# VC3

## Virtual Clusters
## for Community Computation

Douglas Thain, University of Notre Dame
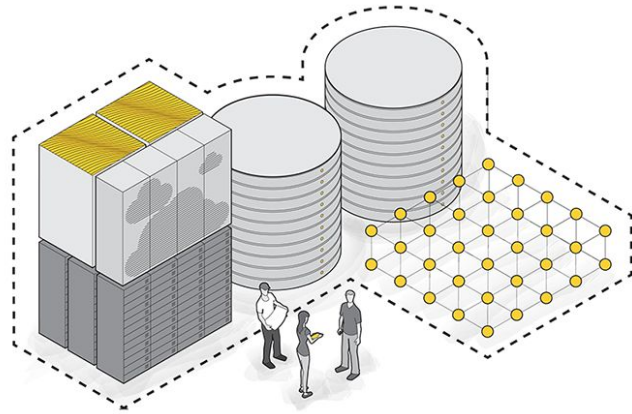Rob Gardner, University of Chicago
John Hover, Brookhaven National Lab

# Outline

- Intro to the Virtual Cluster Concept
- Self-service Provisioning of Middleware Across Heterogeneous Resources
- Dealing with Software Environments
- Clustering Middleware Examples
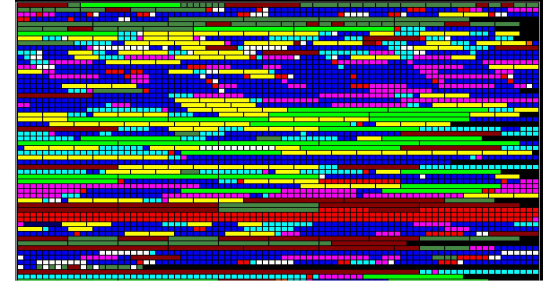- New Applications and Configurations
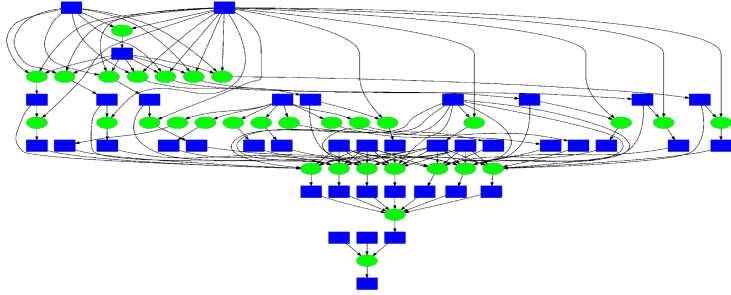- Thoughts and Lessons Learned

# Introduction

VC3: A platform for provisioning cluster frameworks over heterogeneous resources for collaborative science teams

You have developed a complex workload which runs successfully at one site, perhaps your home university.
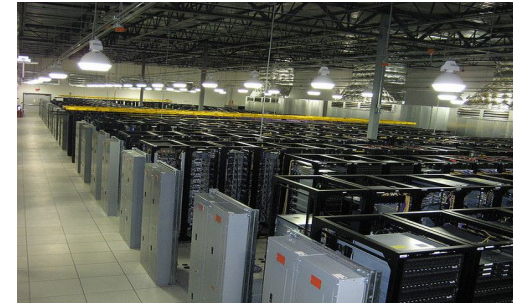


Now, you want to migrate and expand that application to national-scale infrastructure.
And allow others to easily access and run similar workloads.



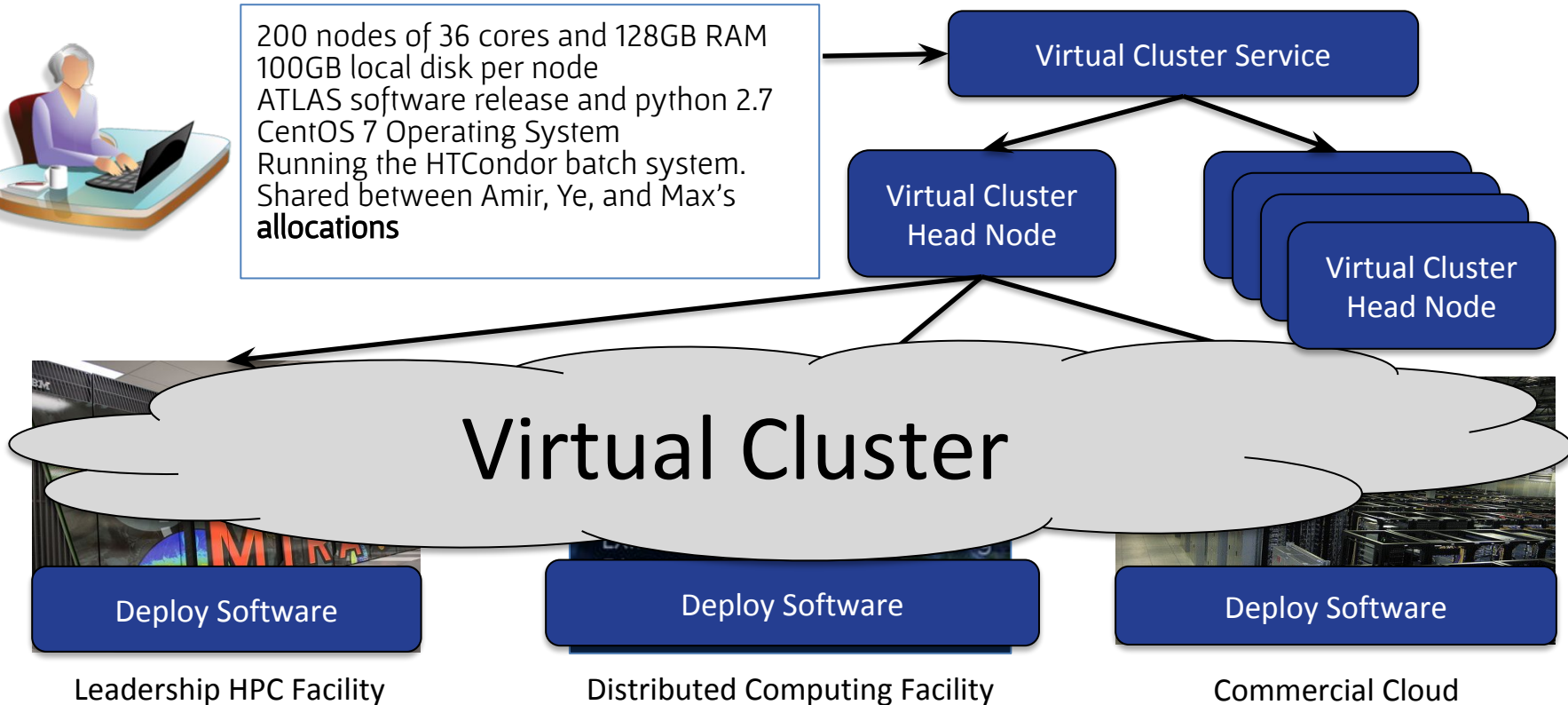Leadership HPC Facility          Distributed Computing Facility          Commercial Cloud
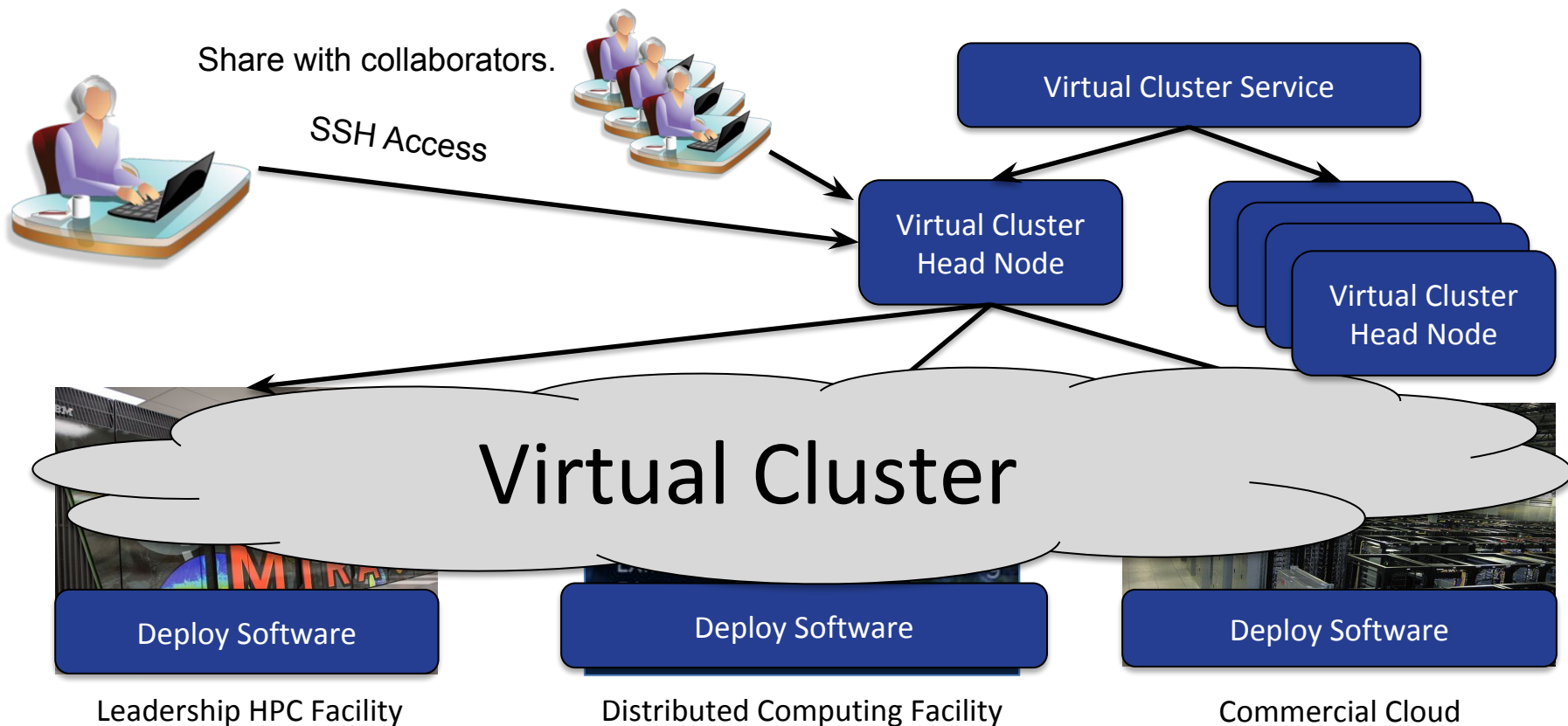
# Concept: Virtual Cluster

200 nodes of 36 cores and 128GB RAM
100GB local disk per node
ATLAS software release and python 2.7
CentOS 7 Operating System
Running the HTCondor batch system.
Shared between Amir, Ye, and Max's
**allocations**

Virtual Cluster Service

Virtual Cluster Head Node

Virtual Cluster Head Node

## Virtual Cluster

Deploy Software

Deploy Software

Deploy Software

Leadership HPC Facility

Distributed Computing Facility

Commercial Cloud

# Concept: Virtual Cluster



Share with collaborators.

