# HEPCloud: provisioning heterogeneous resources using GlideinWMS and HTCondor

Marco Mambelli for the HEPCloud team

HTCondor week , Madison, WI

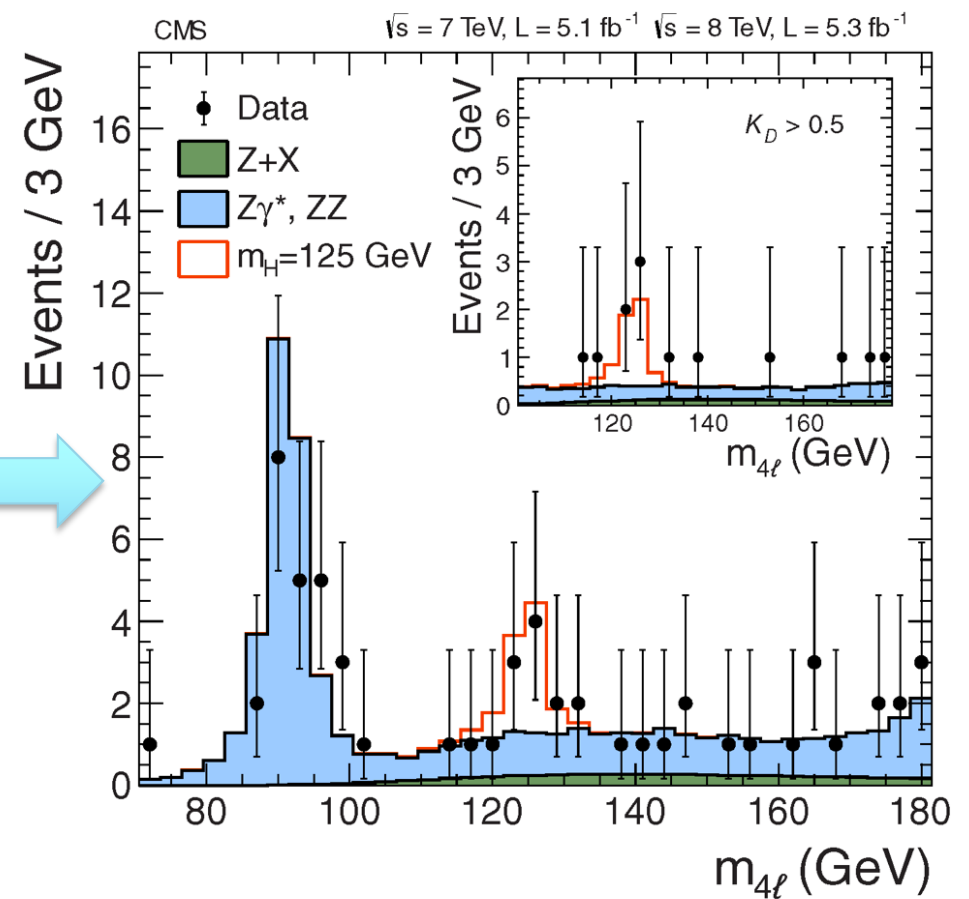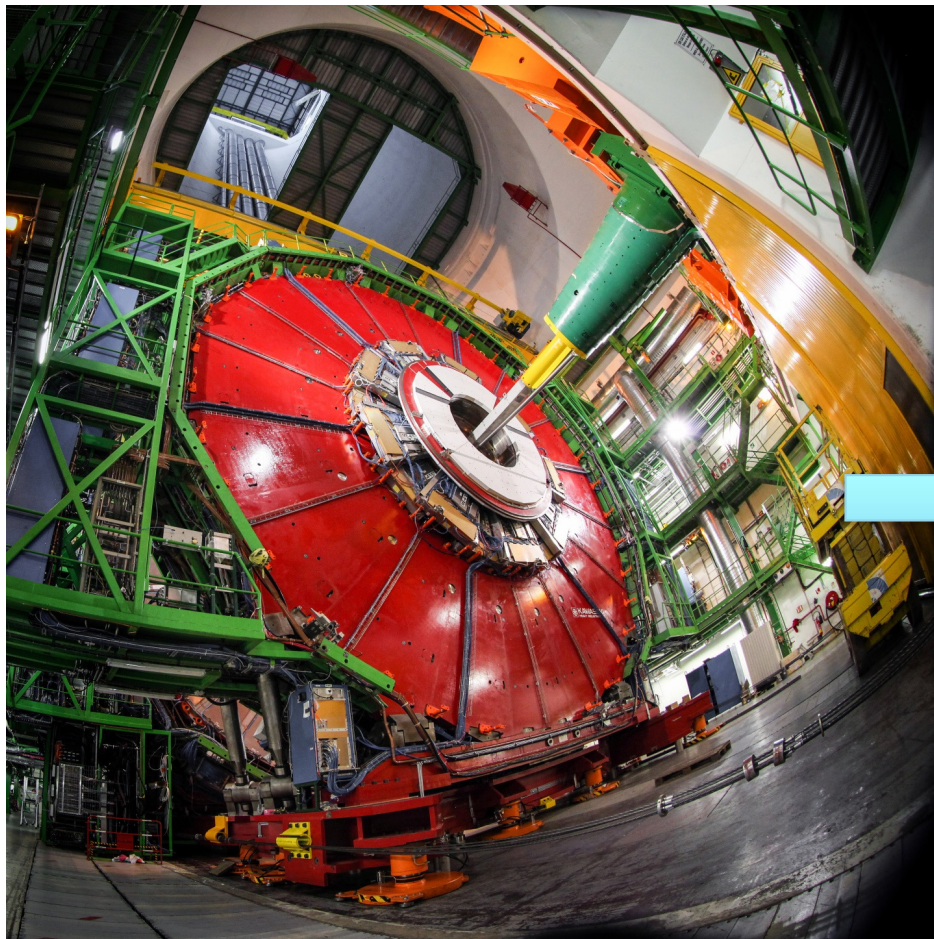24 May 2022

