pool of virtual machine hypervisors presenting a single controller interface
  - Run many instances of one virtual machine configured for ATLAS computing

- Advantages:
  - Isolate complex application software from site administration
  - Minimize dependence on local system
  - Flexible resource allocation

- Examples:
  - OpenStack
  - Nimbus

  Running at labs (e.g., CERN), universities (e.g., Victoria), and research networks (e.g., GridPP)

  - Commercial clouds: Amazon, Google, etc.

# Cloud Scheduler

- Cloud Scheduler is a python package for managing VMs on IaaS clouds
- Users submit HTCondor jobs
  - Optional attributes specify virtual machine properties
- Dynamically manages quantity and type of VMs in response to user demand
- Easily connects to many IaaS clouds, and aggregates their resources
- Provides IaaS resources in the form of an ordinary HTCondor batch system
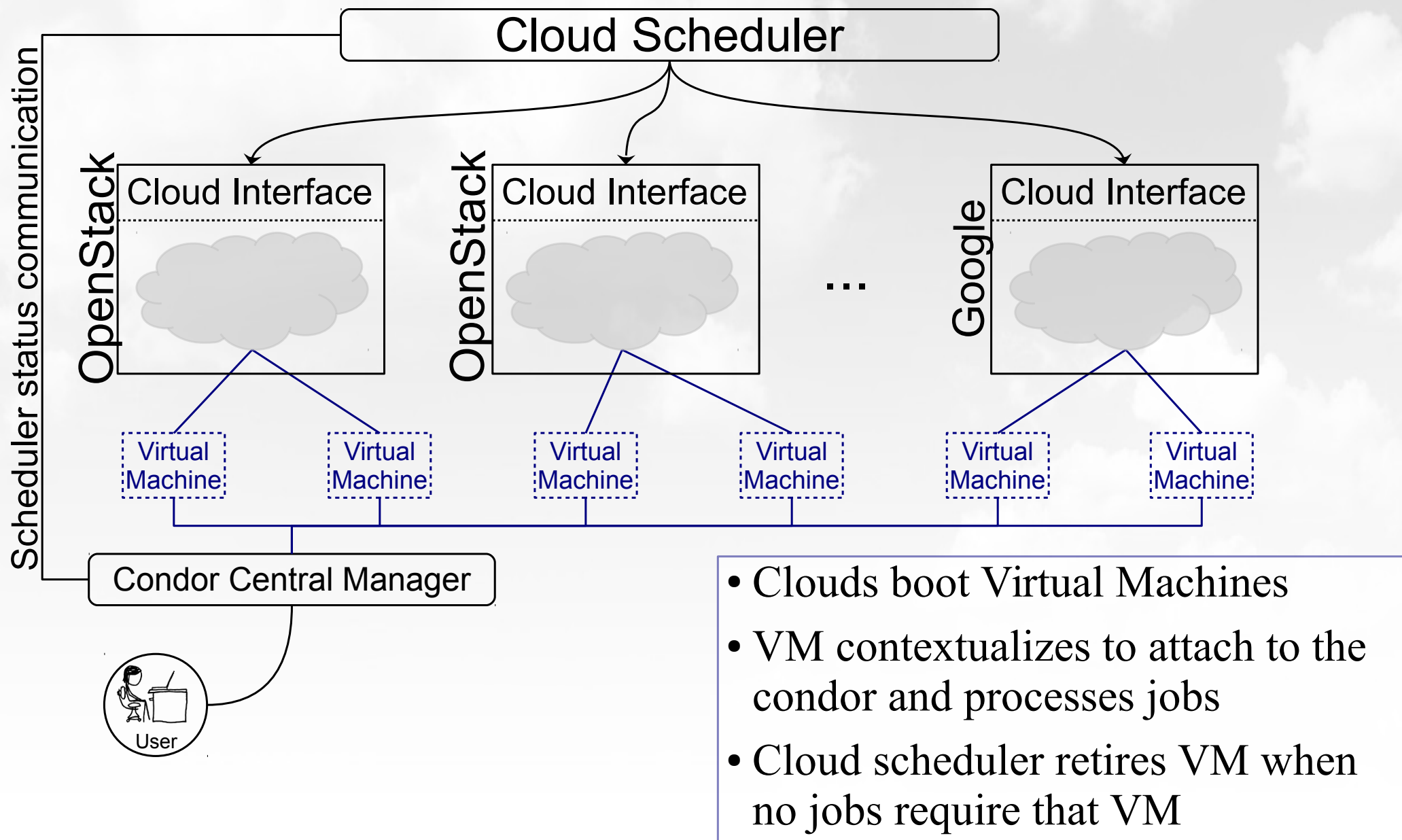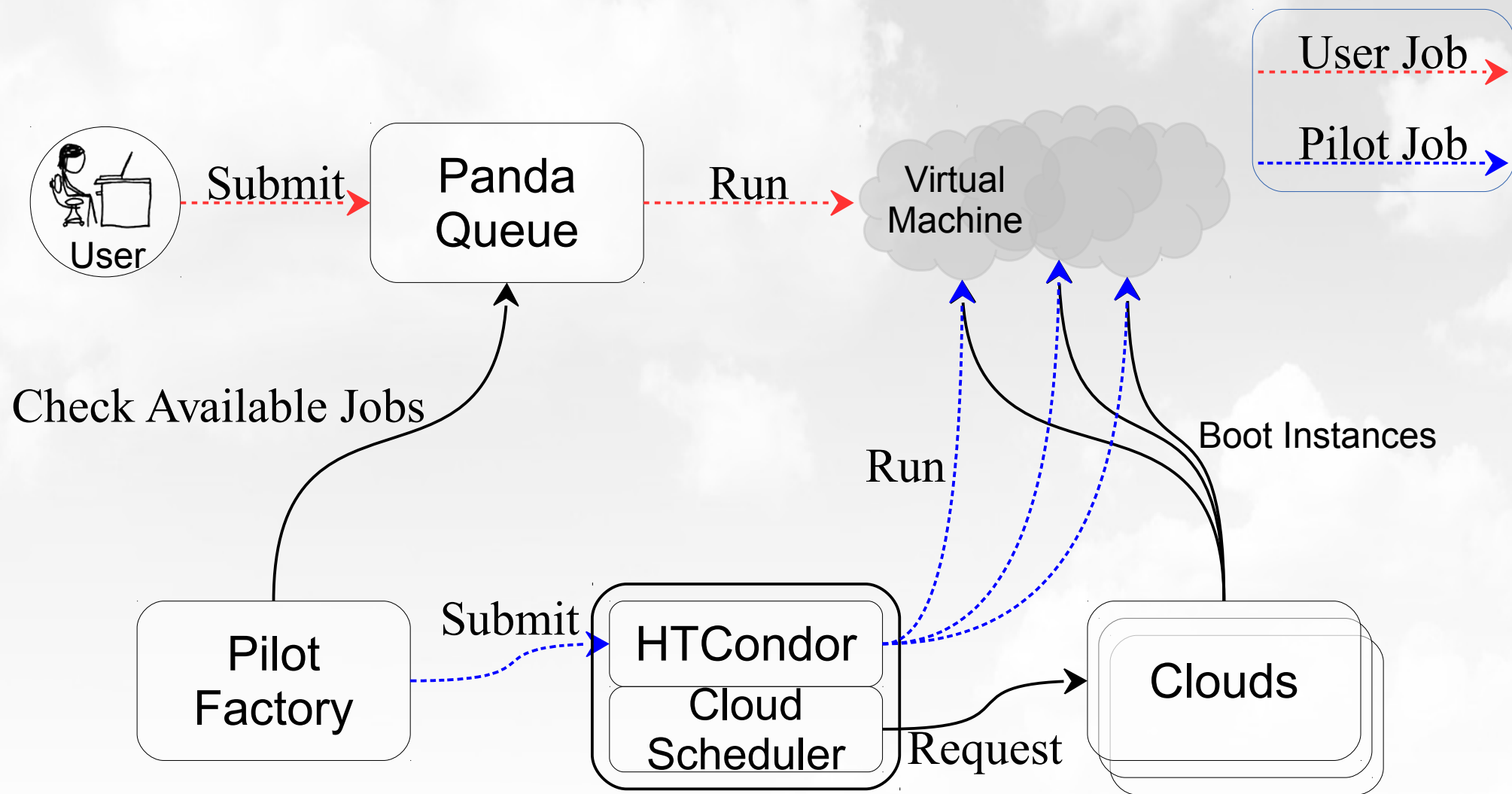- Used by ATLAS, Belle II, CANFAR, and BaBar