SSH Access

Virtual Cluster Service

Virtual Cluster Head Node

Virtual Cluster Head Node

Virtual Cluster

Deploy Software

Deploy Software

Deploy Software

Leadership HPC Facility

Distributed Computing Facility

Commercial Cloud

# VC3: Virtual Clusters for Community Computation

- **VC3 is an interactive service for creating/sharing/using virtual clusters.**
- A virtual cluster consists of:
  - 1 x **head node** for interactive access to the cluster.  (SSH, Jupyter,..)
  - N x **worker nodes** for executing your workload.
  - **Middleware** to manage the cluster.  (HTCondor, Makeflow, Spark, ...)
  - **Application** software to do real work.  (BLAST, CMSSW, etc...)
- A virtual cluster is created using:
  - Your standard accounts/credentials on existing facilities.
  - Plain ssh/qsub/srun access on each facility.
  - Container technology (if available) or user-level software builds (otherwise).
  - (No special privileges or admin access required on the facility.)

# Self-Service Provisioning

# www.virtualclusters.org



login with a vc3 account

# Institutional Identity (CI-Logon)

# Curated Resources

| Name | Organization | Description | Cores | Memory | Storage | Native OS | Features |
|------|-------------|-------------|-------|--------|---------|-----------|----------|
| **Cori** | National Energy Research Scientific Computing Center (NERSC) | Cori Supercomputer at NERSC | 32 | 4000 MB | 10000 MB | suse:v12 | Shifter |
| **MWT2** | Midwest Tier 2 | ATLAS Midwest Tier 2 Center job gateway (UChicago) | 4 | 1000 MB | 1000 MB | scientificlinux:v6.9 | N/A |
| **Midway** | University of Chicago Research Computing Center (RCC) | Midway cluster at the University of Chicago Research Computing Center (RCC) | 64 | 4000 MB | 10000 MB | scientificlinux:v6.7 | N/A |
| **Stampede 2** | Texas Advanced Computing Center (TACC) | Stampede 2 Super Computer | 96 | 2000 MB | 10000 MB | centos:v7.4 | Singularity |
| **CoreOS** | University of Chicago | CoreOS Cluster | 4 | 1000 MB | 1000 MB | scientificlinux:v6.9 | Singularity |
| **UCT3** | University of Chicago | UChicago ATLAS Tier 3 | 4 | 1000 MB | 1000 MB | scientificlinux:v6.9 | N/A |
| **ND CCL** | University of Notre Dame Cooperative Computing Lab | ND-CCL login none | 4 | 1000 MB | 10000 MB | redhat:v7 | Singularity |
| **Bridges** | Pittsburgh Supercomputing Center | Bridges Supercomputer at PSC | 28 | 4000 MB | 35000 MB | centos:v7.3 | Singularity |
| **VC3 Test Pool** | VC3 | VC3 Test Pool | 4 | 1000 MB | 1000 MB | centos:v6.9 | N/A |
| **UCLA Hoffman2** | University of California, Los Angeles | UCLA Hoffman2 | 8 | 1000 MB | 10000 MB | centos:v6.9 | N/A |
| **OSG Connect** | Open Science Grid | Open Science Grid (SL7) | 4 | 1000 MB | 1000 MB | Unknown | N/A |

Resource Profiles

▼ Filter

# Allocations

## Step 1: Log Into Resource

In a terminal, type:

```
ssh btovar@cclvm05.crc.nd.edu
```

## Step 2: Access Resource

Enter your password for `cclvm05.crc.nd.edu` for access

## Step 3: Add Allocation SSH Public Key to Resource

Once the SSH key is generated below, click `Copy to Clipboard` and paste the following line into your SSH session. You will only need to do this once per allocation.

3-IooIwWcJuHKdBKwMtCKgLawLRVX/an5jPwIowcDqgf5c13504NOqg7egk
/GTjhj8YrCyX6UhqG+S3nOxOf+ewxx3RSIMf9lsFZpDNdXwJi1YD1dyRCYy8TwNhBg9GlkCxEKMqfOgo
L6ROpicuUhFY6yT9apKG0x1mPSM
/94ETHxIkBmNK8Ph926fuT+F+QQToSQV0v9o9hWLGiGdNoztW8OUkSFzZ6uZE5zfFpP0xq45a4+FYE
TorljRapgPsjmSjmSB7TeD+qs1ECilwrrg3iJP0RBOEMMeLf7rwgIXXjxtZkBUQ72lkq5ltXTUAyeu0CbGglI
Q7ZHGHrNTyKkSPLl7rXEi7nnz6ofgUJCU3L7hr2VKKy84RcHPSfep64qV3iIOCw1o6SPvu6IwRYeqhfe
A0o
/yKp1IvapyfM7Ptuy+6yWZ7grZlb9AtBolcoBC0lpi964MR8T4D8RKp1960nCG5ltXwC4mmPSgffQ0fOl
WJom7TudG+yTWouWikup0ieObZX5w8SKFc0H

[Copy to Clipboard]

## Step 4: Validate Allocation

# Projects

| Project Profiles | | | | ▼ Filter |
|---|---|---|---|---|
| **Name** | **Members** | **Allocations** | **Description** | |
| **vc3-team** | **Benjamin Tovar (Owner) - btovar@nd.edu**<br>Lincoln Bryant (UChicago)<br>Jeremy Van (UChicago)<br>Robert Gardner (UChicago)<br>Kenyi Hurtado (University of Notre Dame) | btovar-ndccl<br>khurtado-osgconnect<br>lincolnb-midway | Currently no description | |
| **btovar** | **Benjamin Tovar (Owner) - btovar@nd.edu**<br>Benjamin Tovar (University of Notre Dame) | btovar-ndccl | Currently no description | |

# Launching a Virtual Cluster

VIRTUAL CLUSTER NAME

my-virtual-cluster

CLUSTER TEMPLATE *

lincolnb-htcondor-10-workers

shared cluster definition

ENVIRONMENT

btovar-oasis-osg

workers will have this environment installed

ALLOCATIONS*

Nothing selected

Select Allocations for Virtual Cluster ✕

Select All        Deselect All

btovar-ndccl

khurtado-osgconnect

lincolnb-midway

allocations available in this project

# Cluster Status

| | My Virtual Clusters | | | ▼ Filter |
|---|---|---|---|---|
| **Name** | **State** | **Cluster Template** | **Workers** | **Head Node** |
| **my-virtual-cluster** | Running<br>All requested compute workers are running. | lincolnb-htcondor-10-workers | Requested: 10<br>Running: 7<br>Queued: 3<br>Error: 0 | 128.135.158.187 |