# Thank you!

- Development team
  - Brandon White
  - Bruno Coimbra
  - Hyun Woo Kim
  - Kyle Knoepfel
  - Lisa Goodenough
  - Patrick Riehecky
  - Shreyas Bhat
  - Steven Timm (operations/dev)
  - Vito Di Benedetto
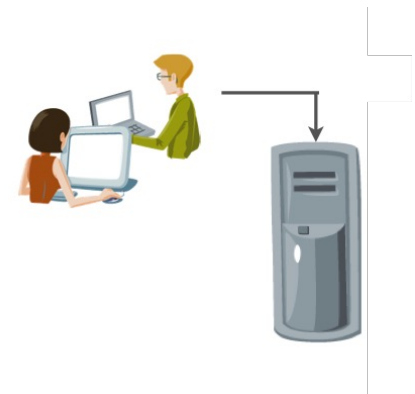  - Andrew Norman (project lead)

# From here … to there

# From dedicated supercomputers…

# ... to computer centers ...



Batch system

Oliver Gutsche l Software, Computing & Analysis - CMS DAS January 2017

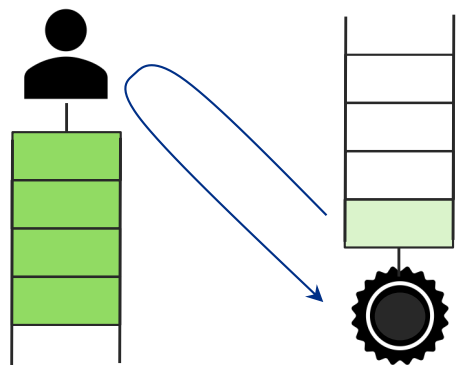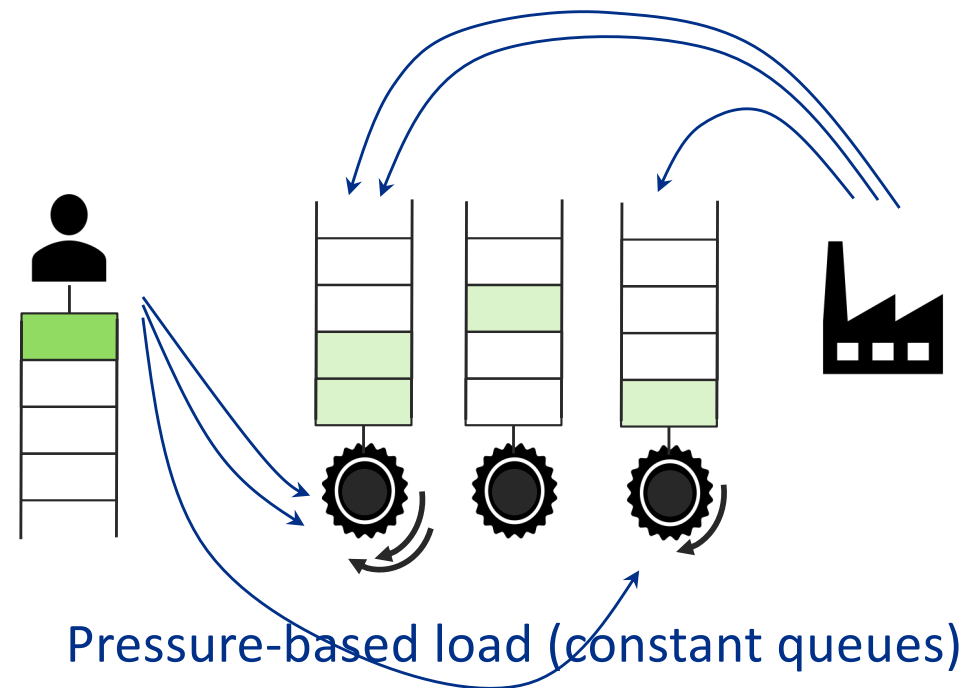09. January 2017