| | |
|---|---|
| Code | https://github.com/hep-gc/cloud-scheduler |
| Website | http://cloudscheduler.org/ |
| Publication | http://arxiv.org/abs/1007.0050 |

# Cloud Scheduler



- Clouds boot Virtual Machines
- VM contextualizes to attach to the condor and processes jobs
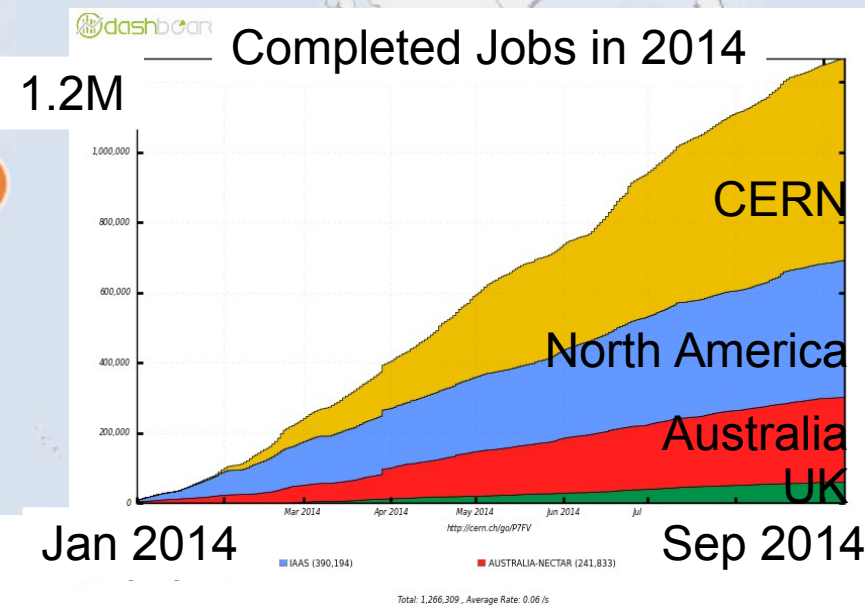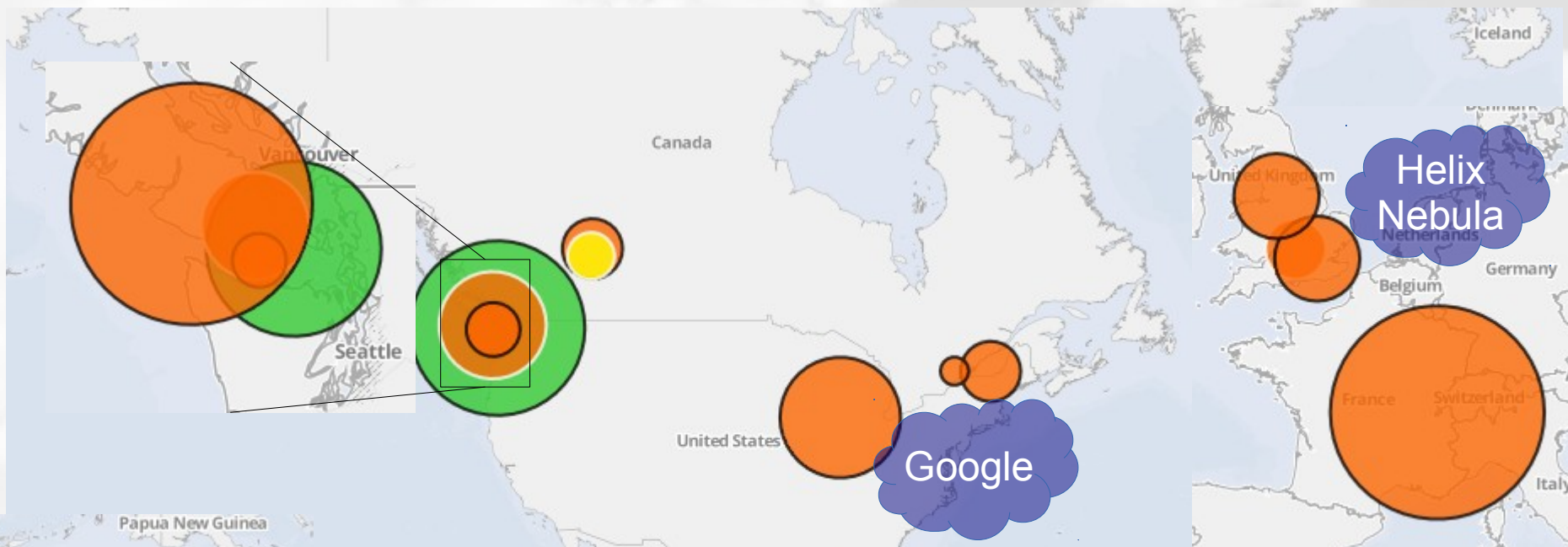- Cloud scheduler retires VM when no jobs require that VM