# Workers from many sites

```
[btovar@btovar-my-virtual-cluster ~]$ ip addr | grep 128.135.158.187
    inet 128.135.158.187/25 brd 128.135.158.255 scope global dynamic eth0
[btovar@btovar-my-virtual-cluster ~]$ condor_status
Name                                          OpSys      Arch    State      Activity LoadAv Mem   ActvtyTime

slot1@glidein_21791@camd01.crc.nd.edu  Notre  LINUX      X86_64  Unclaimed  Idle       5.120  4013  0+00:19:37
slot1@glidein_29106@camd01.crc.nd.edu  Dame   LINUX      X86_64  Unclaimed  Idle       5.120  4013  0+00:19:37
slot1@glidein_91802@camd05.crc.nd.edu         LINUX      X86_64  Unclaimed  Idle       5.260  4013  0+00:19:37
slot1@glidein_39133@iut2-c257.iu.edu   OSG    LINUX      X86_64  Unclaimed  Idle      34.620  3223  0+00:19:48
slot1@glidein_61297@lnxfarm275.colorado.edu   LINUX      X86_64  Unclaimed  Idle       6.990  3002  0+00:14:36
slot1@glidein_28373@midway091.rcc.local       LINUX      X86_64  Unclaimed  Idle       8.170  2013  0+00:19:36
slot1@glidein_71179@midway098.rcc.local       LINUX      X86_64  Unclaimed  Idle       7.480  2013  0+00:19:36
slot1@glidein_46364@midway260.rcc.local UChicago LINUX   X86_64  Unclaimed  Idle       8.040  2013  0+00:19:35
slot1@glidein_39282@midway324.rcc.local       LINUX      X86_64  Unclaimed  Idle       8.750  2013  0+00:19:36
slot1@glidein_39133@uct2-c373.mwt2.org        LINUX      X86_64  Unclaimed  Idle      34.080  2415  0+00:19:33

                    Machines Owner Claimed Unclaimed Matched Preempting  Drain

    X86_64/LINUX        10      0      0       10        0         0        0


         Total          10      0      0       10        0         0        0
[btovar@btovar-my-virtual-cluster ~]$
```
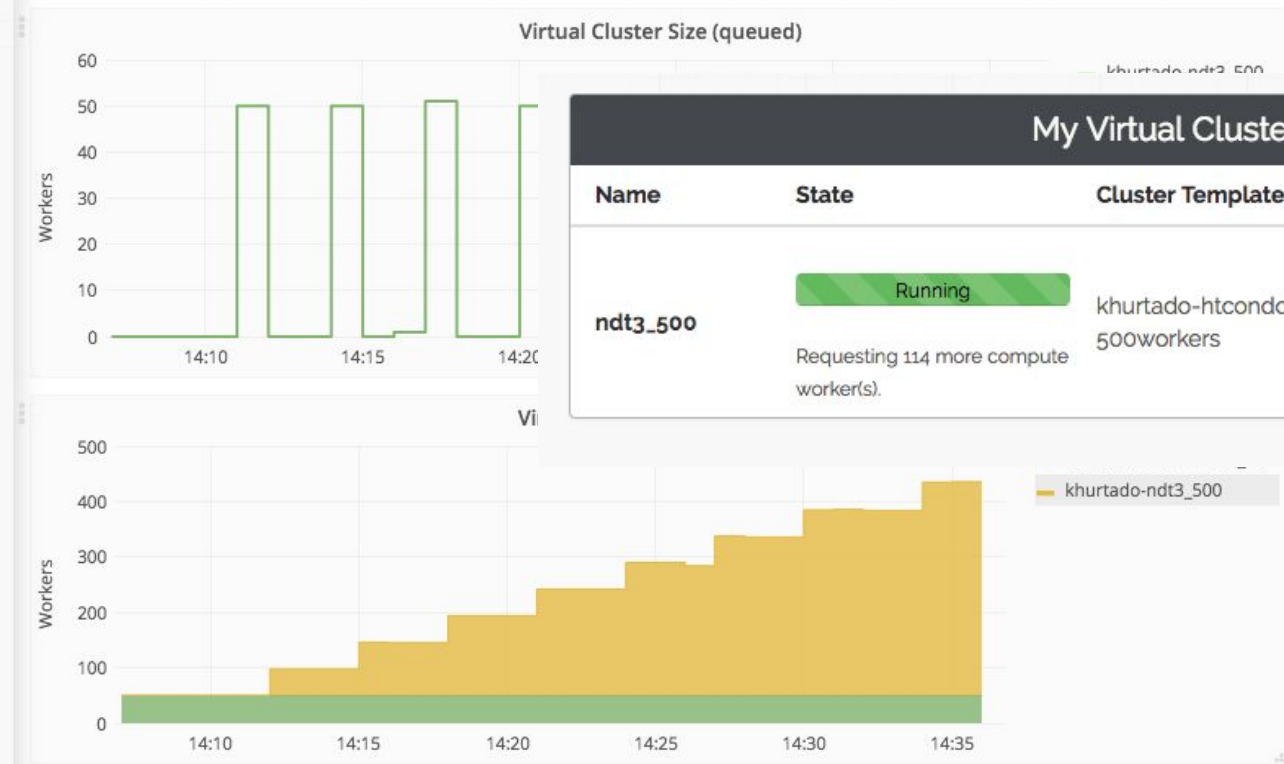
# Dealing with Software Environments

# Deploying Software Environments

The **vc3-builder**, a command-line tool for deploying software environments on clusters.

```
vc3-builder
    --require-os centos:7
    --mount /scratch=/data
    --require /cvmfs
    --require python:2.7 -- myapp ...my args...
```

https://github.com/vc3-project/vc3-builder

# Basic Use Case: Application Software

```
vc3-builder --require ncbi-blast
..Plan:    ncbi-blast => [, ]
..Try:     ncbi-blast => v2.2.28
....Plan:    perl => [v5.008, ]
....Try:     perl => v5.
....could not add any
....Try:     perl => v5.
....could not add any
....Try:     perl => v5.
......Plan:    perl-vc3-
......Try:     perl-vc3-
......Success: perl-vc
....Success: perl v5.2
....Plan:    python => [v2.006, ]
```

**(New Shell with Desired Environment)**

```
bash$  which blastx
/tmp/test/vc3-root/x86_64/redhat6/ncbi-blast/v2.2.28/bin
/blastx

bash$ blastx –help
USAGE
  blastx [-h] [-help] [-import_search_strategy filename]
```
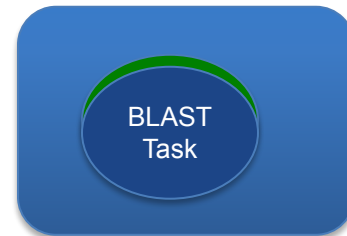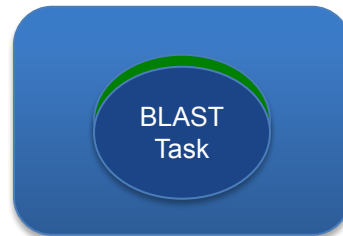
# Problem: Long Build on Head Node
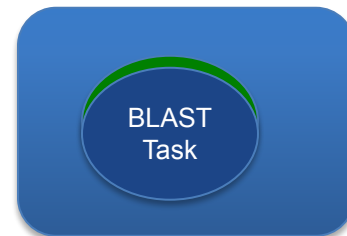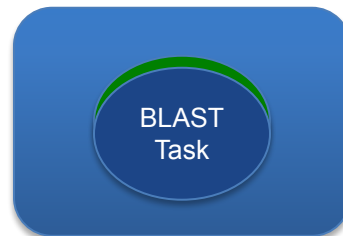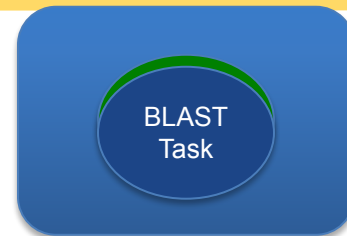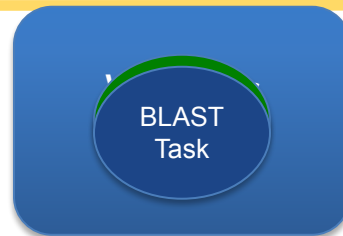
- Many facilities limit the amount of work that can be done on the head node, so as to maintain quality of service for everyone.

- Solution: Move the build jobs out to the cluster nodes. (Which may not have network connections.)

- Idea: Reduce the problem to something we already know how to do: *Workflow!* But how do we bootstrap the workflow software? *With the builder!*

# Bootstrapping Workflows and Apps

# Example Applications

Octave

MAKER



Typically, 2-3x faster overall.  But more importantly, filesystem-intensive jobs run on the cluster resources, not on the head node!

# Controlling Cluster Size

# Dynamic Cluster Resizing

VC3 Supports several mechanisms for setting the number of workers in a given virtual cluster.

- Set static size for cluster upon creation.
- Manually change cluster size via web interface. The provisioning factory automatically adds or removes worker jobs at resources.
- (For HTCondor cluster) Scale the cluster dynamically based on idle jobs on the head node.
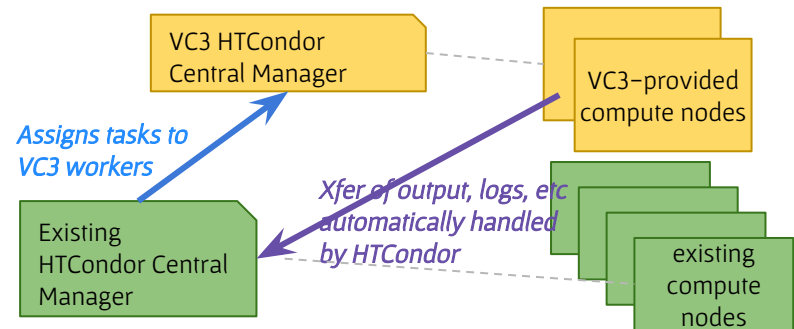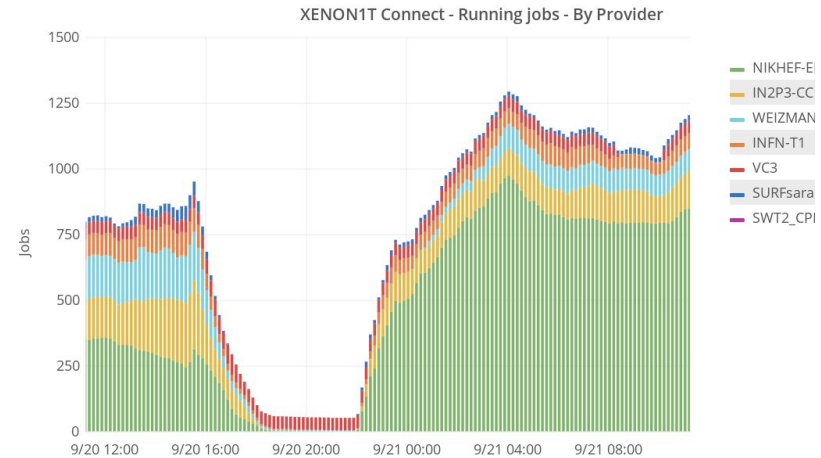
# Current Resizing Cases

| | | Job Removal Setting | |
| --- | --- | --- | --- |
| | | Peaceful = True | Peaceful = False |
| Source of desired number of workers | User (via web interface) | **Increasing:** Works as expected. **Decreasing**: Workers will terminate when they have finished current job. By default, the younger workers go away first. | **Decreasing**: Running jobs may be killed. But cluster can be reduced rapidly. |
| | VC3 (automatic, dynamic, based on user's job pressure on head node.) | **Increasing:** Works as expected. **Decreasing**: Workers will terminate when they have finished current job. By default, the younger workers go away first. | |

# Scaling Out Production Clusters

- Collaborations who have existing HTCondor pools can extend them by adding more worker nodes via VC3
- Add XSEDE resources, Open Science Grid, and campus HPC clusters
- End-users can transparently use additional resources



