



U.S. DEPARTMENT OF
ENERGY

Office of
Science

FermiCloud On-Demand Services: Data-intensive Computing on Public and Private Clouds

Steven C. Timm, Gabriele Garzoglio

Grid and Cloud Services Department, Scientific Computing Division,
Fermilab

HEPiX—Spring 2014—Annecy-le-Vieux, FR

Overview

- Cloud Program of work at Fermilab
- FermiCloud Project and On-Demand Services
- Recent Cloud Research and Development Results
- Neutrino production workflows on FermiCloud and AWS
 - Fermilab neutrino applications
 - Code relocation to CVMFS
 - Workflow management and data handling
 - Running on distributed Grid sites
 - Running on FermiCloud and Amazon AWS
 - Relative performance and cost
 - Best practices
- Current and future research directions
 - Provisioning of complex services
 - VM Batch Submission Service

The Cloud Program Team

- **Fermilab Grid And Cloud Services Department**
 - Gabriele Garzoglio (Dept. Head)
 - Steven Timm (Assoc. Dept. Head for Cloud Computing, FermiCloud Project lead)
 - Joe Boyd (New group leader of Grid and Cloud Services Operations Group)
 - Gerard Bernabeu Altayo (Deputy group leader of Grid and Cloud Services Operations Group and lead of operations for cloud services)
 - Hyun Woo Kim (lead FermiCloud development)
 - Tanya Levshina, Parag Mhashilkar, Kevin Hill, Karen Shepelak (Technical Support)
 - Keith Chadwick (Scientific Computing Division Technical Architect)
- Scientific Server Support Group (Ed Simmonds, Leader) manages cloud servers.
- Enterprise Services Department (Mike Rosier, Head) manages our SAN.
- **KISTI-GSDC team**
 - Haeng-Jin Jang, Seo-Young Noh, Gyeongryoon Kim
- **Funded through Collaborative Agreement with KISTI**
 - Nick Palombo (consultant) (summer 2013)
 - Tiago Pais (2013 summer students, IIT)
 - Hao Wu (IIT PhD. Student on FermiCloud with us since summer of 2013)
 - Francesco Ceccarelli (2013 summer student, INFN Italy)
 - 3 new summer students coming this summer plus another consultant
- **Professors at Illinois Institute of Technology (IIT)**
 - Ioan Raicu, Shangping Ren

Cloud Program of Work at Fermilab

FermiCloud Project—

Delivers high-level cloud-hosted on-demand services for Fermilab scientists, focus of this talk.

Fermilab-KISTI Cooperative R+D Agreement

Focused on pushing development of cloud middleware and federation technology to the benefit of both labs and others.

One main goal for this summer:

– Federated Cloud between FermiCloud (OpenNebula) /gCloud@KISTI(OpenStack) /Amazon (AWS)

Other R+D activities

Non-FNAL communities on FermiCloud

White paper with K. Keahey (ANL) and J. Lauret (BNL)

STAR is first example community.

Policy-based Tiered Provisioning:

Run first locally, then on dedicated grid site, then opportunistic on Grid and Cloud.

Allocation for multiple communities on Amazon Web Services

One AWS account per group?

FermiCloud obtain one account and charge back? Other?

Submitting to native “Nova” API of OpenStack.

Several community OpenStack clouds available in US,

Rackspace working to facilitate participation.

Cloud Computing Environment—next directions for security in the cloud

Recent publications

- Automatic Cloud Bursting Under FermiCloud
 - ICPADS CSS workshop, Dec. 2013.
- A Reference Model for Virtual Machine Launching Overhead
 - CCGrid conference, May 2014
- Grids, Virtualization, and Clouds at Fermilab
 - CHEP conference, Oct. 2013.
- Exploring Infiniband Hardware Virtualization in OpenNebula towards Efficient High-Performance Computing.
 - CCGrid Scale workshop, May 2014.
- SUBMITTED: A Reference Model for Virtual Machine Launching Overhead, submitted to IEEE Transactions on Cloud Computing, May 2014

- Measurements at all frontiers – Electroweak physics, neutrino oscillations, muon g-2, dark energy, dark matter
- 8 major experiments in 3 frontiers running simultaneously in 2016
- Sharing both beam and computing resources
- Impressive breadth of experiments at FNAL



	FY10	FY11	FY12	FY13	FY14	FY15	FY16	FY17	FY18	FY19	FY20	FY21	FY22	FY23
Physics: Particle Experiments				Operations (future experiments with C-0 or higher level approval)			Data analysis continues							
Intensity Frontier	v: LBNE													
	μ: Mu2e													
	e: Mu2e g-2													
	v: NOvA													
	v: MicroBooNE													
	v: MINOS+													
	neutrino: SeaQuest													
	v: MINERvA													
	v: MINOS													
	v: MiniBooNE													
<i>Not included are the ORKA kaon experiment which received the Stage 1 approval from Fermilab, and experiments such as nuSTORM, proton EDM and neutron-antineutron oscillation experiments which are currently developing proposals with the encouragement of Fermilab Physics Advisory Committee.</i>														
Energy Frontier	LHC (14 TeV, Lum upgrade): CMS													
	LHC (14 TeV): CMS													
	LHC (7-8 TeV): CMS													
	Tevatron: CDF/D0													
Cosmic Frontier	DE: LSST													
	DM: Gen 3													
	Dark EnergyMS DESI													
	Dark Matter: Generation 2													
	5/23/2014													
	Dark Matter: Generation 1													
	Dark Energy: DES													

