



CMS experience using VC3

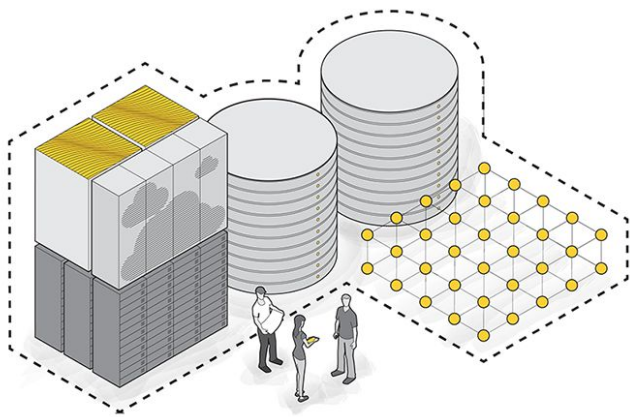
for provisioning spark clusters and
deploying T3s on top of campus resources
at the user level

Kenyi Hurtado
University of Notre Dame
khurtado@nd.edu

Outline

- What is VC3?
- VC3 CMS use case examples
 - Building Spark clusters on top of Global Pool Resources
 - Deploying a Tier 3 on top of campus resources (with no grid-friendly environment or root access level to the resource)
- Conclusions

What is VC3?