XENON1T Connect - Running jobs - By Provider

NIKHEF-E
IN2P3-CC
WEIZMAN
INFN-T1
VC3
SURFsara
SWT2_CPI

VC3 HTCondor Central Manager

VC3-provided compute nodes

*Assigns tasks to VC3 workers*

*Xfer of output, logs, etc automatically handled by HTCondor*

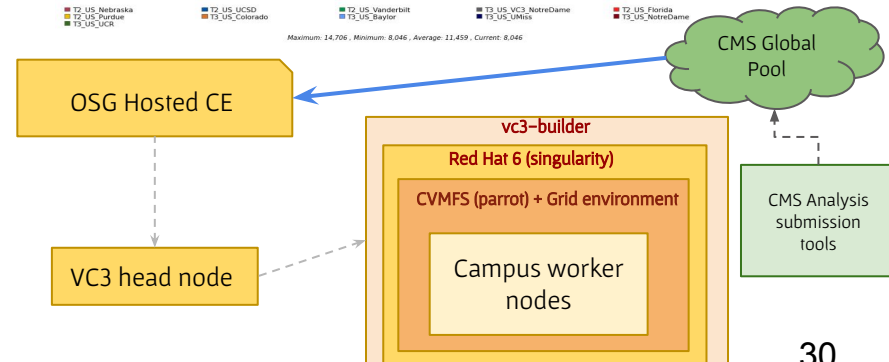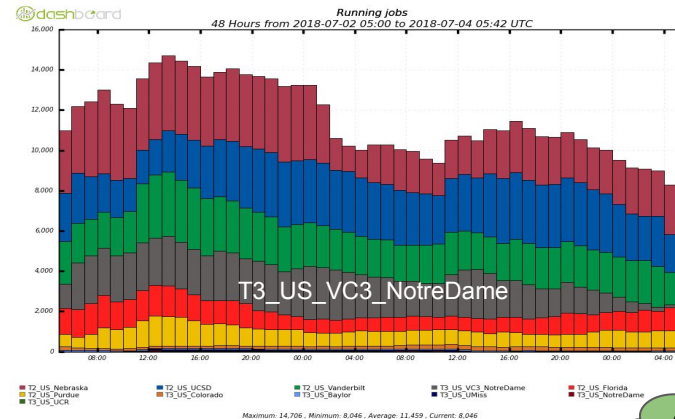Existing HTCondor Central Manager

existing compute nodes

29

# Building Tier 3s on top of campus resources

- Although different computing resources are often available at universities, meeting all requirements to deploy a valid Tier 3 able to run CMS workflows on the grid is not an easy task to achieve without the intervention of a system administrator with root access.
- VC3 allows the provisioning at user-level of:
  - The CERN File System (CVMFS) (via parrot)
  - The OSG grid environment on the worker nodes (via CVMFS)
  - Customized Operating Systems (via singularity)
- The OSG Compute Element (CE) is then integrated with the VC3 submit host, allowing the creation of a CMS Tier 3 using Notre Dame opportunistic campus resources without any root access level.

CMS Analysis activity per Tier Site

*Note T3_US_VC3_NotreDame performing at the scale of the Tier 2s!*
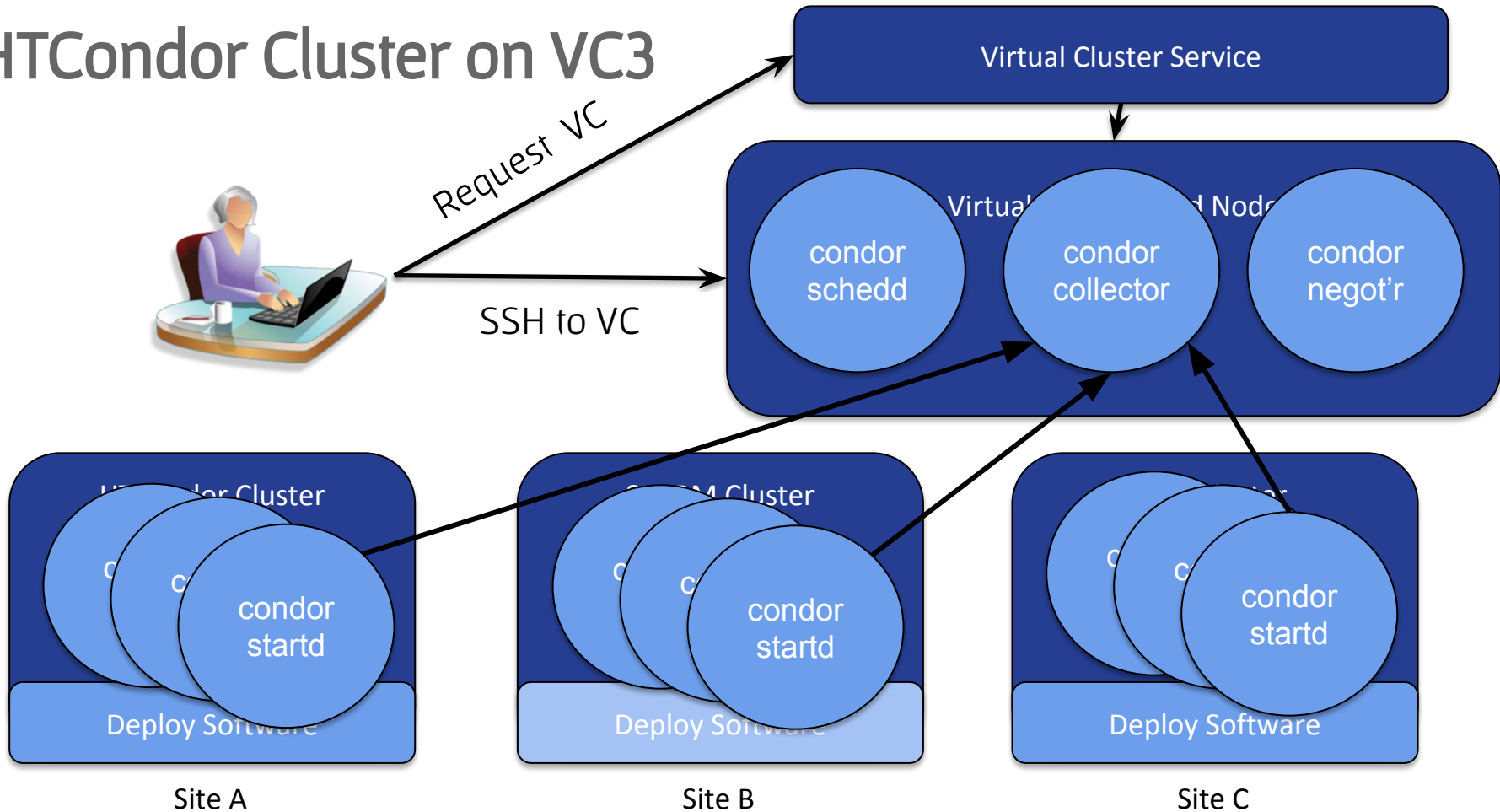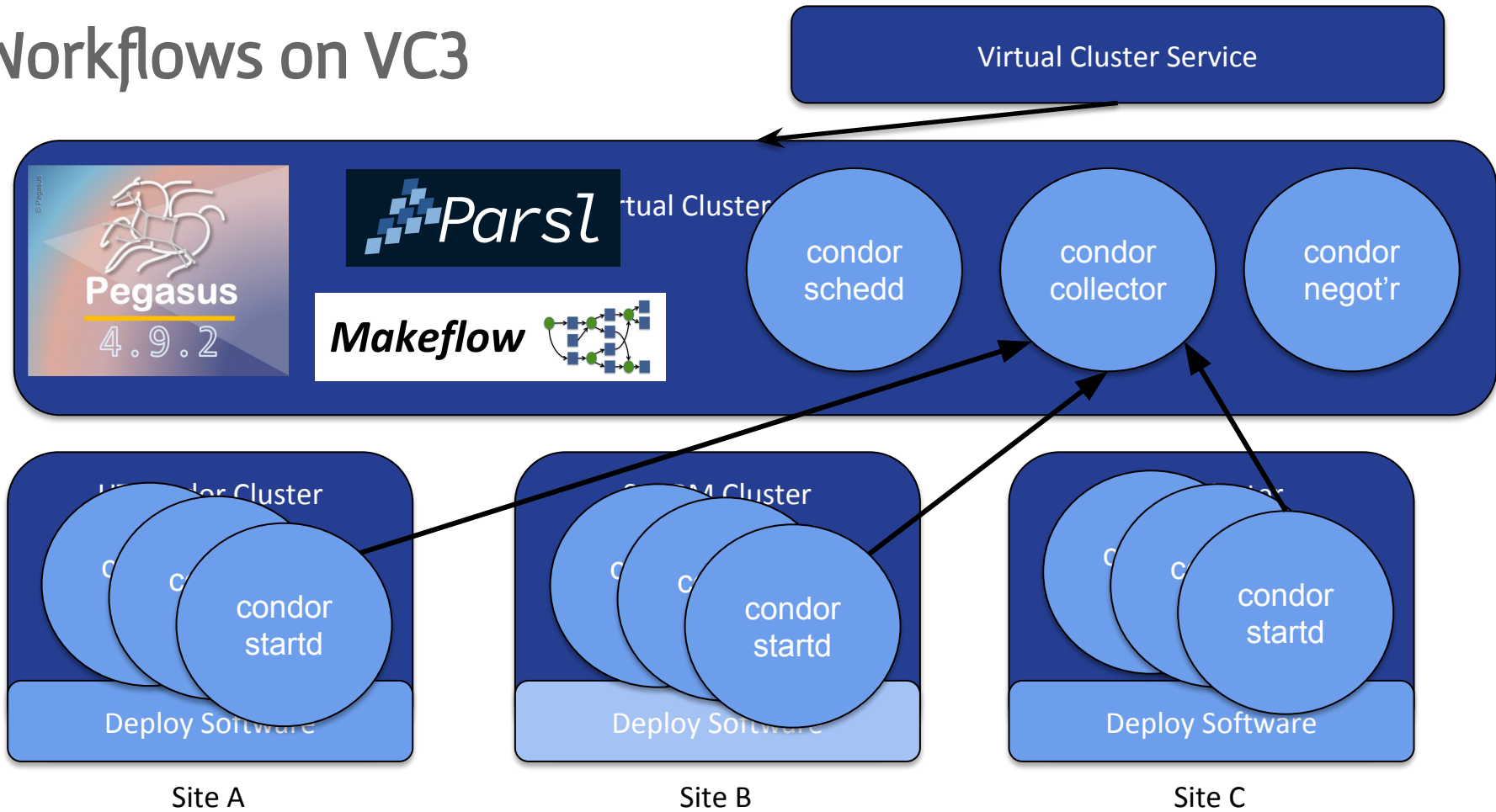
# Examples of clustering middleware