Motivation: Computing needs for IF / CF experiments

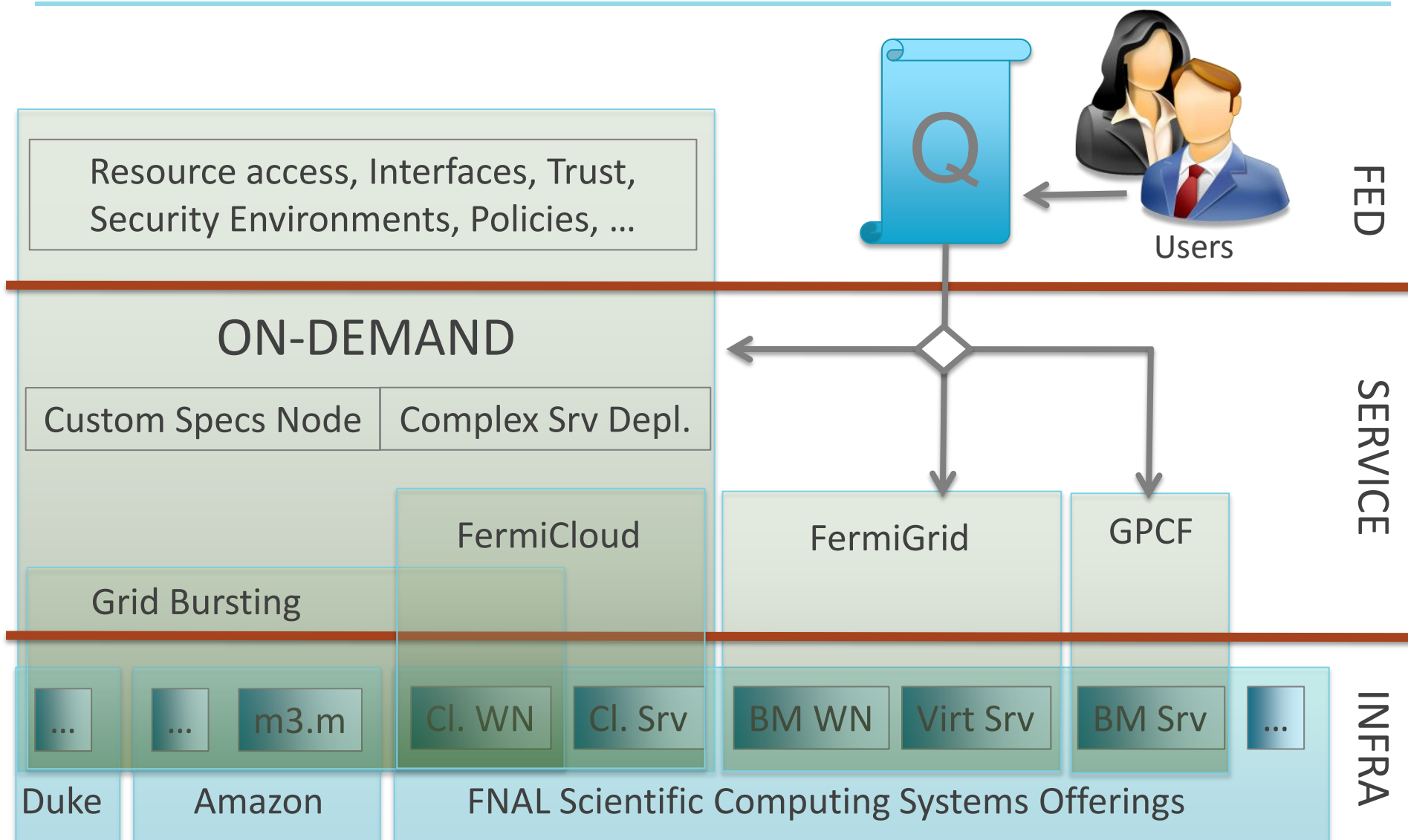
- Requests for increasing slot allocations
- Expect stochastic load, rather than DC
- Focus on concurrent peak utilization
- Addressing the problem for IF / CF experiments as an ensemble under FIFE
- Collaborating with CMS on computing and infrastructure

Fermilab Scientific Computing Review*

Experiment	Allocation (Slots)	Average Utilization (last quarter)		FY14 Request (Slots)
		Average Slots	Peak Slots	
ArgoNeuT	200	75	1712	0
Muon g-2	200	9	195	0
LBNE	500	30	1202	200
MARS	1225	212	1444	100
MicroBooNE	500	39	597	100
Minerva	1600	667	3580	500
Minos	1200	294	2427	0
Mu2e	500	597	2491	500
Nova	1300	410	1799	500
Total Requested	7225	2333	15447	1900