**Fermilab**

# ... to Grids ...

# Submission Infrastructure (up to the Grid)

- distributed High Throughput Computing (dHTC)
- (Semi-)dedicated (shared) resources
- Pilot based systems
  - resource validation
  - late binding
- Pressure based workloads



Pilot system

Pressure-based load (constant queues)

🔹 Fermilab
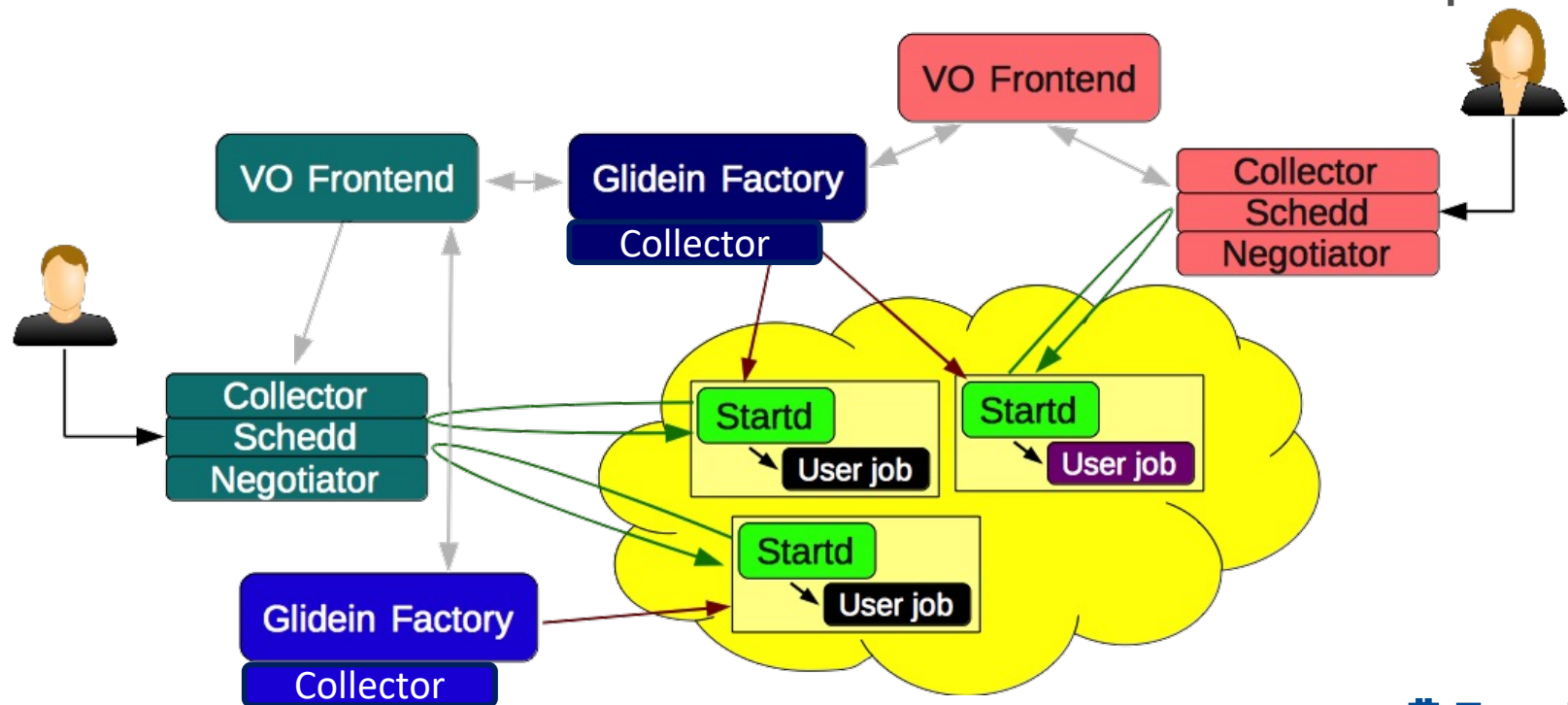
# GlideinWMS pressure-based pilot system