# Use Cases

- **HTCondor Pool**
  - Already demonstrated self−service, fully automated provisioning clusters, binding allocations into projects
- **Workflow Technologies**
  - Makeflow + Work Queue Pegasus + HTCondor,
  - Parsl + iPyParallel
- **Spark**
  - **D**eployment of a Spark cluster for data analysis on top of existing schedulers

- **JupyterLab**
  - **Popular interactive** analytics environment – provisionable by VC3
- **SCAILFIN + REANA**
  - Complex set of REANA services deployed on minikube on head node + HTCondor on Stampede, Blue Waters and PSC.
- **KOTO**
  - Helping a new collaboration with no establishes resources to run on BNL HPC, KEK HPC + OSG HTC + campus resources.
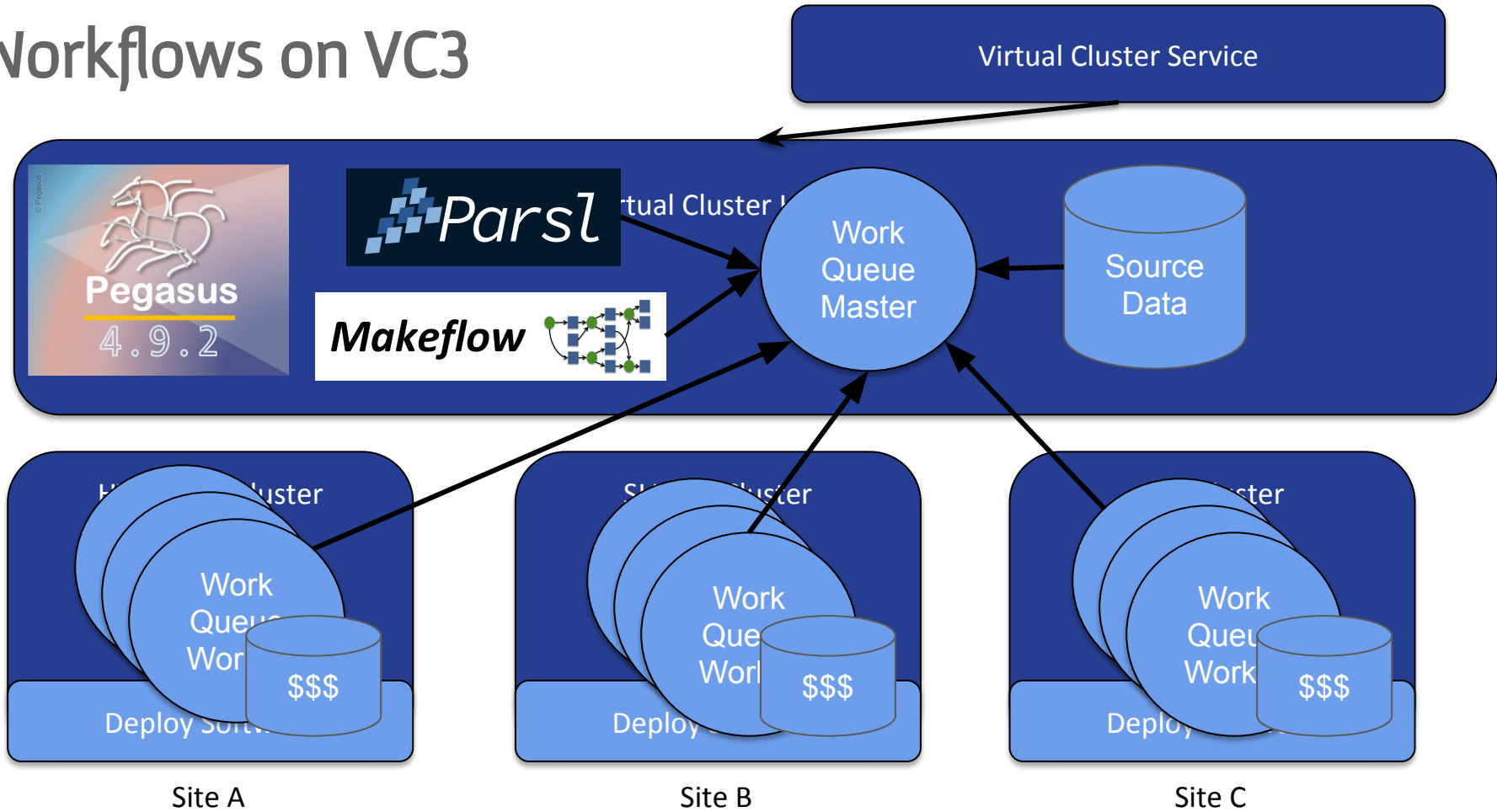
# HTCondor Cluster on VC3

# Workflows on VC3

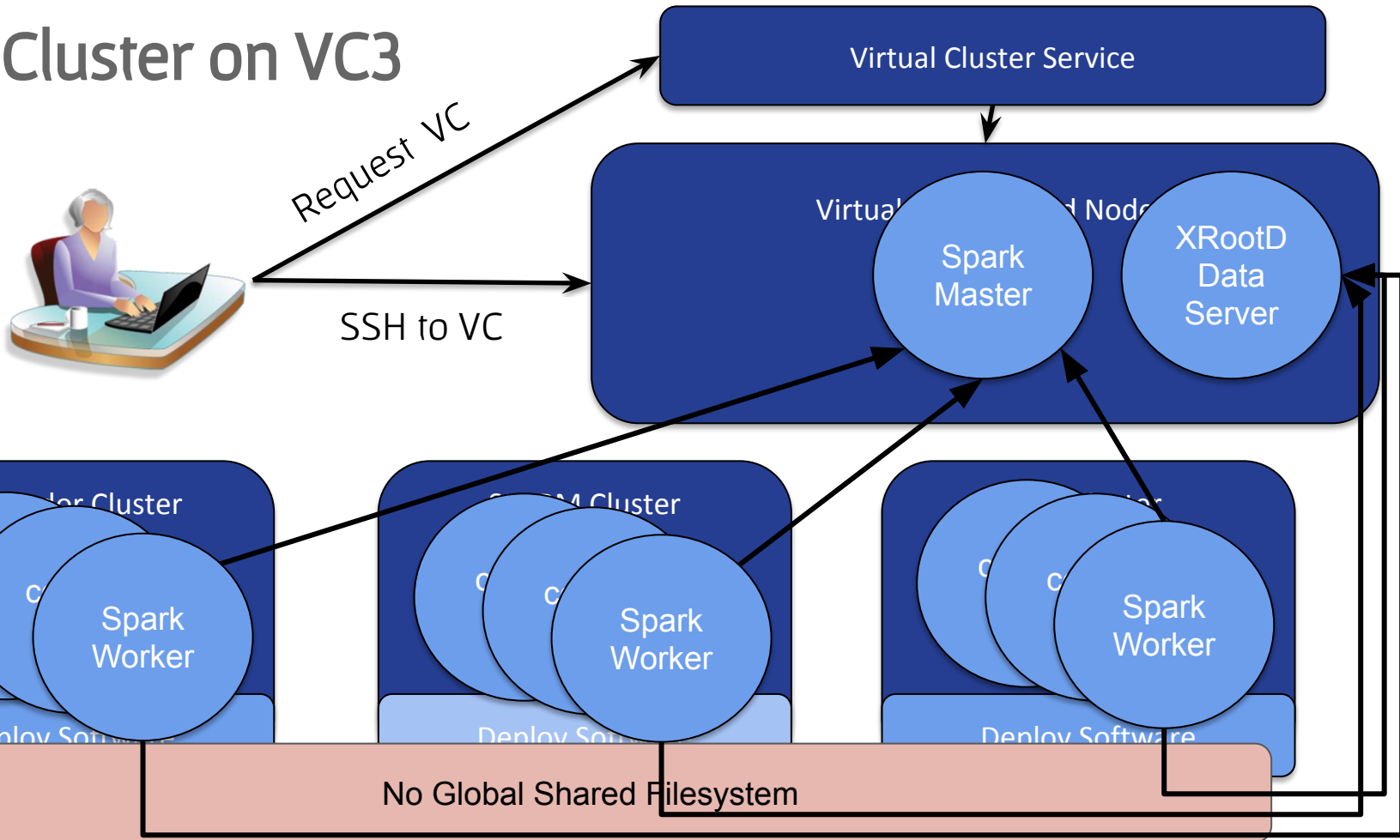# Workflows on VC3

# Spark Middleware Integration

- Apache Spark now a first-class supported middleware in VC3
- Spark master runs on the virtual cluster head-node
- Virtual cluster workers are spark slaves
  - Java JRE, spark, scala, pyspark are deployed as-needed by the target resource.
- A shared secret secures connection to the master from workers and application clients
- No shared filesystem, use of an IO driver is needed.

# Spark Cluster on VC3

Virtual Cluster Service

Request VC

SSH to VC

Virtual Head Node

Spark Master

XRootD Data Server

HTCondor Cluster

Spark Worker

SLURM Cluster

Spark Worker

Cluster

Spark Worker

Deploy Software

Deploy Software

Deploy Software

No Global Shared Filesystem

# Spark & CMS Analysis

- Late-stage custom analysis code processing a Mini-AOD file to produce a histogram.
- Application written in python, using pyspark
- Data IO with the spark-xrootd plugin
  - Automatically installed in headnode and workers
- Ran with 20 workers, 4 cores each.
- Maximum concurrency was 80 tasks.