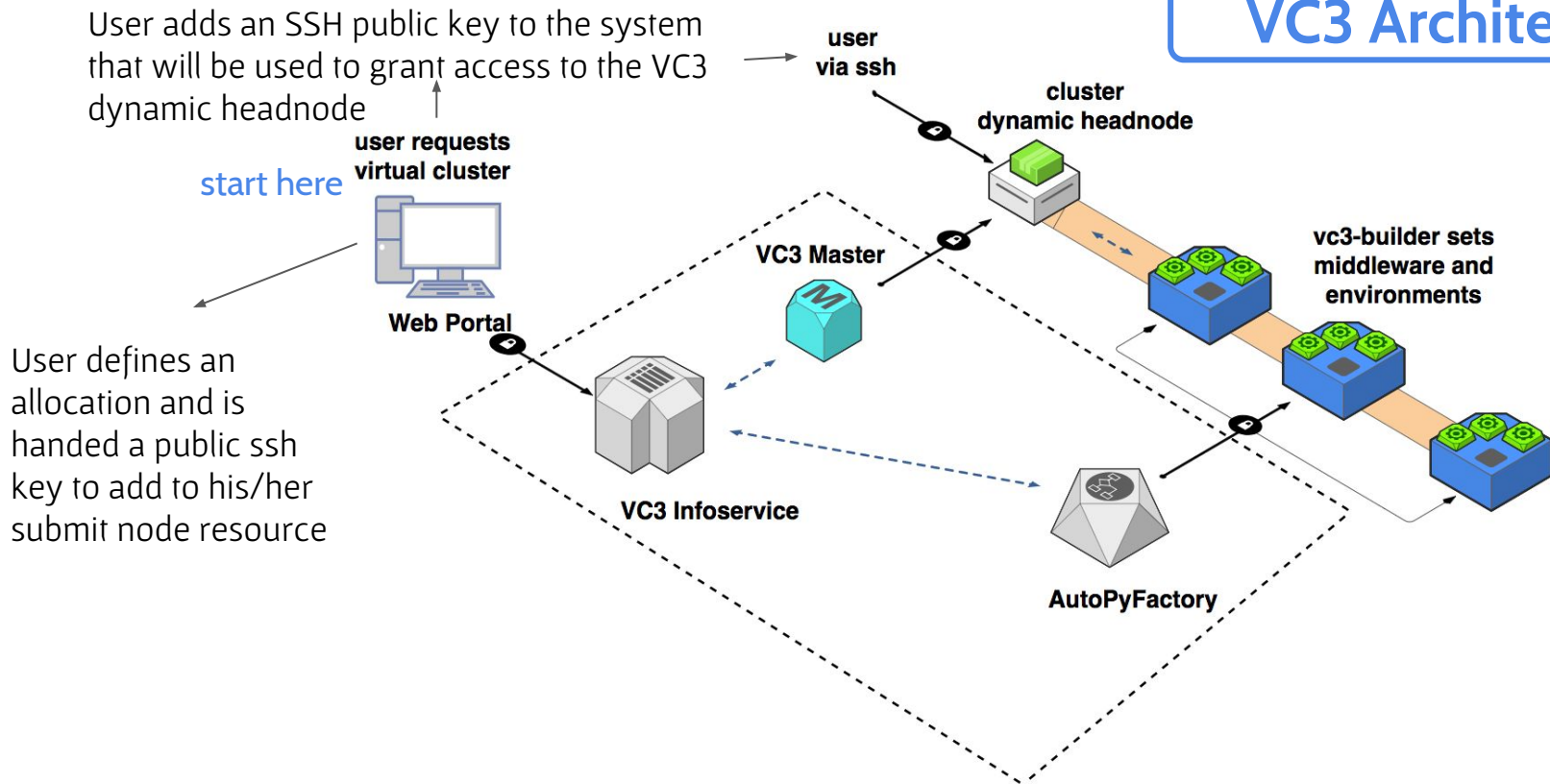


VC3: A platform for provisioning cluster frameworks over heterogeneous resources for collaborative science teams

<https://www.virtualclusters.org>

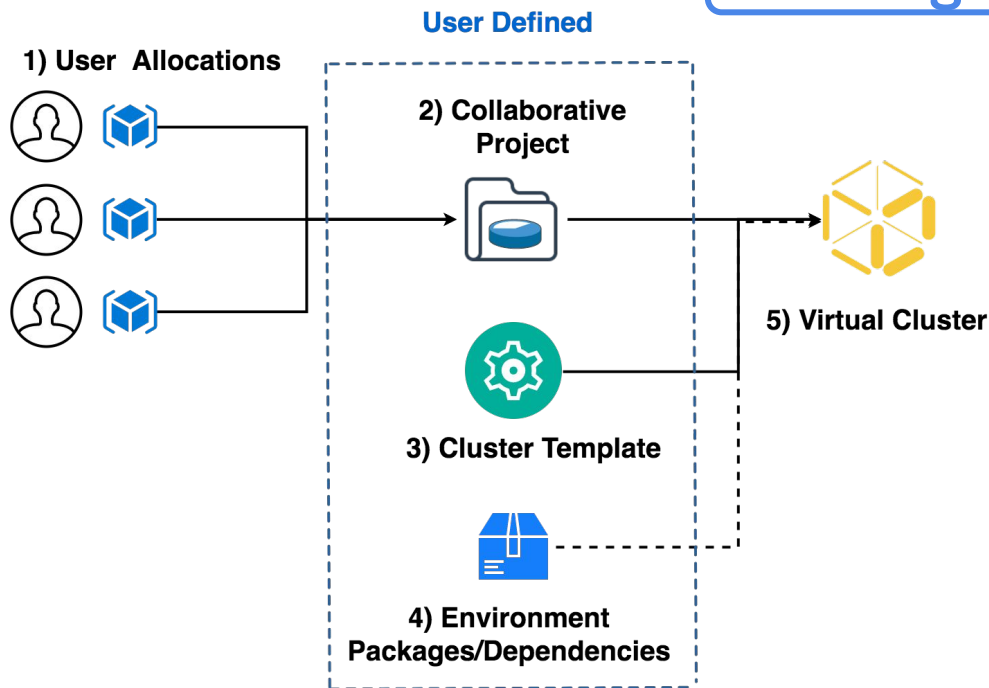
What is VC3?

VC3 Architecture



What is VC3?