# Cloud Job Flow (on the Grid)



- Easy to connect and use many clouds
- Integrated with DIRAC as well as Panda

# ATLAS Cloud Production in 2014



Helix Nebula

Google

CERN

North America

Australia

UK

**Completed Jobs in 2014**

1.2M

1,000,000

800,000

600,000

400,000

200,000

Mar 2014   Apr 2014   May 2014   Jun 2014   Jul

http://cern.ch/go/P7FV

Jan 2014

Sep 2014

IAAS (390,194)   AUSTRALIA-NECTAR (241,833)

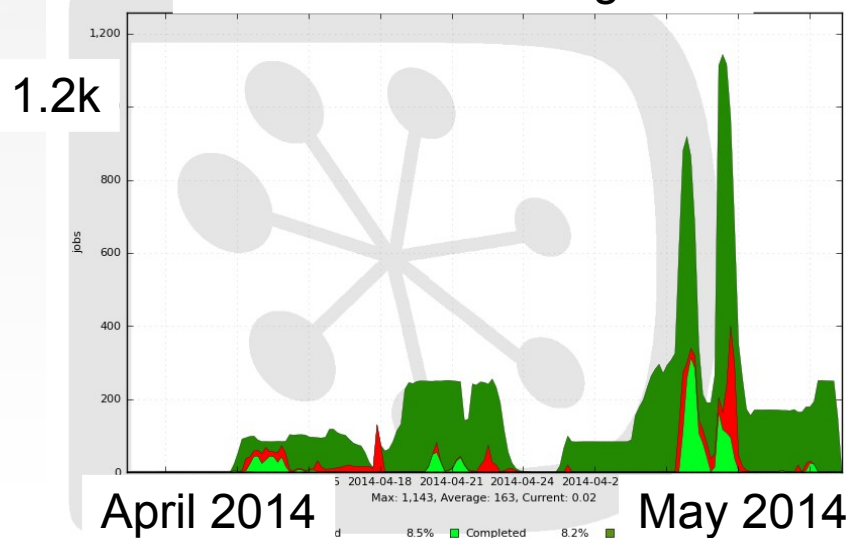Total: 1,266,309 , Average Rate: 0.06 /s

- Over 1.2M ATLAS Jobs Completed
- Mostly Single core production
- New ATLAS requirements:
  - Multi-core
  - High memory