# Spark & CMS Analysis

**Spark** 2.2.1    **Spark Master at spark://128.135.158.221:7077**

**URL:** spark://128.135.158.221:7077
**REST URL:** spark://128.135.158.221:6066 *(cluster mode)*
**Alive Workers:** 17
**Cores in use:** 68 Total, 68 Used
**Memory in use:** 895.3 GB Total, 17.0 GB Used
**Applications:** 1 Running, 0 Completed
**Drivers:** 0 Running, 0 Completed
**Status:** ALIVE

## Workers

| Worker Id |
| --- |
| worker-20180919123112-10.32.83.11-44930 |
| worker-20180919123157-10.32.79.15-40509 |
| worker-20180919123157-10.32.80.39-36998 |
| worker-20180919123209-10.32.79.65-42856 |
| worker-20180919123211-10.32.79.65-41106 |
| worker-20180919123220-10.32.79.45-44809 |
| worker-20180919123232-10.32.79.43-39239 |
| worker-20180919123233-10.32.80.47-45735 |
| worker-20180919123245-10.32.72.4-36789 |
| worker-20180919123250-10.32.77.80-44497 |
| worker-20180919123250-10.32.77.80-46283 |
| worker-20180919123307-10.32.76.62-41504 |
| worker-20180919123346-10.32.80.71-37519 |
| worker-20180919123414-10.32.76.64-35737 |
| worker-20180919123414-10.32.76.64-36614 |
| worker-20180919123456-10.32.76.74-40008 |
| worker-20180919123458-10.32.76.74-44468 |

## Running Applications

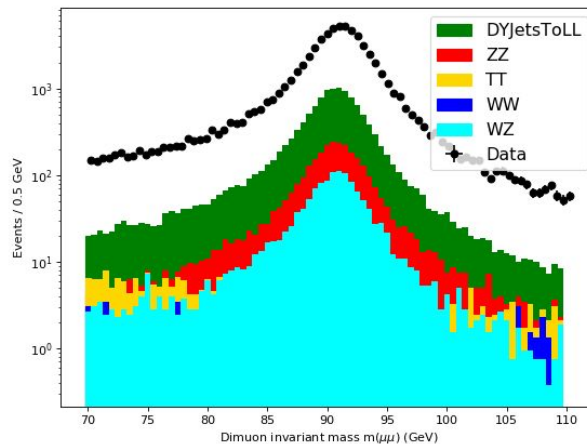| Application ID | Name | | Duration |
| --- | --- | --- | --- |
| app-20180919163039-0000 | (kill) Zpeak_Nanoaod-SF | | 4.5 min |

**Virtual Cluster: btovar-zpeak**

Terminate Cluster

STATE OF VIRTUAL CLUSTER

Running

All requested compute workers are running.

**Status**

BTOVAR-ZPEAK

Cluster Framework:

Requested 20
Running 20
Queued 0
Error 0

Resize Workers



Dimuon invariant mass $m(\mu\mu)$ (GeV); Events / 0.5 GeV; DYJetsToLL, ZZ, TT, WW, WZ, Data

# Provisioning JupyterLab Notebooks

- JupyterLab-based head nodes now launchable in VC3 development branch
- Web-based notebook interface for Python
- Uses Globus identity plugin to login
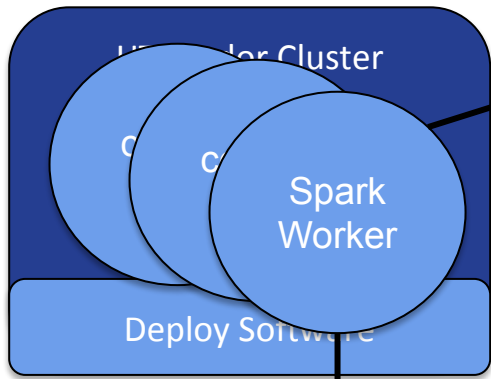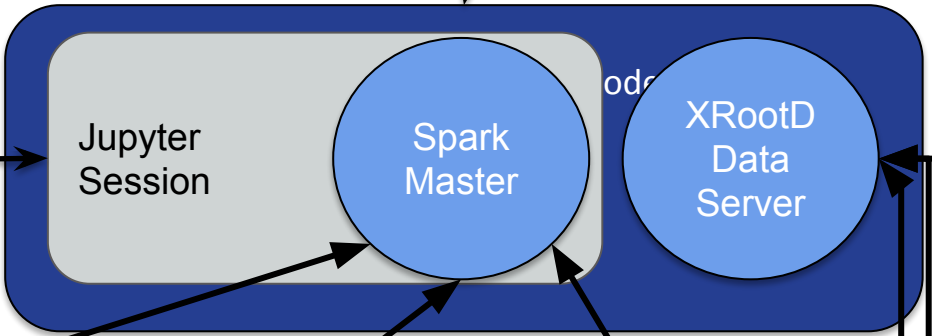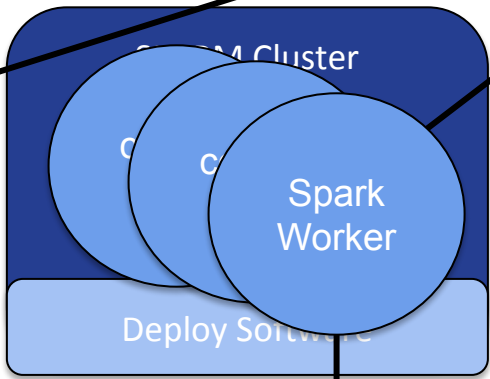- Integration with HTCondor and other middleware types

# Jupyter +Spark on VC3

Virtual Cluster Service

Request VC

WWW Access

Jupyter Session

Spark Master

XRootD Data Server

Spark Worker

Spark Worker

Spark Worker

Deploy Software

Deploy Software

Deploy Software

Site A

Site B

Site C

41

# Provisioning JupyterLab Notebooks

# Credential Delegation

Request VC

WWW Access

**Virtual Cluster Service**

Jupyter Session

Spark Master

XRootD Data Server

**Credential Delegation:**
- Globus Auth to VC3
- VC3 WWW to Jupyter
- VC3 WWW to Site SSH
- User to Jupyter
- Worker to Master
- Worker to XRootD

Each one a different technology!

SLURM Cluster

Spark Worker

Deploy Software

Spark Worker

Deploy Software

Site A

Site B

Site C

# Scalable cyberinfrastructure on HPC facilities

**SCAILFIN: Scalable CyberInfrastructure for Artificial Intelligence and Likelihood Free Inference**

The SCAILFIN project aims to deploy artificial intelligence and likelihood−free inference techniques and software using scalable cyberinfrastructure (CI) that is developed to be integrated into existing CI elements, such as the **REANA** system, to work on **HPC facilities**.



VC3 headnode



REANA is automatically deployed and integrated with HTCondor in VC3

REANA components are started via kubernetes(minikube)

# Simulation-Based Likelihood Free Inference

Symbolically:



Estimation of optimal estimator lends itself to ML methods:

- Training data derived from simulations
- Can be guided by optimal sampling based on phase space density of generator, sensitivity to physics under study

# reana reana-cluster - Simplified Diagram

(CERN Instance)

**REANA Cluster / Service Components**

| **DB /** *SQLAlchemy* TCP 5432 | **Workflow Mon** *ZMQ Proxy* TCP 8666/8667 | **Msg Broker** *RabbitMQ* TCP 5672, 15672 | **Workflow Engine & Controller** |

**CERN OpenStack COE + Kubernetes**

| **Job Controller** | **W.E Yadage+ Serial** | **WDB** TCP 1984, 19840 |

**REANA Cluster / Workers**

| **Yadage-run Docker -> Payload** | **Yadage-run Docker -> Payload** | **Yadage-run Docker -> Payload** |

# SCAILFIN on Blue Waters via VC3



VC3 Factory

Bosco (GSI-SSH)

**VC3 Headnode**

**REANA Components**

Kubernetes/Minikube

Reana-Job-Controller

Condor Schedd collector/CCB

Internet

**BW Submit Node**

Torque

reverse connection from condor startd to CCB/collector.

**MOM Node**

VC3-glidein

aprun -b -- shifter . . .

**Compute Node**

Run Shifter Payload for REANA workflow step
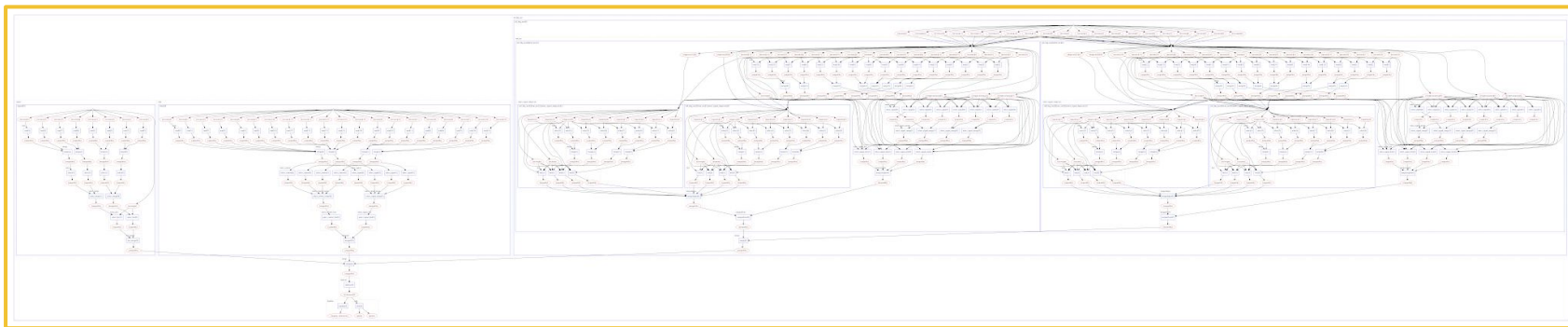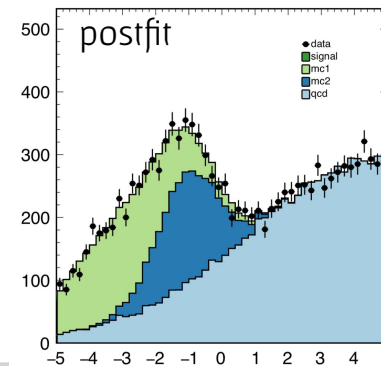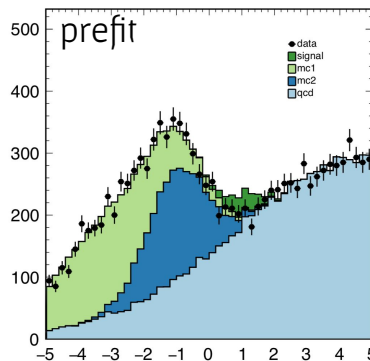
# SCAILFIN on Blue Waters via VC3

Complex Beyond Standard Model example (ran at BW via REANA + HTCONDO Rvia VC3)



Data is generated/emulated according to Standard Model expectations.