* Number of slots have been updated periodically since

Federation->Service->Infrastructure

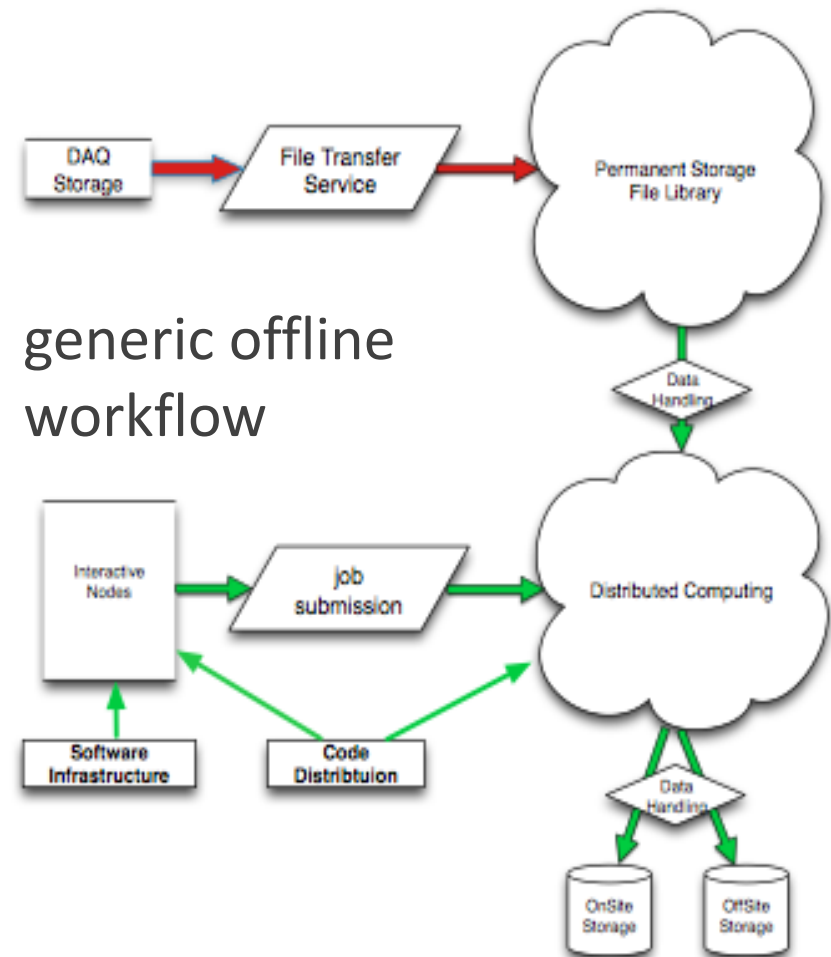


Neutrino Program Applications

- Full program of work
 - Beamline simulation with MARS
 - Detector simulation and calibration
 - Data reconstruction and analysis
- Special characteristics:
 - Flux Files for beam intensity, common to many runs, 1-10GB
 - Shower files for pattern recognition in data, common to many runs
 - Restricted external resources, databases for example.
 - Stochastic load—different experiments have different high demand periods
- Until 2013 all ran locally at Fermilab almost exclusively
 - Bluearc NAS device used for code distribution, log files, data storage, flux files
 - Not scalable in terms of money or for distributed computing.

FIFE—Fabric for Fronter Experiments

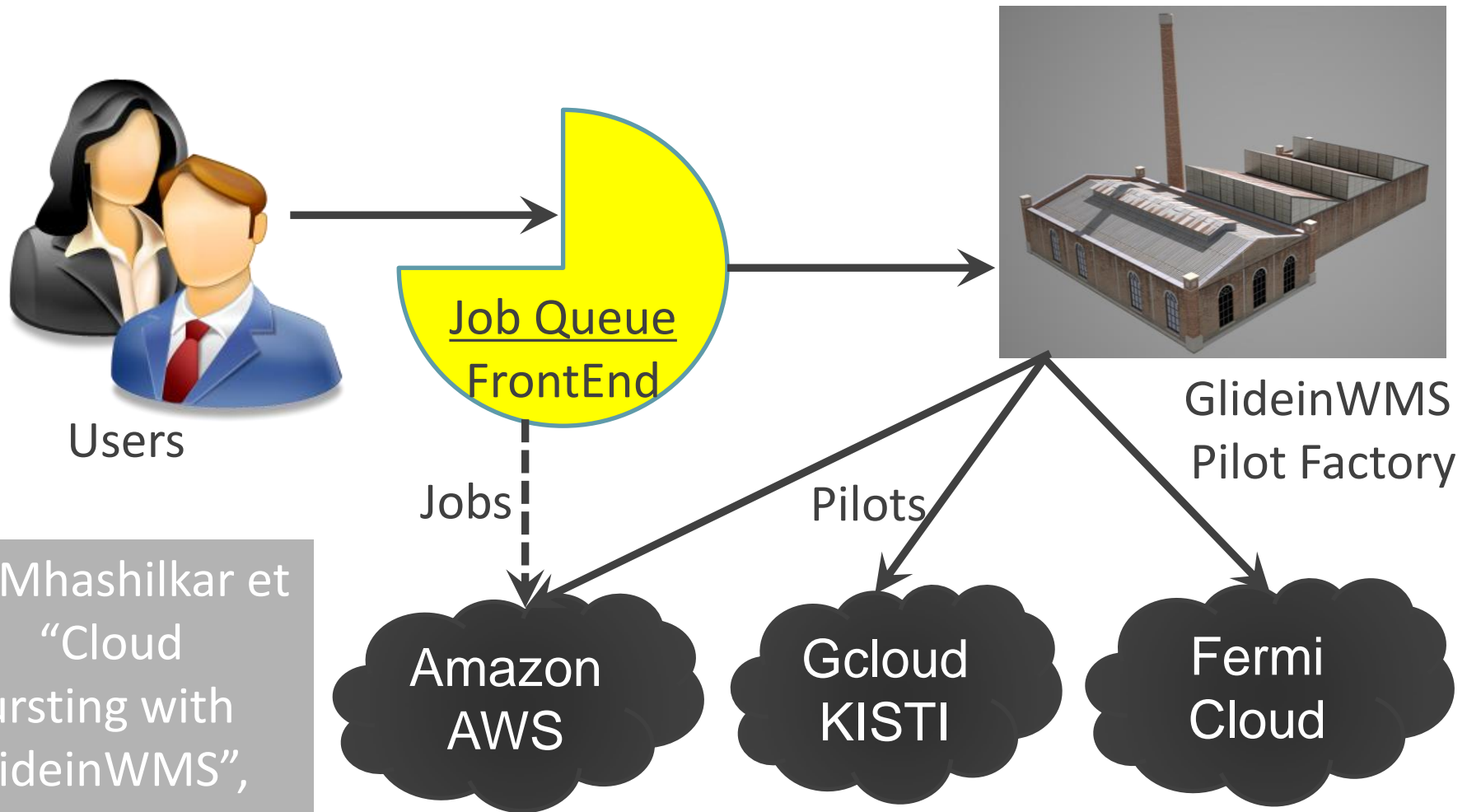
- Software Framework – ART supports integrated services (Data Handling, Geant4, ROOT)
- build environment and distribution
 - git, svn, gmake, cmake and distribution with CERN Virtual Machine File System (CVMFS)
 - build process and machine in development
- data handling - access to file catalog, tape storage, cache disk, and file transfer
- job submission infrastructure – based on GlideinWMS pilot system
- database infrastructure and access
- shared software (LArSoft for Liquid Argon TPC reconstruction)
- additional infrastructure - authentication, electronic control room log book, new user accounts