# Belle II Cloud Production in 2014



Amazon

Concurrent Running Jobs

1.2k

April 2014    May 2014

- Supporting the DIRAC WMS

- Completed 400k Belle II production Jobs

- Using up to 1000 cores on Amazon EC2

# CernVM3Worker Nodes

- Features:

  - Operating system and project software is made available over cvmfs

  - Cloud-init and puppet contextualize images on boot

    https://github.com/berghaus/atlasgce-modules

    - Developed by Frank Berghaus and Henric Öhman

  - Same image works anywhere, on any Hypervisor, and on any Cloud Type

- Old static condor configuration (8-core VM):

```
NUM_SLOTS = 8

SLOT1_USER = slot01

SLOT2_USER = slot02

...

SLOT8_USER = slot08

DEDICATED_EXECUTE_ACCOUNT_REGEXP = slot[0-9]+
```

# CernVM3Worker Nodes

- Features:

  - Operating system and project software is made available over cvmfs

  - Cloud-init and puppet contextualize images on boot

    https://github.com/berghaus/atlasgce-modules

    – Developed by Frank Berghaus and Henric Öhman

  - Same image works anywhere, on any Hypervisor, and on any Cloud Type

- Dynamic Condor Slot Configuration:

```
NUM_SLOTS_TYPE_1 = 1

SLOT_TYPE_1 = cpus=100%

SLOT_TYPE_1_PARTITIONABLE = True

SLOT1_1_USER = slot01

...

SLOT1_8_USER = slot08

DEDICATED_EXECUTE_ACCOUNT_REGEXP = slot[0-9]+
```

# Dynamic Batch Jobs

- **Jobs** need to specify resource requirements:

```
request_cpus    =  2

request_memory  =  5000     # in mbytes

request_disk    =  10000000 # in kbytes
```

- Condor creates a dynamic slot of appropriate size on a worker node with sufficient resources

- Problem:

  - Jobs with small resource needs **fragment** worker nodes into slots with small resources

  - Jobs with large resource requirements can not run on fragmented worker nodes

- Solution: The **DEFRAG** daemon

# Dynamic Batch Slots

- **Defrag** daemon cleans up unused dynamic slots (*testing*):
  - Thanks to Andrew Lahiff from RAL for their configuration details

```
DAEMON_LIST = DEFRAG

DEFRAG_INTERVAL = 600

DEFRAG_DRAINING_MACHINES_PER_HOUR = 30.0

DEFRAG_MAX_CONCURRENT_DRAINING = 60

DEFRAG_MAX_WHOLE_MACHINES = 300

DEFRAG_SCHEDULE = graceful

DEFRAG.SETTABLE_ATTRS_ADMINISTRATOR =
DEFRAG_MAX_CONCURRENT_DRAINING,DEFRAG_DRAINING_MACHINES_PER_HOUR,DEFRAG_MAX_WHOLE_MACHINES

ENABLE_RUNTIME_CONFIG = TRUE

DEFRAG_RANK = ifThenElse(Cpus >= 8, -10, (TotalCpus - Cpus)/(8.0 – Cpus))
```

# Defining Target Shares

- Use condor groups to prioritize job types

```
GROUP_NAMES = group_analysis, group_production
GROUP_QUOTA_DYNAMIC_group_analysis = 0.05
GROUP_QUOTA_DYNAMIC_group_production = 0.95
GROUP_ACCEPT_SURPLUS = True
```

- In the job definition add:

```
AccountingGroup = "group_production" or
AccountingGroup = "group_analysis"
```