Creating a Virtual Cluster











Use case example 01

Provisioning Spark clusters on top of Global Pool resources

Creating a Spark Cluster – Step 1

← → ↺ <https://www.virtualclusters.org/resource>

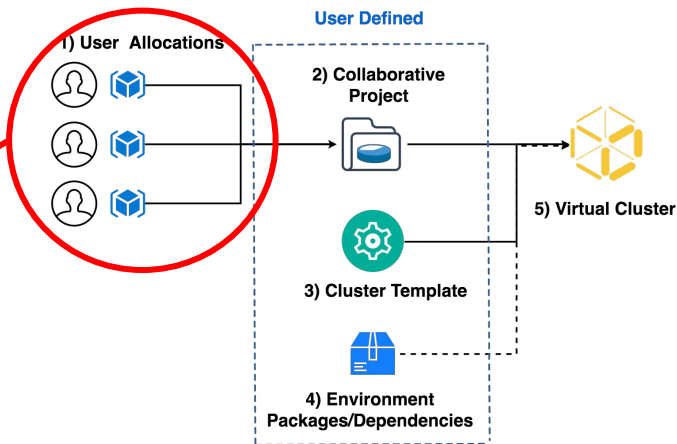
VC3  News Community Documentation

-  Resources
-  Allocations
-  Projects
-  Environments (beta)
-  Virtual Clusters
-  Monitoring
-  Admin

Platform	Resource	Virtual Cluster
Research Computing Center (RCC)	Chicago Research Computing Center (RCC)	
Stampede2	Texas Advanced Computing Center (TACC)	Stampede2 Super Computer
CMS Connect	CMS	CMS Connect
CoreOS	University of Chicago	CoreOS/Kubernetes Cluster with HTCondor Overlay
UCT3	University of Chicago - Enrico Fermi Institute	UChicago ATLAS Tier 3
ND CCL	University of Notre Dame Cooperative Computing Lab (CCL)	Notre Dame CCL Job Gateway

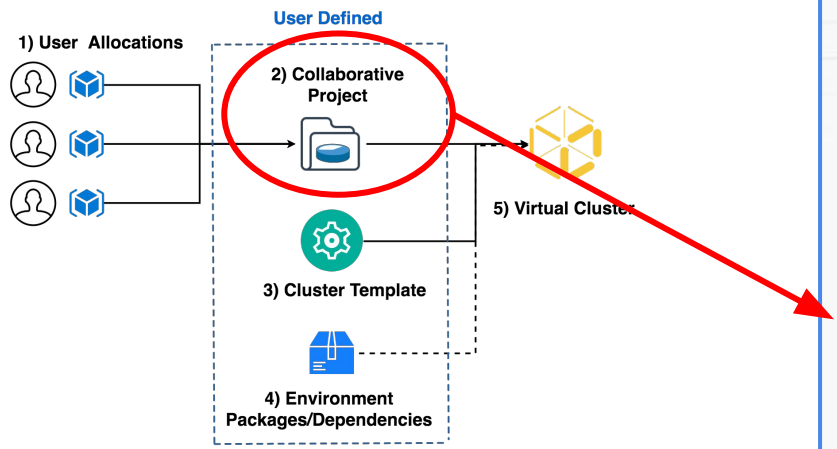
VC3 website:

<https://www.virtualclusters.org/>



Using CMS Connect to access the CMS Global Pool

Creating a Spark Cluster – Step 2



Create a New Project

A project connects your allocations to your team members

▪ - INDICATES REQUIRED FIELD

PROJECT NAME ▪

bigdatacms

PROJECT MEMBERS

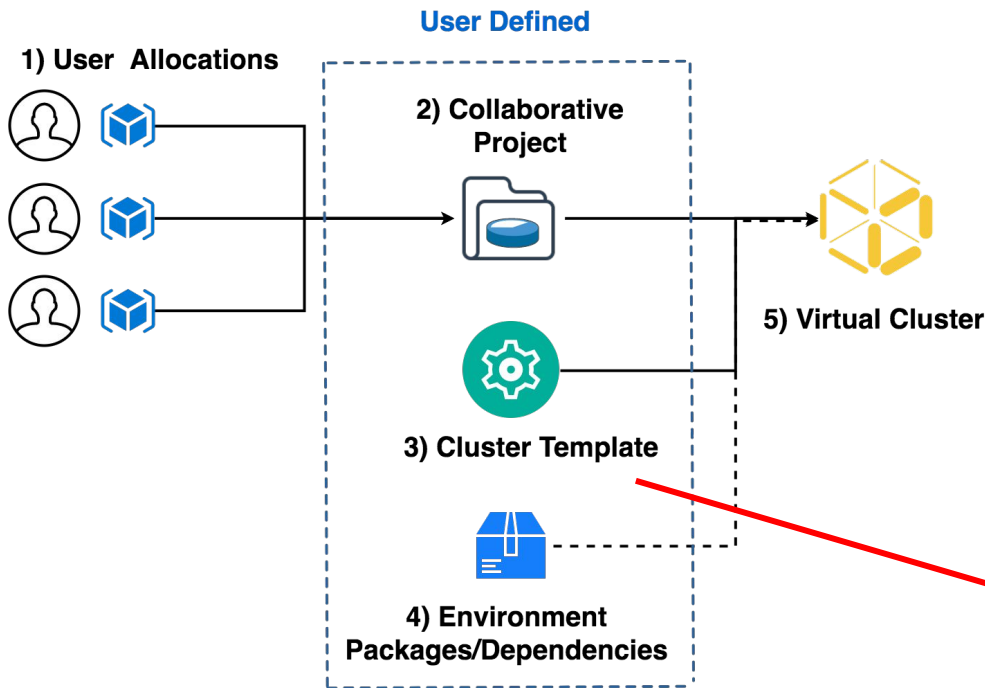
Benjamin Tovar, Benedikt Riedel, Suchandra Thapa ▲