NOvA Code Port to CVMFS

- Central CVMFS server for Open Science Grid—
oasis.opensciencegrid.org
- NOvA put in all dependent packages, their own code, and some data (flux) files
- Included copies of system libraries missing at some grid sites, especially X11-related
- Spent a lot of time to make all the code relocatable so the same binaries could be used in CVMFS and locally.
- Spent effort to ensure all files readable by “other”

GlideinWMS – Grid Bursting



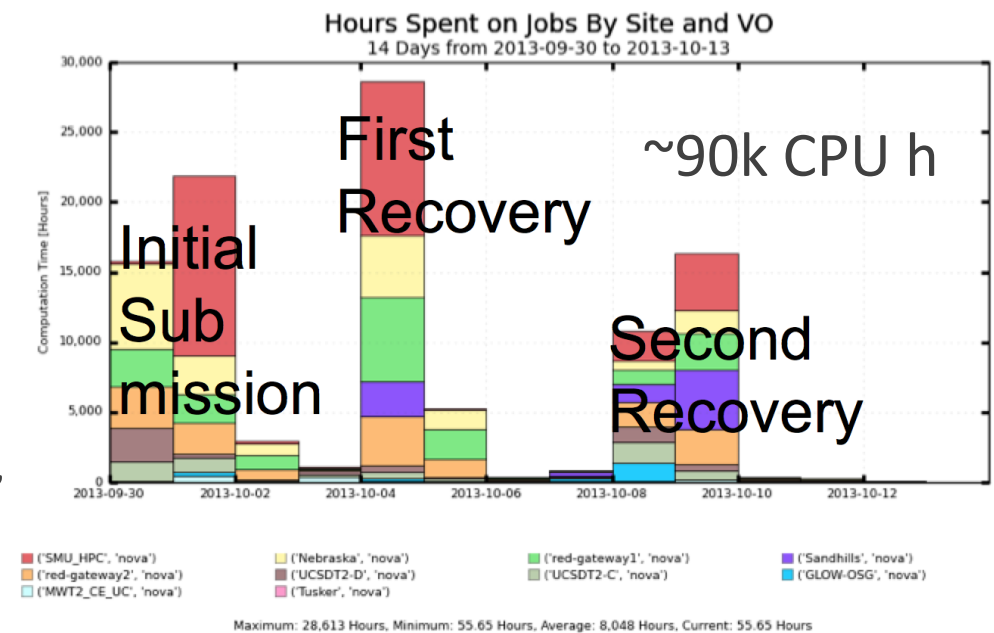
P. Mhashilkar et al, "Cloud Bursting with GlideinWMS", CHEP 2013

Workflow Management and Data Handling

- GlideinWMS used for workflow management
- Three dedicated and four opportunistic sites accessed through Open Science Grid
- In each site, some challenges at first to make sure we could write the local storage element, understand network topology, etc.
- Output from simulation stored first on local SE
- “ifdh” wrapper software surrounds srmcp / gridftp
- One job stays alive to push all data back to Fermilab, the others exit.

