



# Managing a Dynamic Sharded Pool

Anthony Tiradani

HTCondor Week 2019

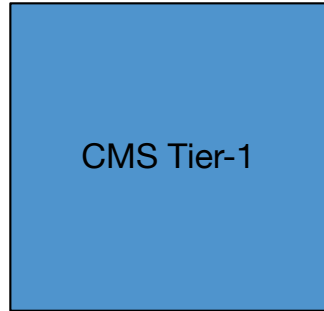
22 May 2019

# Introduction

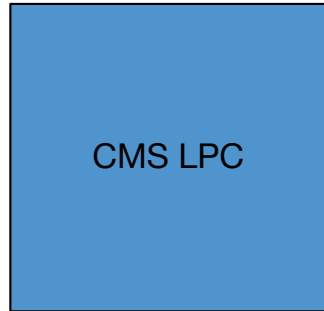
- Some archaeology from my time at Fermilab
  - Earliest archived Fermilab talks at HTCondor Week – 15 years ago!
  - My earliest HTCondor Week talk in 2012
- Describe the current state of the cluster(s)
- Along the way, I hope to:
  - Show some (maybe) unique uses of HTCondor
  - Explain why we did what we did
  - Give a peek into some future activities

# In the Beginning... (At least for me)

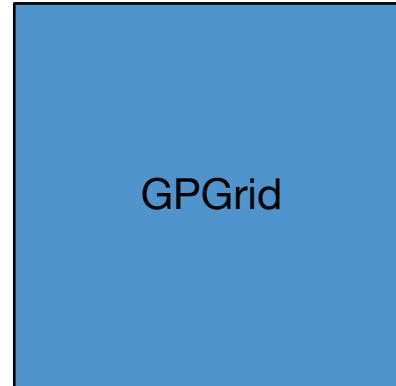
- There was HTCondor! And it was Good.
  - When I started, the silent “HT” hadn’t been added to the name yet



- Single VO
- Grid-enabled
- Priorities
- CMS + OSG



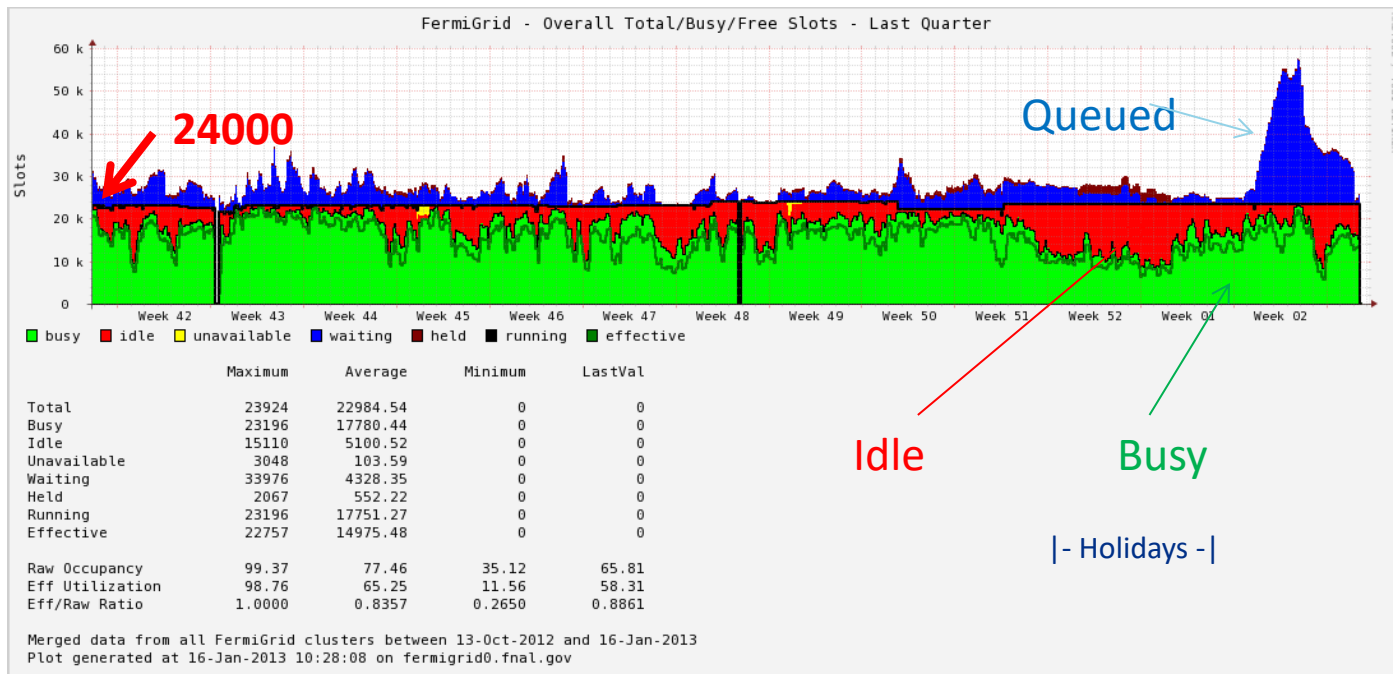
- Single VO Pool
- Local Analysis only
- Priority based scheduling