  - Thanks to Joanna Huang and Sean Crosby from the Australian ATLAS group

# Relevant Tools for Distributed Computing

# Shoal: Dynamic Squid Discovery

## List of Active Squids

**4 active in the last 180 seconds**

| # | Hostname | Public IP | Private IP | Bytes Out | City | Region | Country | Latitude | Longitude | Last Received | Alive | Verified | Access Level |
|---|----------|-----------|-----------|-----------|------|--------|---------|----------|-----------|---------------|-------|----------|--------------|
| 1 | atlascaq3.triumf.ca | 142.90.110.68 | | 25 kB/s | Vancouver | | Canada | 49.2765 | -123.2177 | 9s | 29h34m39s | ✔ | Global |
| 2 | chrysaor.westgrid.ca | 206.12.48.3 | 172.22.5.2 | 16034 kB/s | Vancouver | | Canada | 49.2836 | -123.1041 | 19s | 29h33m34s | ✔ | Global |
| 3 | atlas-squid.cern.ch | 128.142.200.105 | | 0 kB/s | Geneva | | Switzerland | 46.1956 | 6.1481 | 23s | 29h34m37s | ✘ | Global |
| 4 | t2software02.physics.ox.ac.uk | 163.1.5.127 | | 2 kB/s | Oxford | | United Kingdom | 51.75 | -1.25 | 29s | 29h33m34s | ✔ | Global |

- Ready for larger scale deployment: installation instructions

- Current server: http://shoal.heprc.uvic.ca/

- Connected squids: UVic, TRIUMF, Oxford, CERN Cloud

- Included with CernVM since release 3.2

- Meets the requirements of the squid discovery task force

# EMI Dynamic Federation
## Thanks to Fabrizio Furano and Ryan Taylor

- High-performance

  - aggressive metadata caching in RAM

  - maximal concurrency

  - scalable

    - ~ 10^6 hit/s per core

    - 24 GB of RAM for metadata cache is enough for ~100 PB of data in end-points

- Well-designed

  - stateless, no persistency

  - standard components and protocols, not HEP-specific

  - general-purpose solution; could be adopted by multiple experiments

  - trivial to add endpoints; **no site action needed!**

- Data access

  - automatically download from nearest endpoint, or

  - download from all endpoints simultaneously (metalink + aria2)


EUROPEAN MIDDLEWARE INITIATIVE

# Summary & Outlook

- ATLAS and Belle II Production is running on IaaS clouds
  - Over 1.2M ATLAS jobs completed in 2014
  - Dynamically allocating resources for single and multi core job requirements
  - Planning to test high memory jobs
- Dynamic resource allocation allows quick creation of necessary resources
- Aggregating many computing resources into few batch queues
- Share resources between projects
- Using micro-kernel CernVM3
- Automated squid discovery for cvmfs
- Deploying Dynamic Federation as data access solution

# Backup

# Federator Deployment

- Simple, lightweight. Easy to set up (a few person-days of effort)