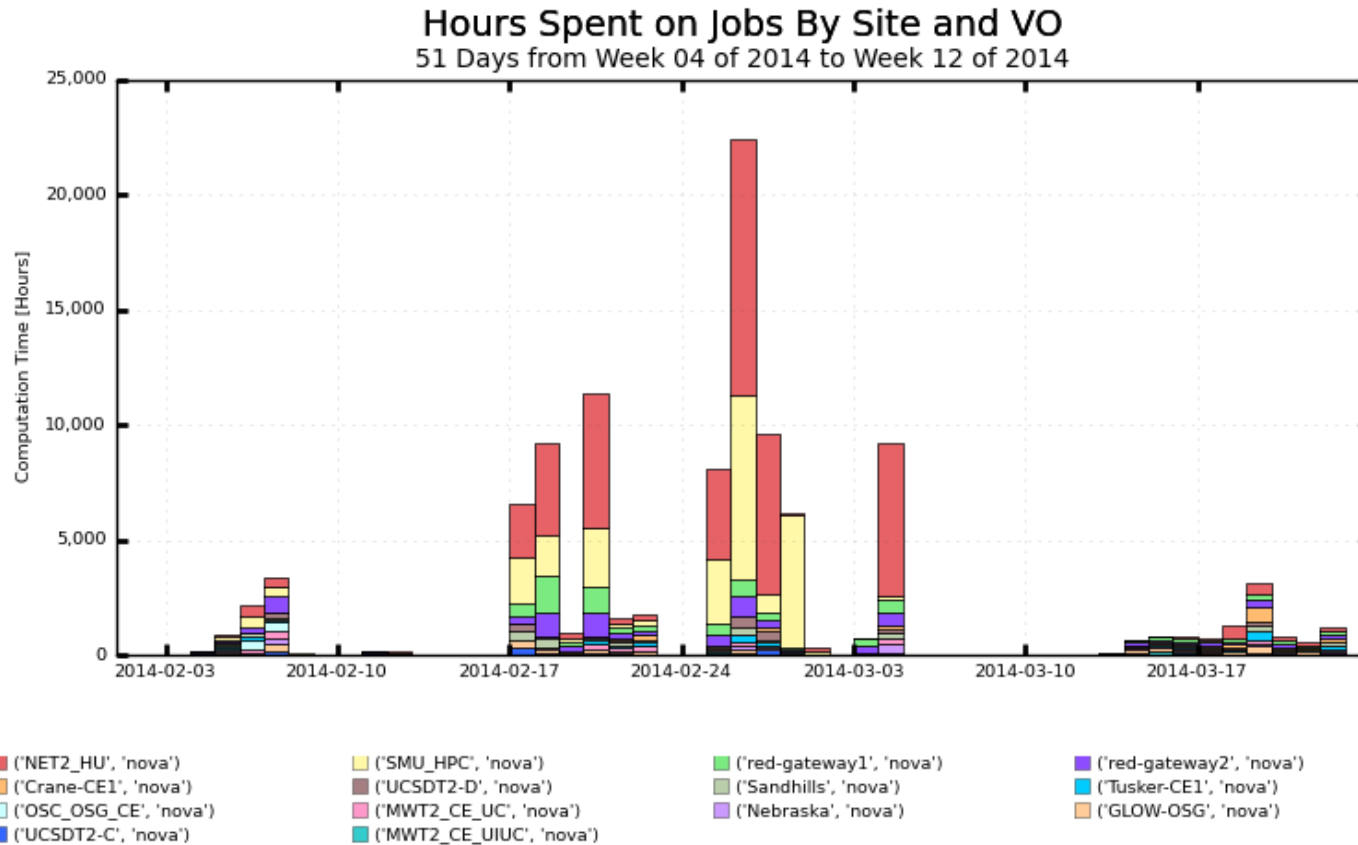
§ 1 – Integration success: early NOvA MonteCarlo on OSG

- Initial contained computing campaign with a single well-understood application
- 1,000,000 ev generated with ~10k jobs for 88,535 CPU h in 2 weeks of ops and 2 TB of data
- Run on OSG at SMU (dedicated), UNL, Uwisc, UC, UCSD and O(100) jobs at FermiCloud
- Operations consisted of 1 submission + 2 recoveries
- Spent about 10% more resources than expected due to preemption



- using same submission command from user perspective, launched jobs “on-demand” at FermiCloud to provide proof-of-principle for Cloud bursting grid jobs

§ 2 – Integration success: Ongoing NOvA computation on OSG ...



Maximum: 22,459 Hours, Minimum: 0.00 Hours, Average: 2,106 Hours, Current: 1,212 Hours

...and Amazon Web Services

§ 3 – Integration success: early experience of MicroBoone on OSG

- LArSoft: package common to multiple Intensity Frontier experiments – Liquid Argon Simulation of Time-Projection Chambers
- Code with dependencies were made fully portable and deployed through CVMFS.
- Demonstrated the ability to run interactively at OSG sites (Fermilab and SLAC)
- Successfully submitted jobs to OSG sites – Nebraska, UCSD, SLAC

On-demand Grid-bursting – NOvA MC on AWS (I)

Number of jobs 1088:

- Successes 1047
 - Relies on extensions to GlideinWMS
- Failures 41
 - File upload (ifdh) 1
 - File download (SAM) 1
 - nova executable 33
 - Art non-0 exit code 1
 - Geant4 failure 4
 - Hang job (100% cpu): killed

Wall
Time
Hours

Total hours: 1135 (56 min/job)