SELECT ALLOCATION:

khurtado-cms-connect ▲

Creating a Spark Cluster – Step 3

Use CMS Connect to access the CMS Global Pool



Launch New Virtual Cluster

Project: vc3-team

* - INDICATES REQUIRED FIELD

VIRTUAL CLUSTER NAME (A-Z, 0-9, _ AND -)*

CLUSTER TEMPLATE FRAMEWORK *

- ☒ Please Select Framework
- ☐ HTCondor
- ☐ WorkQueue
- ☐ JupyterLab
- ☒ Spark

Creating a Spark Cluster – Step 4



Spark Master at spark://128.135.158.246:7077

URL: spark://128.135.158.246:7077

REST URL: spark://128.135.158.246:6066 (*cluster mode*)

Alive Workers: 10

Cores in use: 20 Total, 0 Used

Memory in use: 74.2 GB Total, 0.0 B Used

Applications: 0 Running, 0 Completed

Drivers: 0 Running, 0 Completed

Status: ALIVE

Workers

Worker Id	Address
worker-20181024210235-169.228.132.112-39013	169.228.132.112
worker-20181024210254-169.228.132.112-33352	169.228.132.112
worker-20181024210842-169.228.131.229-37426	169.228.131.229
worker-20181024212048-169.228.132.142-39969	169.228.132.142
worker-20181024212233-169.228.130.186-30117	169.228.130.186
worker-20181024212246-169.228.130.186-5423	169.228.130.186
worker-20181024212423-169.228.131.46-42113	169.228.131.46
worker-20181024212423-169.228.131.46-43409	169.228.131.46

Virtual Cluster: khurtado-spark_ucsd

[Terminate Cluster](#)

STATE OF VIRTUAL CLUSTER

Running

All requested compute workers are running.

Owner

Kenyi Anampa

Project

sparkcms

Status

KHURTADO-SPARK_UCSD
Cluster Framework: spark

Requested 10

Running 10

Queued 0

Error 0

Using the Spark Cluster

https://github.com/SiewYan/PadovaBIGDATA/blob/docker_dev/Docker_dev/notebooks/Zpeak_Nanoaod-SPARK.ipynb



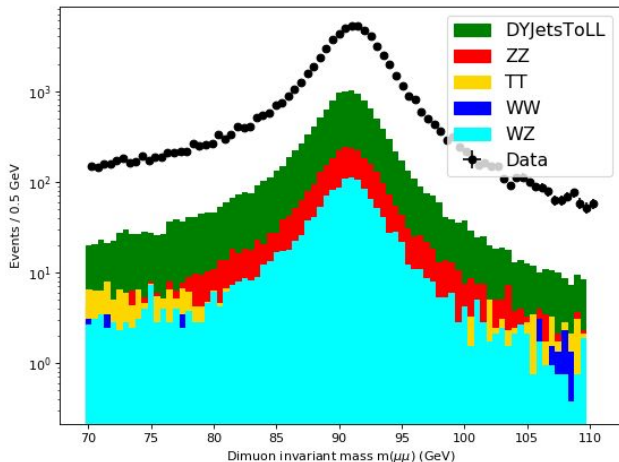
Application: Zpeak_Nanoaod-SPARK

CMS BIG-DATA group example

ID: app-20181024213924-0000
Name: Zpeak_Nanoaod-SPARK
User: khurtado
Cores: Unlimited (20 granted)
Executor Limit: Unlimited (10 granted)
Executor Memory: 7.0 GB
Submit Date: 2018/10/24 21:39:24
State: RUNNING
Application Detail UI