After processing, a statistical model involving both signal and control regions is built and the model is fitted against the observed data.

The signal sample is scaled down significantly to fit the data, which is expected since the data was emulated in accordance with a SM−only scenario

# KOTO experiment
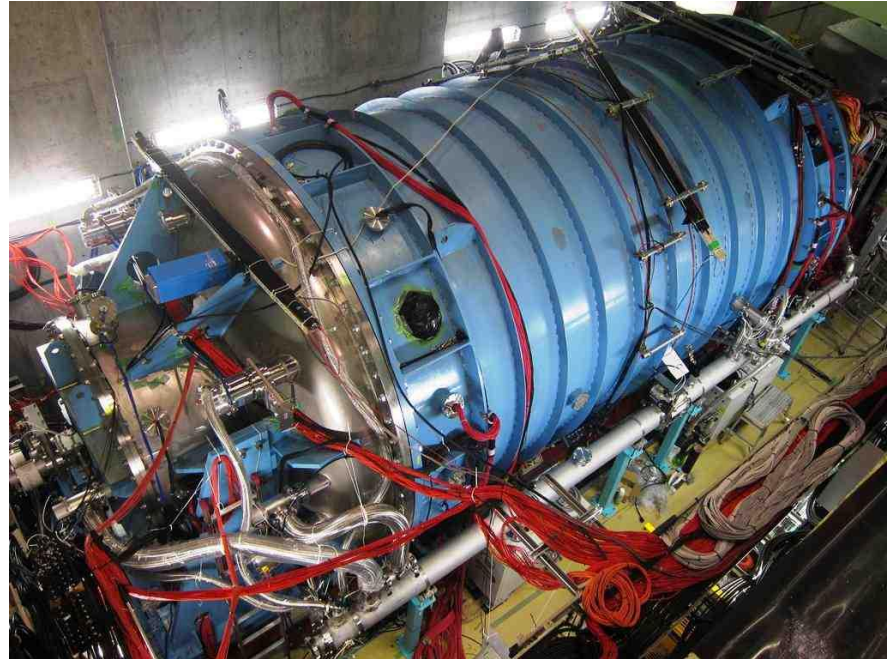
High energy physics experiment at the J-PARC (Japan Proton Accelerator Research Complex) Facility in Japan - http://koto.kek.jp/

Experiment measures the decay rate of neutral K-mesons (kaons) into neutral pi-mesons and a pair of neutrinos

# Requirements

Storage: Estimated storage for accommodating the MC and MC/Analysis portions of the pipeline for a full experiment cycle is in the order of 200 TB per experiment cycle
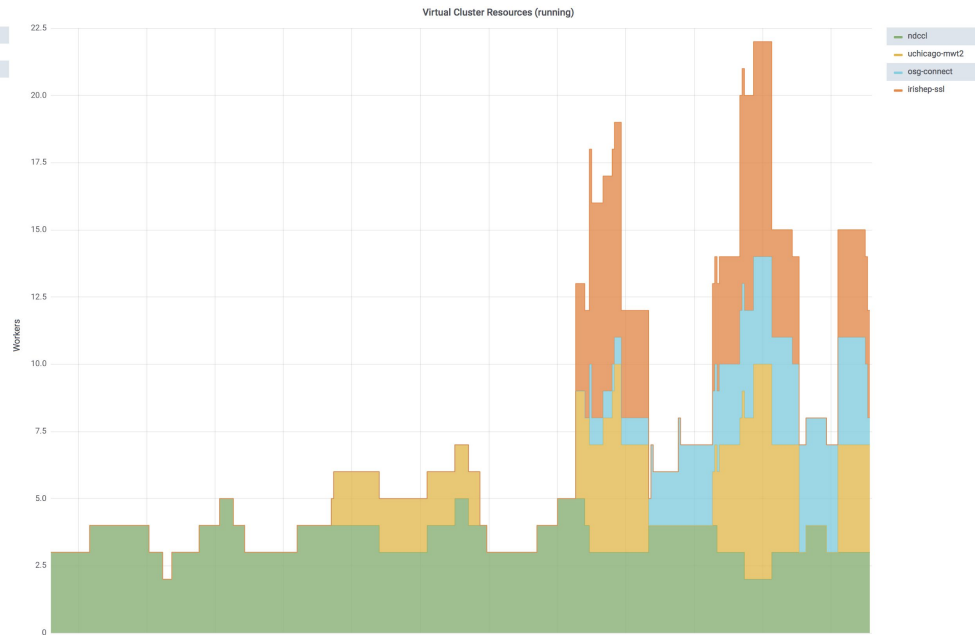
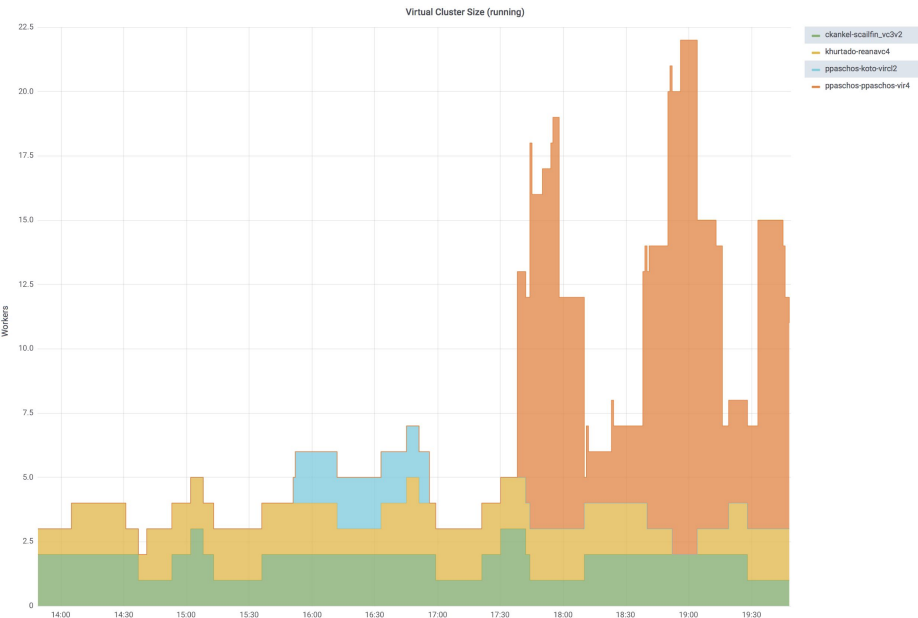Software stack: There are two applications running on the KEKCC machine that needed to be rebuilt on the OSG side under RHEL7. A KOTO GEANT package and an analysis collection of tools. Most of the effort so far is to provide the environment for the software to be built and run on remote OSG sites

Submission scripts: KEKCC uses the LSF scheduler (bsub). Submission scripts needed to be modified to HTCondor
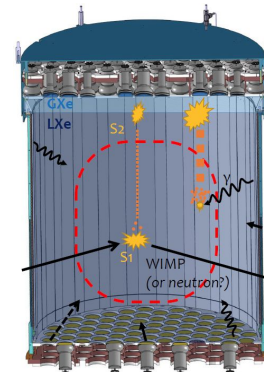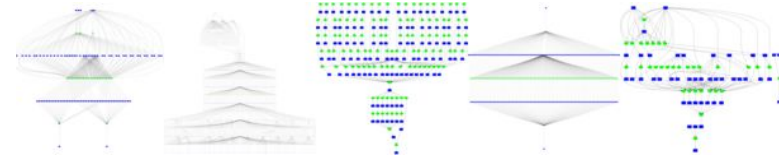
# VC3 job submission for KOTO

- OSG Connect (blue), UC (yellow), IRIS−HEP−SSL (orange)
- VC3 evenly distributes the load across resources

# Working Middleware and Applications

- Various Bioinformatics Workflows
  - Makeflow + HTCondor + BWA, Shrimp, BLAST.
- Lobster CMS Data Analysis
  - Work Queue + Builder + CVMFS
- Lifemapper
  - Work Queue + Builder
- Spark CMS Data Analysis
- South Pole Telescope (SPT−3G) Analysis Framework
  - HTCondor Jobs + Docker/Shifter + CVMFS
- XENON1T Analysis Framework
  - Pegasus + HTCondor + CVMFS
- MAKER Bioinformatics Pipeline
  - Work Queue + Builder
- IceCube Simulation Framework
  - HTCondor



52

# Thoughts and Lessons Learned

# Successful Ideas that Worked

- Automated Provisioning of Cluster Frameworks Across Heterogeneous HPC Resources
- Automated Deployment of User-Level Software
  - (Build itself may require substantial cluster resources)
- Internal Architecture for Robust Service Deployment
  - (Surprisingly complex state machine.)
- Multi-Layer Configuration Enables Collaborations
  - Credentials go deep into the framework

# Things That were Harder than Expected

- New Security Requirements (MFA) Spread Very Quickly
  - Working Idea: Let user login and "pull" the cluster blueprint to the local site where it can be deployed.
- Difficulty of Curating Site/Cluster Details due to Churn
  - Choice: Either invest more in active curation, or involve users in the configuration process.
- Debugging Complex Service Deployments
  - Working Idea: Test before using, expose stages of success to user, give concrete feedback for fixing.