- Used by CMS, DUNE and Fermilab experiments
- Factories use HTCondor (vanilla or grid -batch/ec2/gce-universe) to submit Glideins, pilot jobs
- Frontends trigger the Glidein submissions
- Glideins start startds for a distributed HTCondor virtual pool
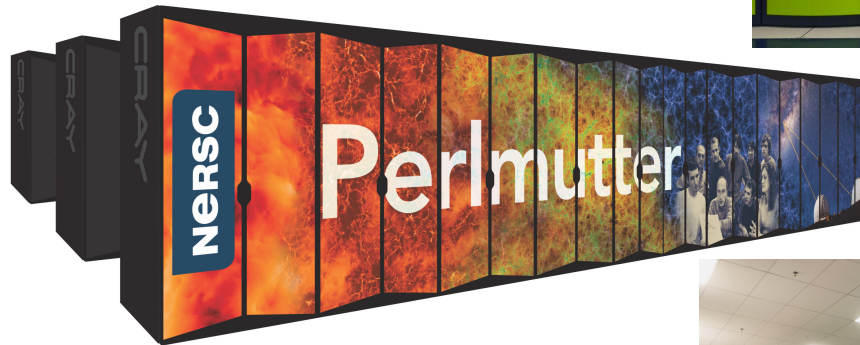
# ... to also Clouds and many supercomputers

# Increasing heterogeneity

- Clouds
  - On-demand
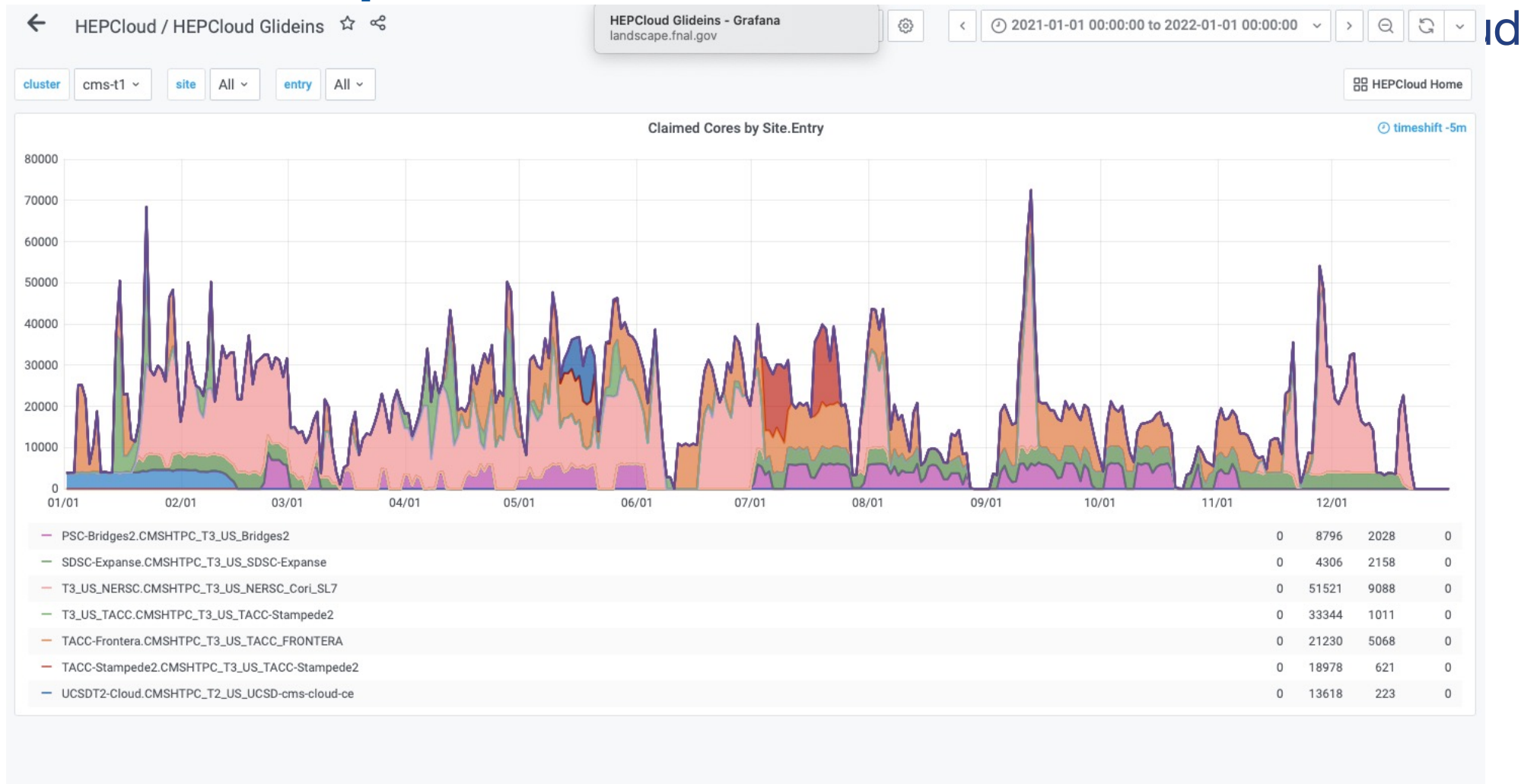  - Budget
  - Cost optimization
  - Services