Executor Summary

ExecutorID
7
8
0
4
9
3
5
2
6
1

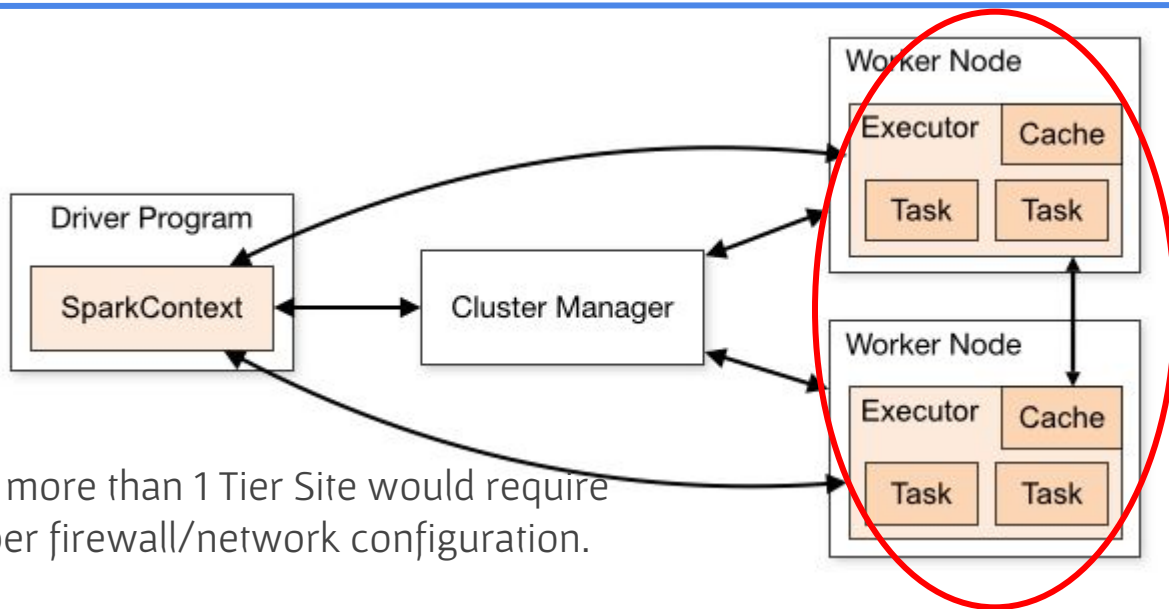


```
vc3-spark-submit --executor-memory 7G --conf spark.sql.caseSensitive=true --conf
spark.serializer=org.apache.spark.serializer.KryoSerializer --packages
org.diana-hep:spark-root_2.11:0.1.16,org.diana-hep:histogrammar-sparksql_2.11:1.0.4
Zpeak_Nanoaod-SPARK.py
```

Cores	Memory	State	Logs
2	7168	RUNNING	stdout stderr
2	7168	RUNNING	stdout stderr
2	7168	RUNNING	stdout stderr
2	7168	RUNNING	stdout stderr
2	7168	RUNNING	stdout stderr
2	7168	RUNNING	stdout stderr
2	7168	RUNNING	stdout stderr
2	7168	RUNNING	stdout stderr
2	7168	RUNNING	stdout stderr
2	7168	RUNNING	stdout stderr

Limitations

1. Spark workers need to communicate between each other for some operations (e.g.: shuffling/repartition)



That means, spark clusters using more than 1 Tier Site would require workers with public IPs and proper firewall/network configuration.

Due to the above, only 1 Tier Site per spark cluster is used in this use case.

This is currently done by exporting an environment variable in `login.uscms.org` (`CONDOR_DEFAULT_DESIRED_SITES`)

Limitations

2. We use singularity containers in the Global Pool. There is a bug in singularity when using the "--contain" (used in CMS) option that affects running spark workers/slaves.