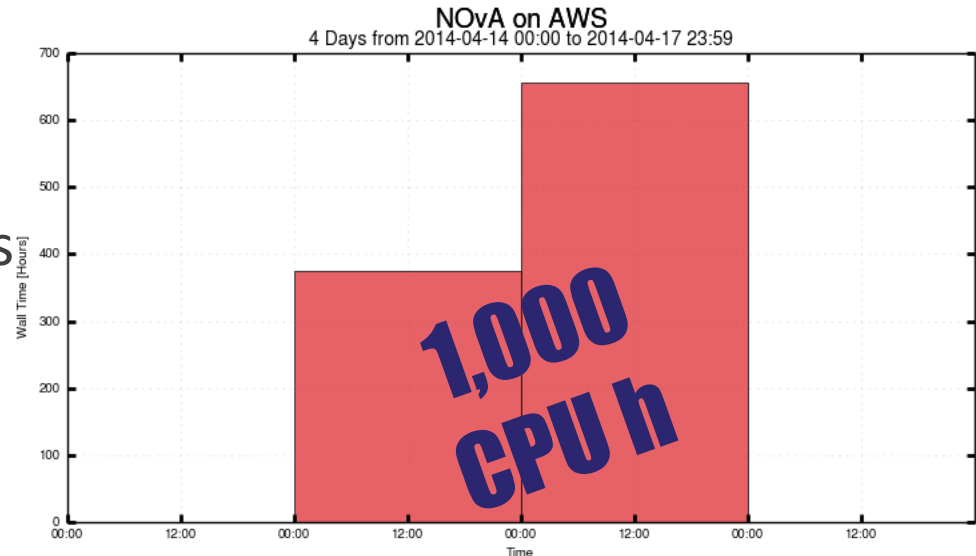
- overhead: 6 min/job

Output data: 326 (incl. 4GB tests)

- data: 152GB (BA via BeStMan)
- log files: 170 GB (on fifebatch1)

Input data: 72 GB (free)

Cosmics Background in the Near Detector



■ NOVA , Amazon_AWS

Maximum: 655.56 Hours, Minimum: 0.00 Hours, Average: 343.20 Hours, Current: 655.56 Hours

m1.medium spotpricing bid \$0.07 →
got “blended” price \$0.066 , instead: need to
investigate

Total cost: \$125

Data transfer: \$39

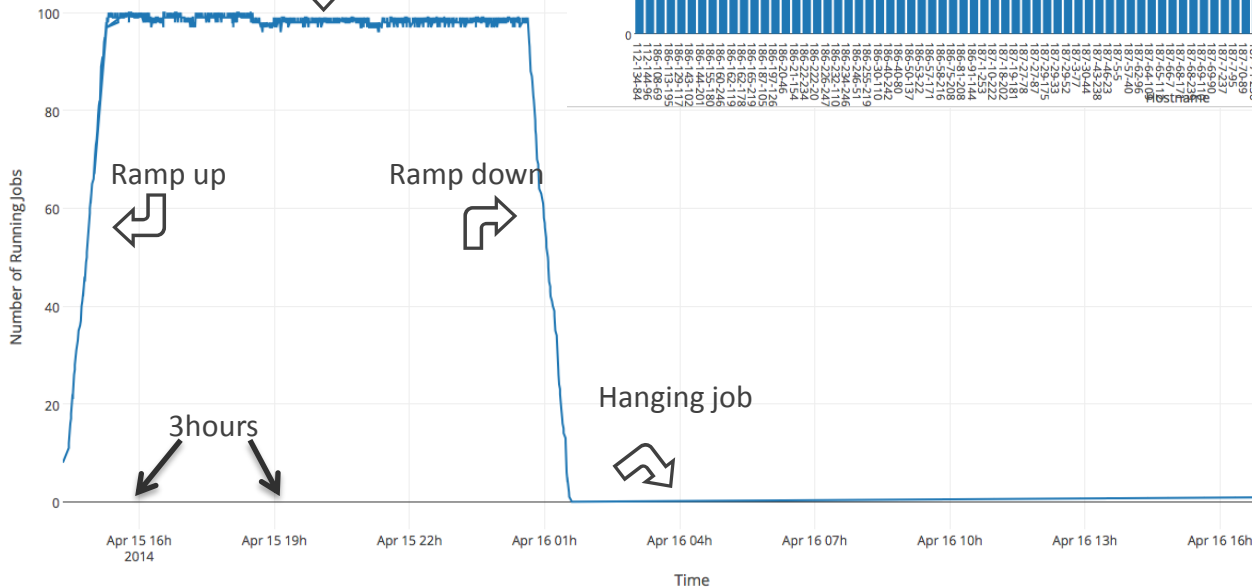
CPU: \$86

On-demand Grid-bursting – NOvA MC on AWS (II)