- Contains all SEs in CA and AU

## /myfed/atlas/

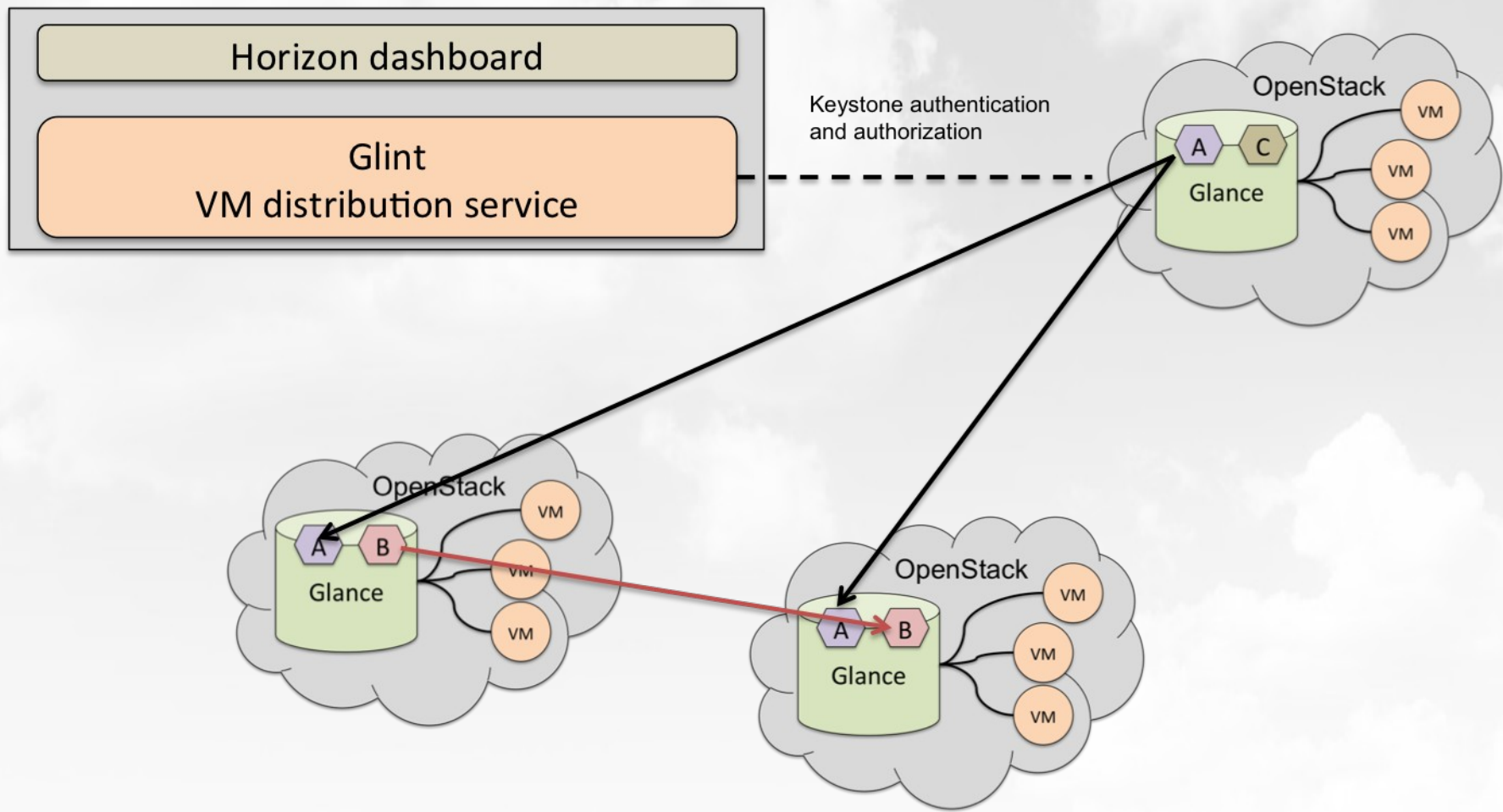| Mode | UID | GID | Size | Modified | Name |
|------|-----|-----|------|----------|------|
| drwxrwxrwx | 0 | 0 | 0 | Fri, 13 Jan 2012 09:48:30 GMT | /atlasdatadisk/ |
| drwxrwxrwx | 0 | 0 | 0 | Wed, 03 Sep 2014 00:37:31 GMT | atlas/ |
| drwxrwxrwx | 0 | 0 | 0 | Fri, 13 Jan 2012 09:48:30 GMT | atlasdatadisk/ |
| drwxrwxrwx | 0 | 0 | 0 | Wed, 22 Jan 2014 09:53:40 GMT | atlasdatadiskrucio/ |
| drwxrwxrwx | 0 | 0 | 0 | Fri, 13 Jan 2012 09:48:30 GMT | atlasgroupdisk/ |
| drwxrwxrwx | 0 | 0 | 0 | Fri, 13 Jan 2012 09:48:31 GMT | atlashotdisk/ |
| drwxrwxrwx | 0 | 0 | 0 | Fri, 02 Aug 2013 21:18:18 GMT | atlaslocalgroupdisk/ |
| drwxrwxrwx | 0 | 0 | 0 | Tue, 18 Jun 2013 19:38:45 GMT | atlasmcdisk/ |
| drwxrwxrwx | 0 | 0 | 0 | Fri, 13 Jan 2012 09:48:33 GMT | atlasproddisk/ |
| drwxrwxrwx | 0 | 0 | 0 | Sun, 01 Dec 2013 02:16:12 GMT | atlasscratchdisk/ |
| drwxrwxrwx | 0 | 0 | 0 | Tue, 31 Mar 2009 13:35:48 GMT | atlasuserdisk/ |
| drwxrwxr-x | 0 | 0 | 0 | Tue, 24 Nov 2009 10:59:12 GMT | au/ |
| -rwxrwxrwx | 0 | 0 | 1000.0M | Fri, 19 Nov 2010 21:00:01 GMT | file1Gc1 |
| drwxrwxrwx | 0 | 0 | 0 | Mon, 24 Mar 2014 16:08:14 GMT | generated/ |
| drwxrwxrwx | 0 | 0 | 0 | Fri, 18 Dec 2009 13:50:46 GMT | install/ |
| -rwxrwxrwx | 0 | 0 | 998 | Tue, 22 Jul 2014 20:24:24 GMT | junk.weiyang |
| drwxrwxr-x | 0 | 0 | 0 | Mon, 22 Apr 2013 21:01:19 GMT | lucien/ |
| drwxrwxrwx | 0 | 0 | 0 | Tue, 22 Jul 2014 20:34:11 GMT | rucio/ |
| -rwxrwxrwx | 0 | 0 | 1.0M | Tue, 22 Jul 2014 20:31:09 GMT | snderitu:user.ivukotic.xrootd.ca-mcgill-clumeq-t2-1M |