- Supercomputers
  - Allocation
  - One-of-a-kind resources
  - Built for specific scopes

# HEPCloud (Facility)

- Built on top of dHTC (GlideinWMS and HTCondor)
- Portal, job routing, resource provisioning
- Decision Engine
  - Business rules
  - Figure of Merit: multidimensional optimization
- Commissioning of new resources

🔷 **Fermilab**

# HEPCloud operations: CMS 2021



Used all NERSC quota plus 90M NERSC-hours bonus
Used all XSEDE and FRONTERA quota 6 months before expiry
Used new UCSD/Azure source

2/2/22  S. Timm | HEPCloud Operations

🔷 Fermilab

# HEPCloud CMS 2021



160M CoreHr from successful jobs on NERSC, PSC, SDSC and TACC
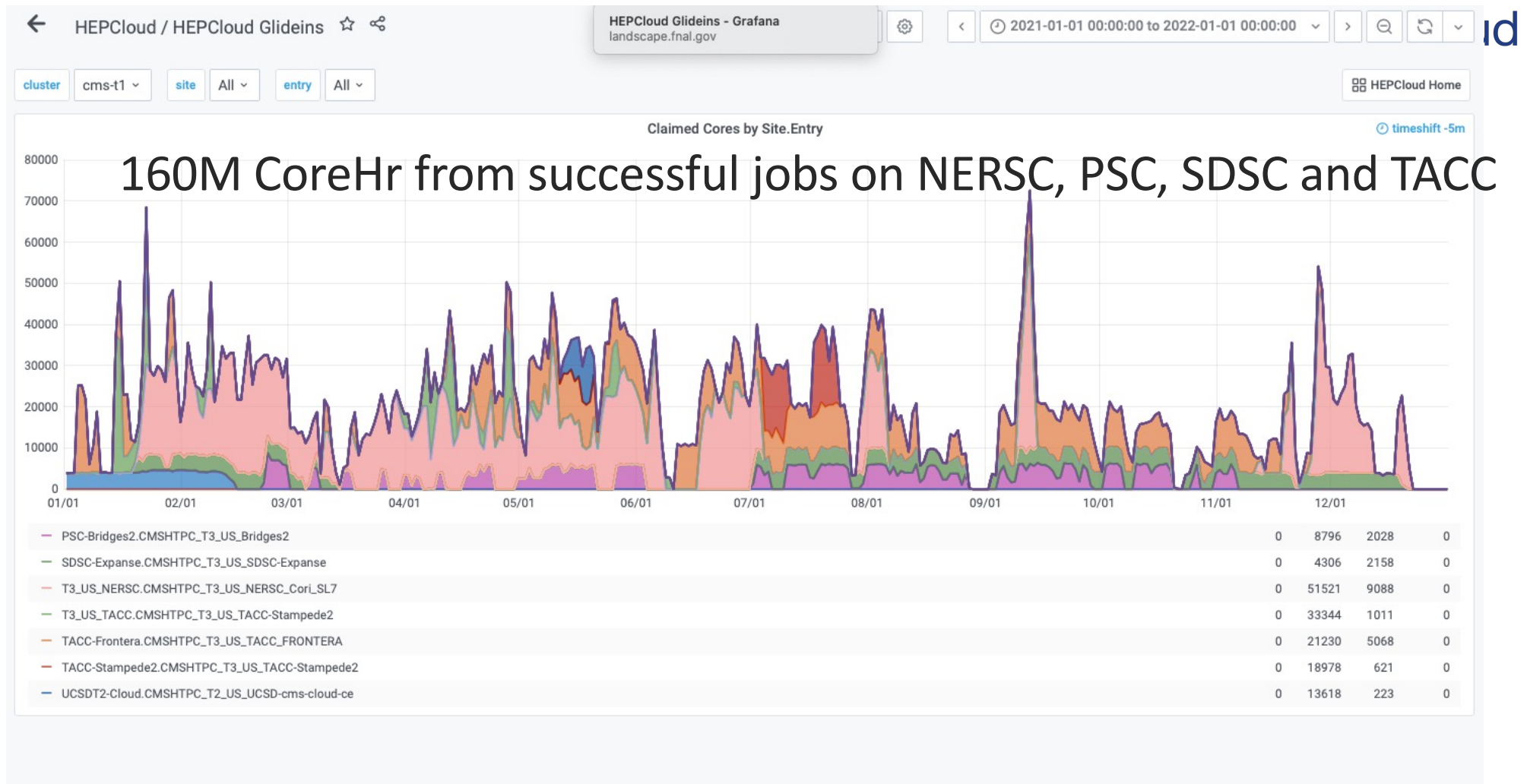
Used all NERSC quota plus 90M NERSC-hours bonus
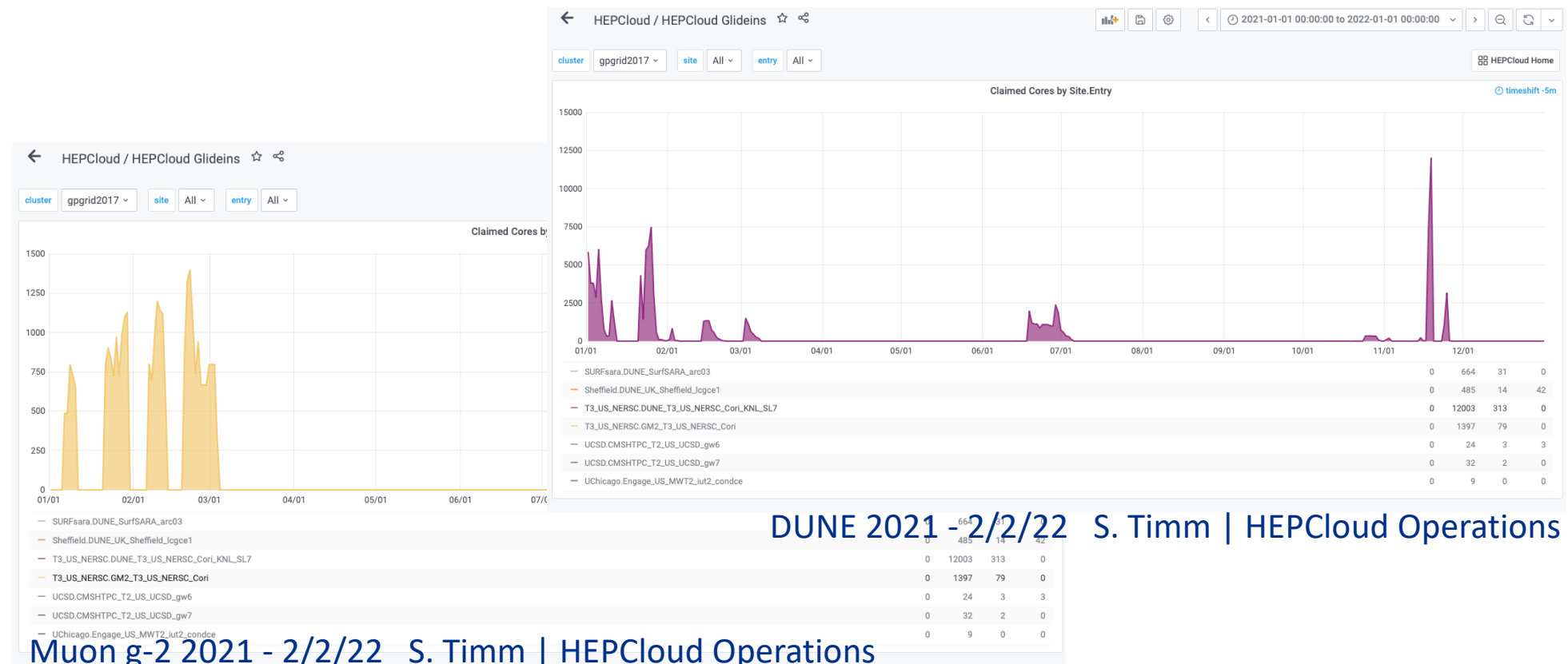Used all XSEDE and FRONTERA quota 6 months before expiry
Used new UCSD/Azure source
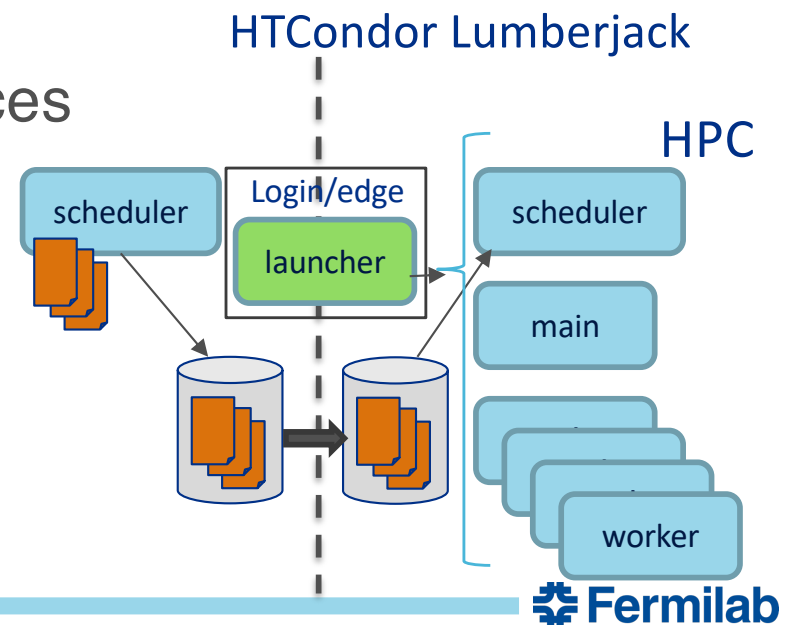
2/2/22   S. Timm | HEPCloud Operations