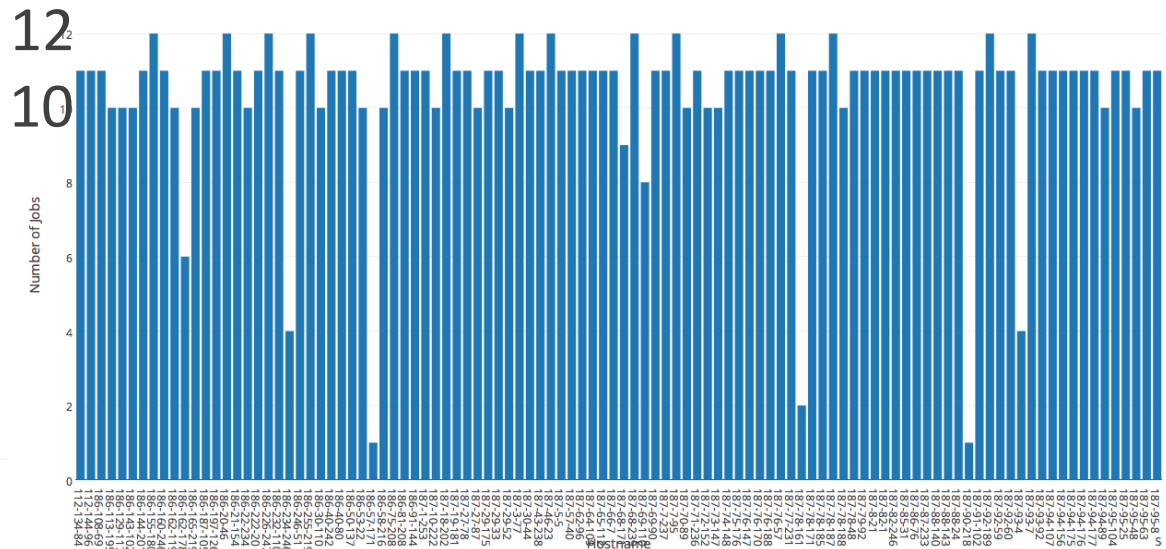
100 VM Limit:
FNAL Squid



100 VMs limit



Number of User Jobs per VM



AWS Issues and workarounds-> Best Practices

- Some DH software depends on “hostname” inside the VM matching the external IP
 - Worked around for now by changing our VM’s hostname command
- Can’t use active mode GridFTP (multiple stream)
 - Still investigating what to do here
- CVMFS requires autofs and fuse
 - We supply a pvops-compatible SLF6 kernel with our image and run from that.
 - Stored on EBS volume for now
- Don’t have Squid proxies on AWS, using a squid at FNAL
- Will revisit some of these workarounds this summer and make a best practices document
- Also working on a automated service to go from FermiCloud “raw” and “qcow2” images to AWS-compatible format and back.

Summary: Analysis on Cloud and Grid

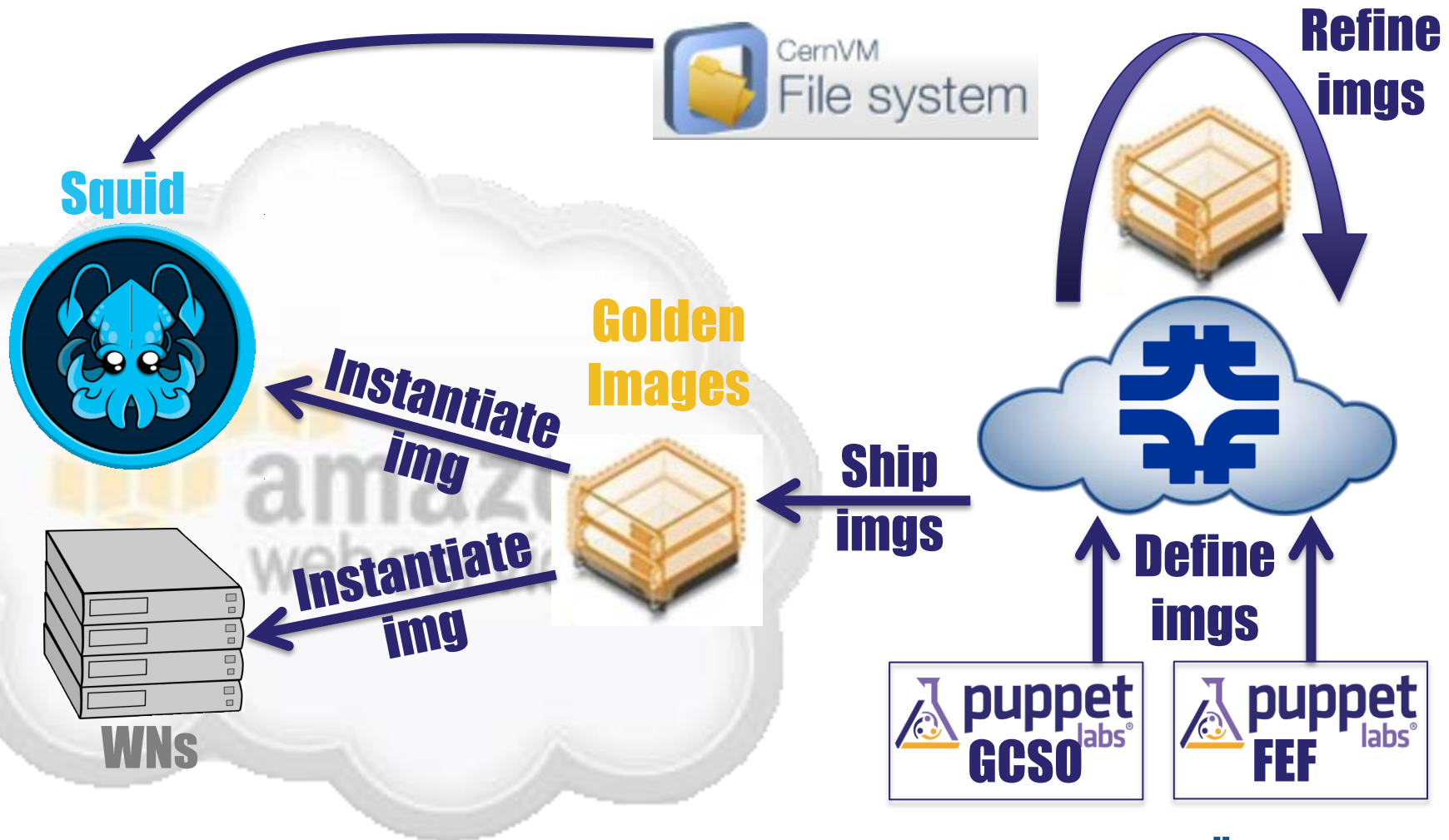
- Fermilab Physics program is extremely active at all three frontiers
- Successful demonstration of FIFE project architecture for small experiments to take advantage of distributed resources on grids and clouds
- Two more experiments (MicroBooNe and LBNE) in the process of on-boarding to running on OSG.
 - They share the LArSoft software package for simulation/analysis of liquid argon neutrino detectors
- Four more experiments (MicroBooNe, LBNE, MINERvA, Mu2e) are requesting cloud capacity for non-standard offline workflows.
- For certain sets of applications, AWS is a good alternative of similar reliability, performance, and ease of use to running opportunistically on grid.