<https://github.com/sylabs/singularity/pull/1420>

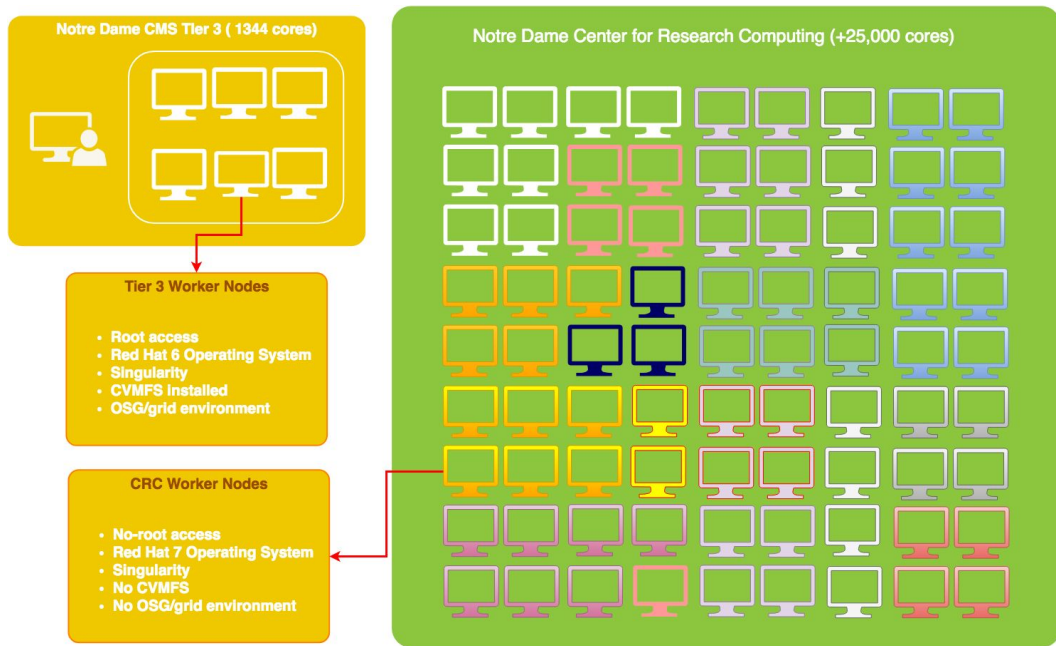
```
09/18/2018 07:17:13 PM [ERROR] There
was an error executing job:
/srv/x86_64/redhat6.10/spark/v2.2.1/bin/spark-class: line 80: /dev/fd/62: No
such file or directory
```

Due to the above, only Sites with Singularity 2.5.2+ will work.
Good news is, 2.6.0 is already available in e.g.: OSG repositories.

Use case example 02

Building a Tier 3 on top of campus resources

The Notre Dame campus cluster



The Notre Dame CMS group operates a Tier 3 with about 1,300 cores for local and grid analysis jobs.

In addition to this, the Center for Research Computing (CRC) provides an opportunistic campus cluster with over 25K cores of computing power researchers have access to, but these resources lack the software components and environment needed by CMS analysis workflows.

The Notre Dame campus cluster

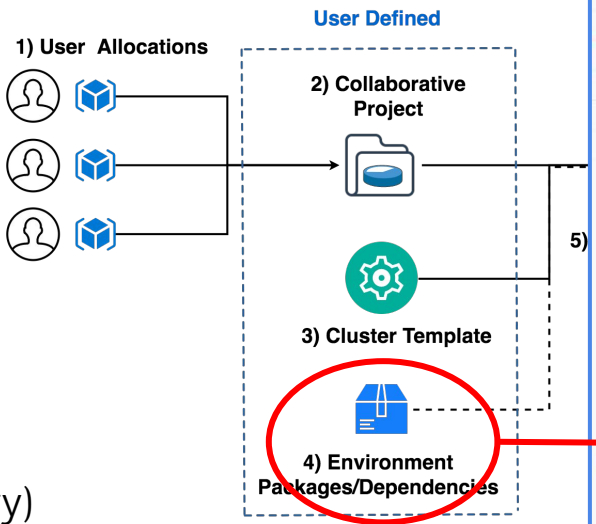
ND CRC cluster characteristics:

- ✓ – HTcondor cluster with +25K cores
- ✓ – Singularity available
- ✗ – **No CVMFS** → CMS glideins use singularity containers located via CVMFS. CMSSW also located here.
- ✗ – **No grid environment (voms, gfal tools, xrootd, etc)** → Can't use VOMS proxies, needed to submit to the Global Pool and access data.
- ✗ – **The ND CMS has no root access to it, just user allocations.** → Can't use yum to install the above using the OSG repositories
- ✗ – **No Compute Element** → Can't install/maintain HTCondor-CE for this resource

Building a Tier 3 on top of campus resources – Step 1

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)



Create New Environment

• INDICATES REQUIRED FIELD

ENVIRONMENT NAME *

PACKAGE LIST

OPERATING SYSTEMS LIST (OPTIONAL)

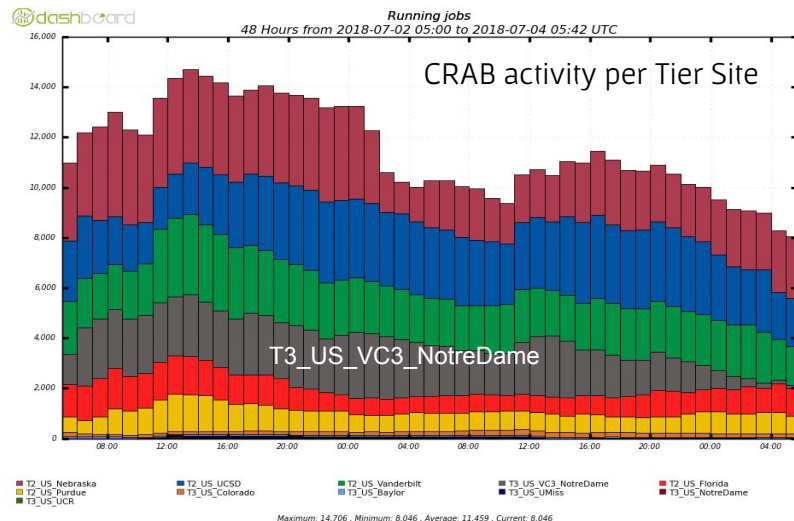
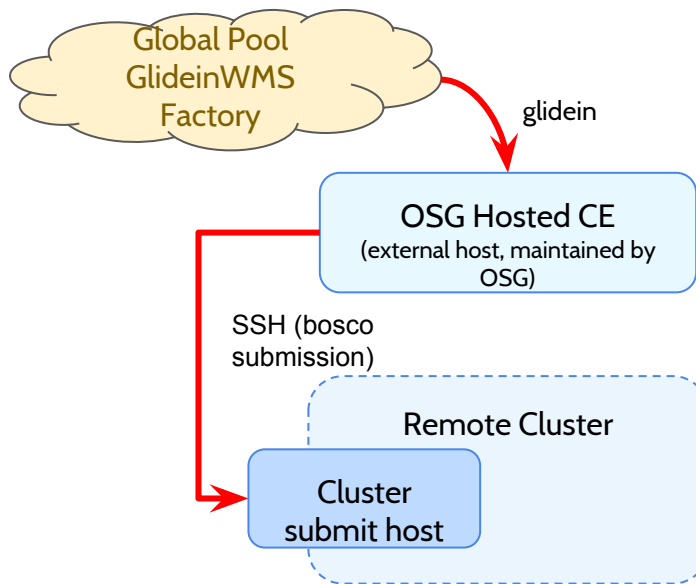
Cancel
Create New Environment

Also, oasis-wn-vc3ndcms links
[/cvmfs/cern.ch/SITECONF/local](https://cvmfs.cern.ch/SITECONF/local) to:
https://gitlab.cern.ch/SITECONF/T3_US_VC3_NotreDame