# Portable Ideas to Keep and Use Elsewhere

- Importance of the Head Node as Infrastructure
  - Allocatable, Personal, Shareable, Configurable
- Self-Service Deployment Model
  - Generalization of Infrastructure-as-Code:
  - OpSys + Software + Middleware + Entry Point
- A Model for Sharing Institutional Resources
  - Classic: Submit My Work to Remote Environments
  - VC3: Deploy My Environment on Remote Resources

# Technology Evolution 2016-2020

Evolution from VMs to Containers:

- Swap out OpenStack in favor of Kubernetes, which is lightweight and better supports process automation.
- Docker is being supplanted by Singularity, and other container environments.

Evolution from credentials to capabilities:

- ssh-keygen -> Upload Pubkey -> RSA-SSH  (Standalone)
- Globus Auth -> GSI Token -> GSI-SSH         (DCDE Fed.)
- Globus Auth -> Web Tokens -> WebTokens-SSH  (Web Native)
- SciTokens? (scitokens.org)

# Towards Federated Ops with SLATE

- Remotely manage edge services at sites by **expert teams** from **trusted organizations**

- Deploy updates more quickly & introduce new services more easily

- Save time and effort for local site admins -- towards **OSG NoOps**

- Edge federation via lightweight server/client overlay using **Kubernetes**, the industry leading container orchestration platform

- Software catalog, with push button deploy using vetted **Helm** charts



```
$ slate instance list
$ slate instance delete <instance name>
$ slate app install --group atlas-xcache --cluster uchicago-prod
--conf MWT2.yaml xcache
```
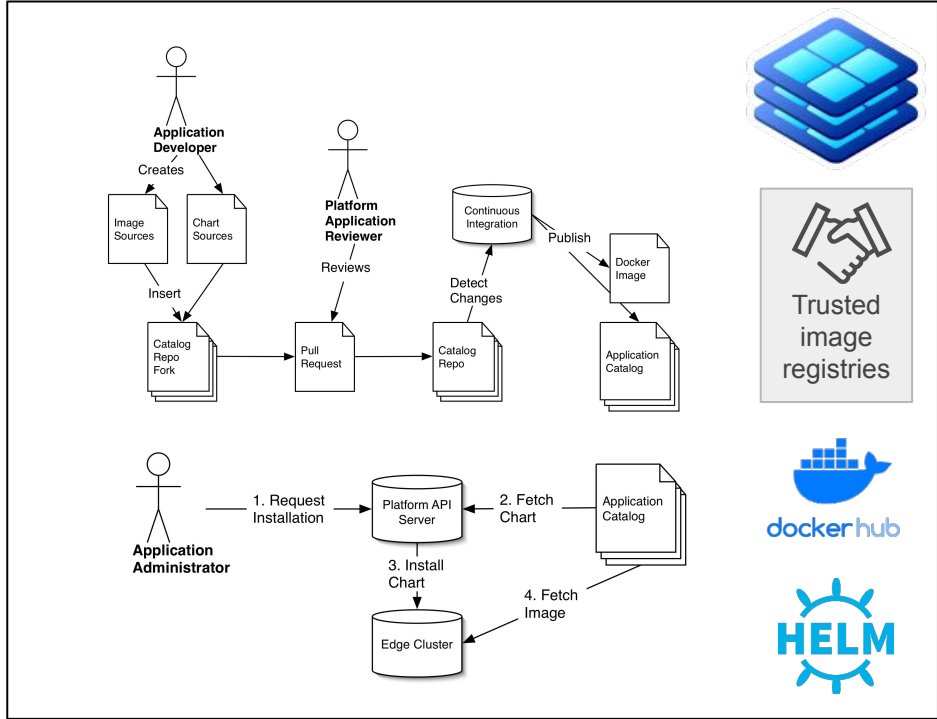
- SLATE - a value added K8s distribution
  - Support for CVMFS, ingress controller (multi-tenant, scoped privileges), Prometheus monitoring, **curated application catalog w/ Jenkins CI**
- Site security & policy conscious
  - SLATE works as an unprivileged user
  - Single entrypoint via **institutional identity**
  - Site owner controls group whitelists & service apps; **retains full control**
- With OSG, WLCG, trustedci.org & others working to establish a "CISO compliant" security posture and **new trust delegation model**



SLATE creates secrets and XCache deployment on cluster

XCache Container Download

Kubernetes objects instantiated

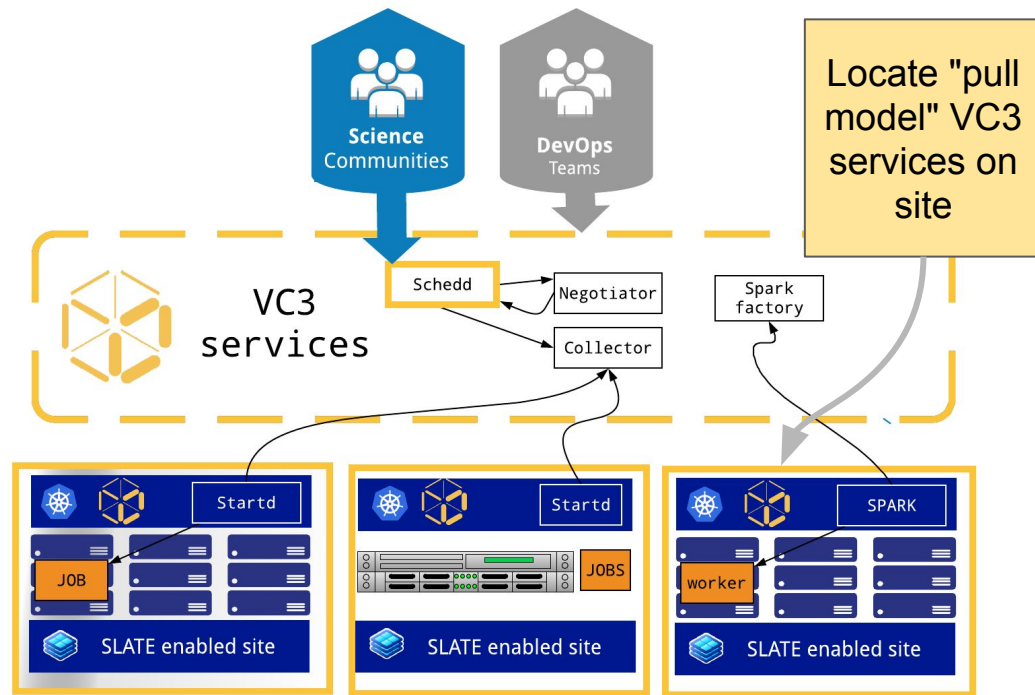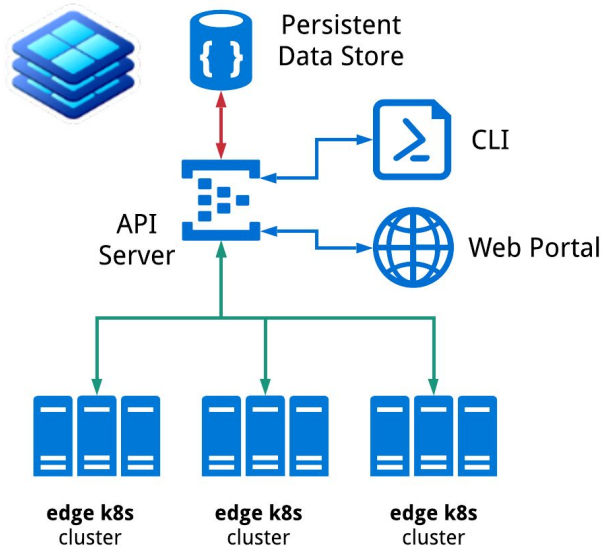Pod starts up, registers itself in AGIS

| < 5 min | 5-10 min | < 5 min | 5-10 min |

**A data caching network deployed in less than 20 minutes.**

*Upgrades are as simple as re-deploying.*

# VC3 and SLATE synergies

## SLATE: federated service orchestration in the SciDMZ
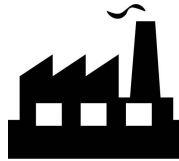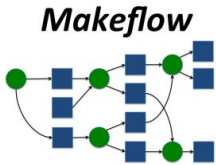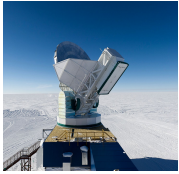
# Open Problems in Federated Systems

- Troubleshooting Across Systems
  - Many stakeholders, system churn, lack of evidence.
  - Troubleshooting posed as a distributed database query?

- State Management and Garbage Collection
  - Remote creation of containers for local sites, how long must they be kept around?

- Socio-Technical Understanding of Delegation
  - User X submits jobs to user Y's account on site S. Ok?
  - Technical ability to represent non-local account?

# Collaborators and Connections

# VC3

Virtual Clusters for Community Computation

https://www.virtualclusters.org

@virtualclusters

**New users signup: http://bit.ly/vc3-signup**

**Register your HPC: http://bit.ly/vc3-new-resource**

# VC3 Funding and Team

Funded by DOE Office of Advanced Scientific Computing Research (ASCR) and NSF Next Generation Networking Services (NGNS)

PIs: Rob Gardner (UC), Douglas Thain (ND), and John Hover (BNL)

co-PIs:  David Miller (UC), Paul Brenner (ND), Mike Hildreth (ND), Kevin Lannon (ND)

dev-team: Lincoln Bryant (UC), Benedikt Riedel (UC), Suchandra Thapa (UC), Jeremy Van (UC), Kenyi Hurtado Anampa (ND), Ben Tovar (ND), Jose Caballero Bejar (BNL).

**U.S. DEPARTMENT OF ENERGY** | Office of Science