# HEPCloud serving different load types

- CMS workload has been steady through the year

- DUNE has campaign bursts

- Muon g-2 is also computing specific numbers in bursts



Muon g-2 2021 - 2/2/22   S. Timm | HEPCloud Operations

DUNE 2021 - 2/2/22   S. Timm | HEPCloud Operations

# HPC Onboarding

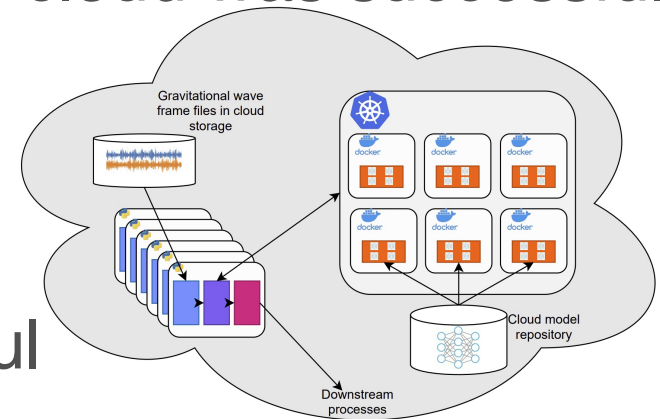- Sites with no network connectivity from the worker nodes
  - Theta integration and HTCondor Split-start and Lumberjack testing — *In progress see M.Acosta talk*

- Heterogeneous sites (CPUs, GPUs, large memory) — *In progress*
  - Evaluation of NERSC Perlmutter and Purdue Anvil