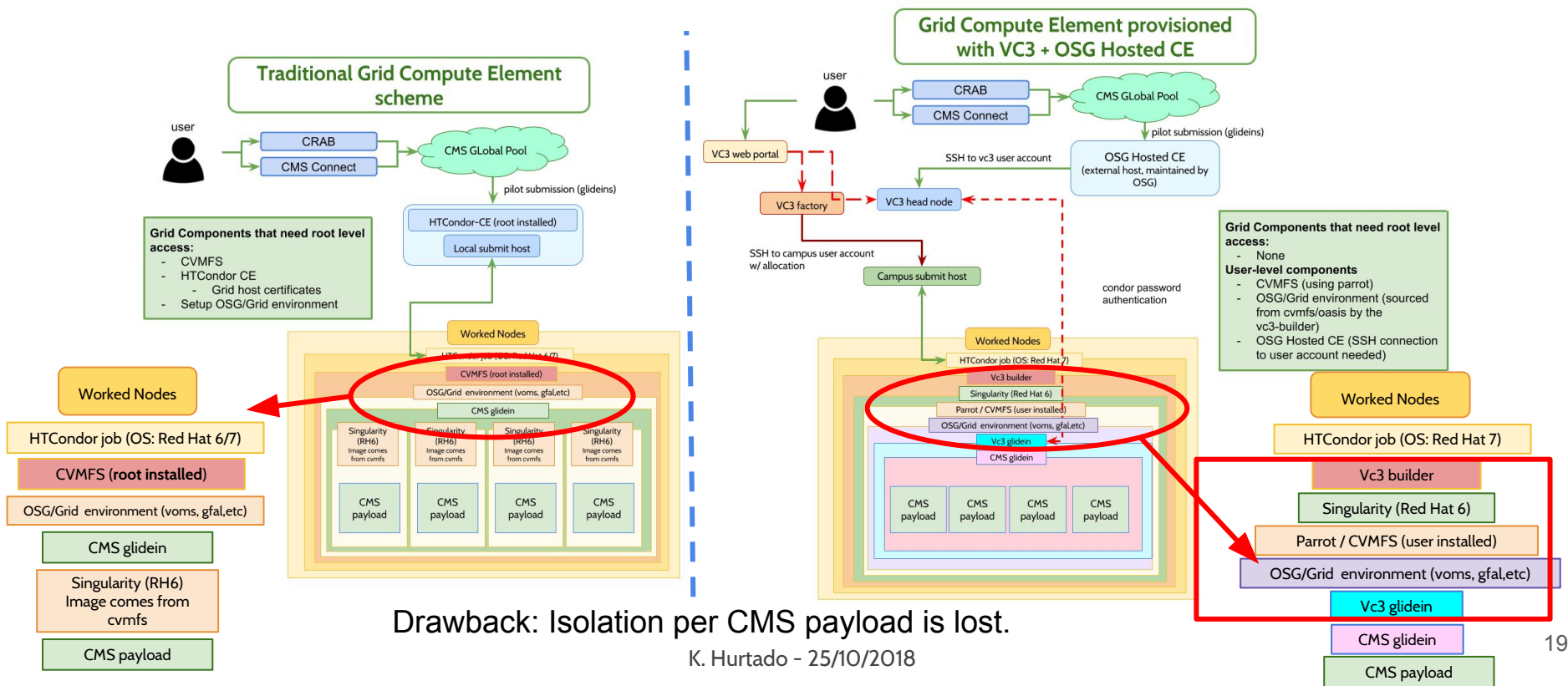
Building a Tier 3 on top of campus resources – Step 2

- 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.

Hosted HTCondor CE over SSH



Traditional vs VC3+OSG Hosted CE T3 components



Conclusions

- In these use case examples, we have successfully used VC3 to:
 - Use resources already available in CMS Global Pool, but that can't be used with Spark as-is.
 - ... by provisioning the spark cluster middleware on top of condor.
 - Add and use resources we normally can't use with CRAB
 - ... by installing and setting up the grid required environment without root access level and plugging the resource into the Global Pool using an OSG Hosted CE

Feel free to contact me at khurtado@nd.edu If you are interested in using VC3 for any of these two use-cases.

VC3 ad slide

Virtual Clusters for Community Computation

<https://www.virtualclusters.org>
@virtualclusters

Limited beta signup: <http://bit.ly/vc3-signup>



Supported by the Department of Energy Office of Advanced Scientific Computing Research and Next Generation Networking Services, Solicitation DE--FOA--0001344 (DDRM), Proposal 0000219942.

Backup slides

CMS Analysis workflows and submission services

CRAB and CMS Connect give analysis users access to the Global Pool