| -rwxrwxrwx | 0 | 0 | 2.4K | Tue, 22 Oct 2013 18:47:16 GMT | test1copy |
| -rwxrwxrwx | 0 | 0 | 2.4K | Tue, 22 Oct 2013 18:48:16 GMT | test1copy2 |
| -rwxrwxrwx | 0 | 0 | 2.0M | Wed, 02 Jul 2014 22:19:05 GMT | test2 |
| -rwxrwxrwx | 0 | 0 | 0 | Wed, 02 Jul 2014 22:11:53 GMT | test3 |
| drwxr-xr-x | 0 | 0 | 0 | Mon, 05 Aug 2013 11:14:36 GMT | testWebDAV/ |
| drwxrwxrwx | 0 | 0 | 0 | Mon, 03 Mar 2014 20:50:03 GMT | test_belle/ |
| -rw-rw-r-- | 0 | 0 | 20 | Tue, 09 Mar 2010 02:11:42 GMT | testfile-put-1268100565-35d990bee619.txt |
| -rw-rw-r-- | 0 | 0 | 20 | Tue, 09 Mar 2010 02:25:28 GMT | testfile-put-1268101507-57176c95c28b.txt |
| -rw-rw-r-- | 0 | 0 | 20 | Tue, 09 Mar 2010 02:55:45 GMT | testfile-put-1268103324-fa99a035fb3e.txt |
| -rw-rw-r-- | 0 | 0 | 20 | Tue, 09 Mar 2010 03:55:51 GMT | testfile-put-1268106930-908e9758ee79.txt |
| -rw-rw-r-- | 0 | 0 | 20 | Fri, 02 Jul 2010 09:17:05 GMT | testfile-put-1278062170-957d70604fb5.txt |
| drwxrwxr-x | 0 | 0 | 0 | Tue, 24 Nov 2009 10:57:55 GMT | users/ |

# Federator Software Components

- Uniform Generic Redirector (UGR)

  - Core component containing all federation logic

  - Integrated as a plugin of Apache

- Apache HTTP server

  - Frontend to clients

  - Handles client redirection

- Memcached

  - for 2nd-level shared metadata caching in RAM

- DMLite

  - Name translations for unifying grid storage endpoints

# Glint: Image Distribution Service



- Addition to OpenStack (see November OS summit)
- Relies on **glance** for image management and **keystone** for authentication
- User interface in horizon dashboard
- Works on OpenStack, Amazon EC2, Google's GCE, and Nimbus