Present and Near Future Cloud Services Work

- Virtual Infrastructure and Provisioning:
 - 1000-VM workflows in production
 - Provisioning algorithms for commercial and private clouds
 - Add more features to Idle VM Detection
 - Deployment of complicated groups of virtual machines on cloud
 - Provisioning of unconventional “worker node” virtual machines
- Interoperability
 - AuthN/AuthZ—future of X.509?
 - VM Image Conversion Service (raw, qcow2 <-> AWS)
 - More commercial and community clouds (Google, Rackspace)
 - Object-based storage R+D (Gluster, Ceph, Swift)

On-demand provisioning of complex services in the Cloud

- Use case: CVMFS / Squid + 1000s WN for NOvA



Projects: VM Batch Submission Facility

- Have done demo of up to 100 simultaneous worker node VM's on FermiCloud, ran NOvA workflows, soon Mu2e.
- Plan to scale up to 1000 over summer (using old worker nodes to expand the cloud). Also burst to AWS and KISTI GCLOUD
- Problem:
 - VM's take lots of IP addresses, we don't have that many.
 - So simulate private net structure of the public cloud
 - Start with Network Address Translation
 - Consider private (Software-defined) net that is routable only on-site
 - This could be key technology if you also need access to storage.
 - Very good way to find out if there is anything depending on NFS access, direct database access, etc.

Acknowledgments

- Significant work done by several NOvA collaborators: Gavin S. Davies, Nate Mayer, Andrew Norman, Eric Flumerfelt, Raphael Schroeter
- Help from the GlideinWMS project at Fermilab to debug remote site problems and add AWS spot-pricing (B. Holzman, P. Mhashilkar) and the OSG glidein factory operators at UCSD led by Igor Sfiligoi
- Help from site admins at Nebraska, UCSD, MWT2, OSC, Harvard, Wisconsin, SMU.
- OSG Operations and the OASIS project
- Collaborators at IIT and KISTI

Extra slides

On-demand deployment of VM with custom-specs

Experiment	Request	Short Description	Status
Minerva / FIFE	RITM0095773	Test batch job on std WN w/ fast turn-around	Planning
Mu2e	RITM0092720	Batch jobs with 8 GB RAM	Testing
CDF	RITM0096667	Interactive 32-cores custom system	Closed
DES	—	2 VM on 2 full nodes for large-memory large-disk CPU-intensive app (co-add)	Req. in prep.
NOvA	—	Grid bursting w/ std WN – current scalability on FC at 150VM	Req. TBD
μBooNE	—	WN with custom specs	Req. TBD
LBNE	—	WN with custom specs	Req. TBD

FermiCloud Project Accomplishments to Date

- Identified OpenNebula Technology for Infrastructure-as-a-Service,
- Deployed a **production** quality Infrastructure-as-a-Service facility,
- Made technical breakthroughs in authentication, authorization, fabric deployment, accounting and high availability cloud infrastructure,
- Continues to evolve and develop cloud computing techniques for scientific workflows,
- Supporting both **development/integration** and **production** user communities on a common infrastructure.
- Project started in 2010.
- OpenNebula 2.0 cloud available to users since fall 2010.
- OpenNebula 3.2 cloud up since June 2012
- **FOCUS on using these technologies to deliver On-Demand Services to scientific users**

Launch of Complex Workflow

- Coordinated Launch of complex workflow (compute + caching + storage)
- Launch a set of virtual machines
 - Let them know where all the others are
 - What web caching server, what database server, etc.
- Investigating using this for data preservation in batch systems
- Investigate benefits of far-end temporary storage element:
 - Is it better to have 400 machines waiting to send back their I/O (and charging you while they run)
 - Or should you queue it all on a few machines (incurring S3 storage charges)? How many is a few?
 - Current strategy relies on locking mechanism, one of many processes takes out a lock and attempts to send back all data queued at the time.

Winston the “Cloud Dog” is all grown up



In 2009 the cloud was an eager puppy ready to fetch and help



Now a full-grown Champion and ready to be put to work!