- HEPCloud provides
  - Single, uniform access to all resources
  - Expertise
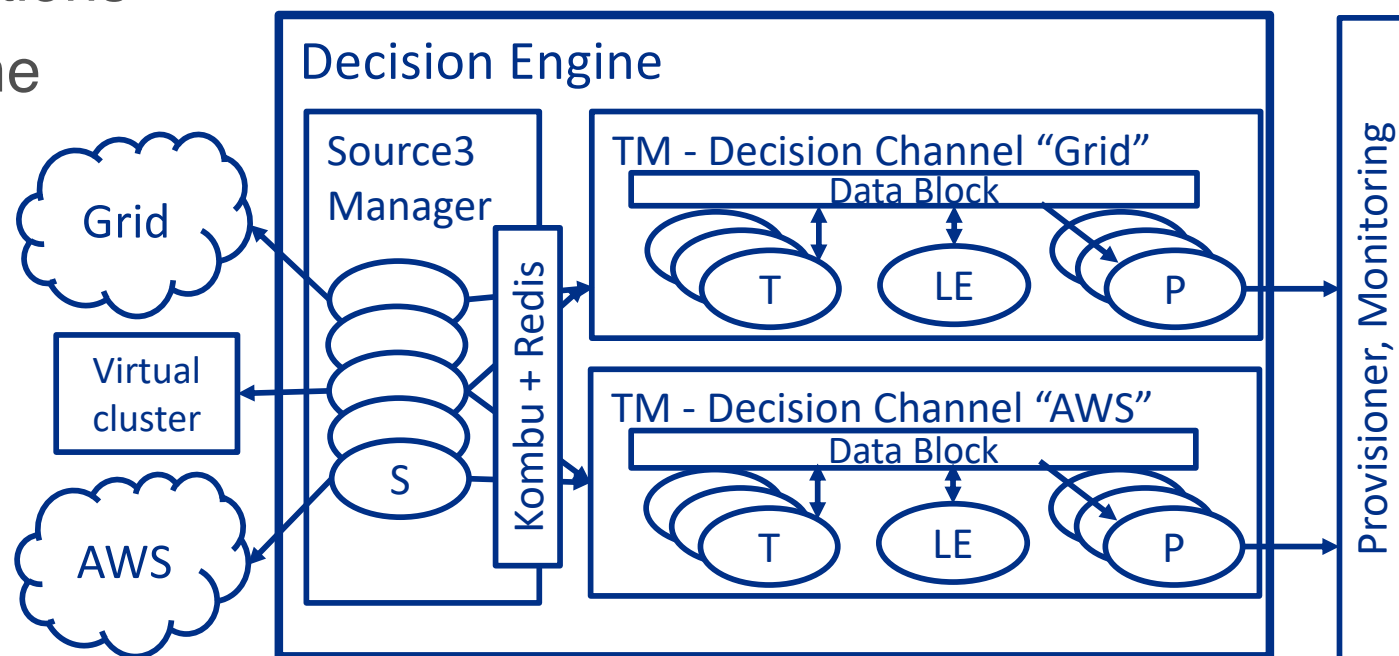  - Solutions to resources constraints

# Other activities

- MIT Inference server testing on Google cloud was successful last summer
  - Paper accepted by Nature Astronomy

    https://arxiv.org/pdf/2108.12430.pdf

- Early Rigetti Computing tests successful
  - Public company building superconducting quantum processors
  - Aspen-10, 32 qbit QPU, available as a service, QCS
  - Running Quantum applications (Quil) on real QPUs via cloud-hosted HTCondor

- Collaboration with RADICAL Cybertools (RCT)  In progress
  - More HPC resources
  - MPI and heterogeneous workflows

- Expanding to more applications and more users

# Decision Engine

- Sensing environment, taking provisioning decisions
- Sources
- (Decision) Channels
    - Task Manager
    - Transformations
    - Logic Engine
    - Publishers

# Integration test/Example configuration

- Test a basic HTC configuration running jobs on Factory provided pilots
  - DE starts and listens to the system (blue)
  - A job is submitted
  - This triggers a pilot running on the CE and registering in the user pool (red)
  - The job runs successfully on the pilot (green)
  - The DE and the pool shut down
- Runs on 3 virtual machines on Fermicloud

🔷 **Fermilab**

# Challenges where HTCondor could help

- Running Parallel/MPI jobs
  - Evaluate HTCondor Parallel universe
  - Use resources that do not have HTCondor as scheduler (most have SLURM)

- Expand HPC support
  - Limited network connectivity
  - Two factor authentication and other complex authentication schemas

See Maria Acosta's talk

- Credentials renewal
  - Glideins use tokens to authenticate w/ the pool and to access resources
    - Should install credmon?
    - Mechanism to update an input file

🔷 Fermilab

# Challenges where HTCondor could help (cont)

- Reserve jobs
  - Park jobs associated to a resource until a timeout is met
  - Give the jobs time to complete on the resource, then reclaim and match with other resources

- Well supported and stable BLAH and tarball distribution
  - Recent changes and incompatibility required extra work to run again on HPC resources

- Easier to debug HTCondor-CE and configuration
  - E.g. less quotes and backslash, or JSON or YAML

🔷 **Fermilab**

# Using HEPCloud and collaborating

- Open Source project on GitHub:
  https://github.com/HEPCloud
  - https://github.com/HEPCloud/decisionengine
  - https://github.com/HEPCloud/decisionengine_modules
- RPMs available (DE and DEM 2.0.0)
  - https://zenodo.org/record/6485889#.YmfuxvPML7E
  - https://zenodo.org/record/6485937#.YmfuyPPML7E
- Instructions:
  https://hepcloud.github.io/decisionengine/install.html#
- Simple installation for HTC pressure-based submission
  - https://github.com/HEPCloud/decisionengine/wiki/Decision-Engine-integration-test

🔷 **Fermilab**