- Multi-VO Pool
- Grid-enabled
- Quotas
- Many experiments + OSG

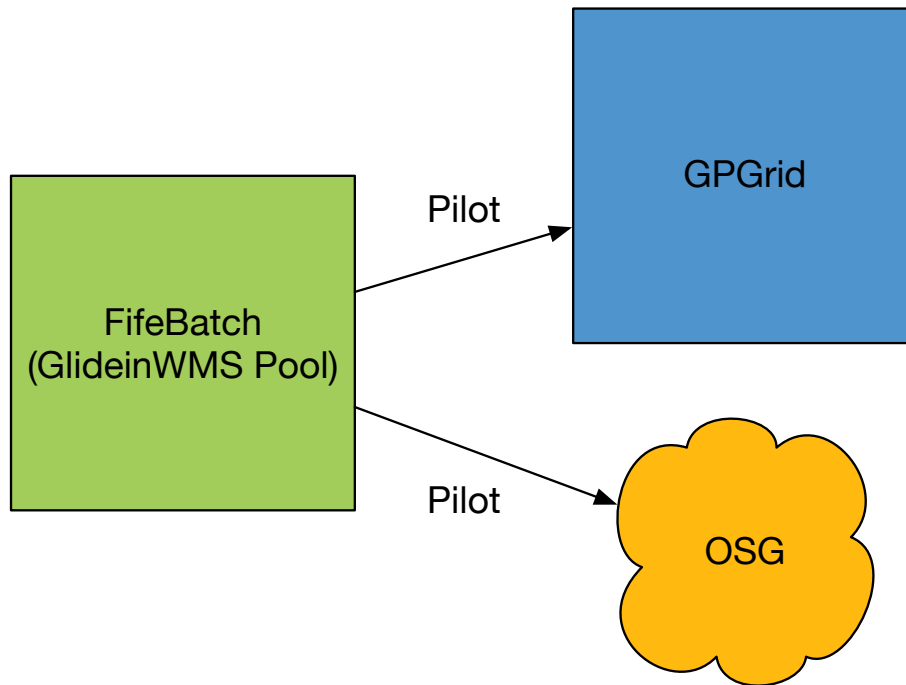
# Net Batch Slot Utilization – 2013 Scientific Computing Portfolio Review

Last 3 months



# FIFEBatch

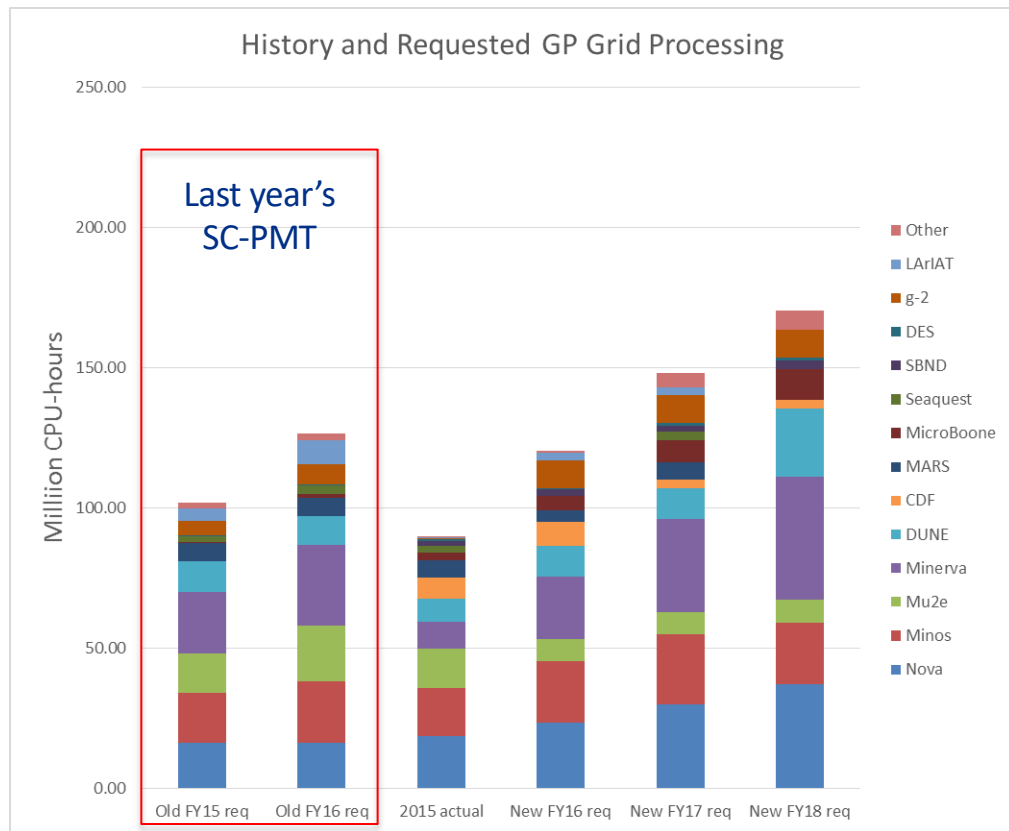
- FifeBatch was created using GlideinWMS
  - Main motivation was the desire to use OSG resources seamlessly.



# FIFEBatch

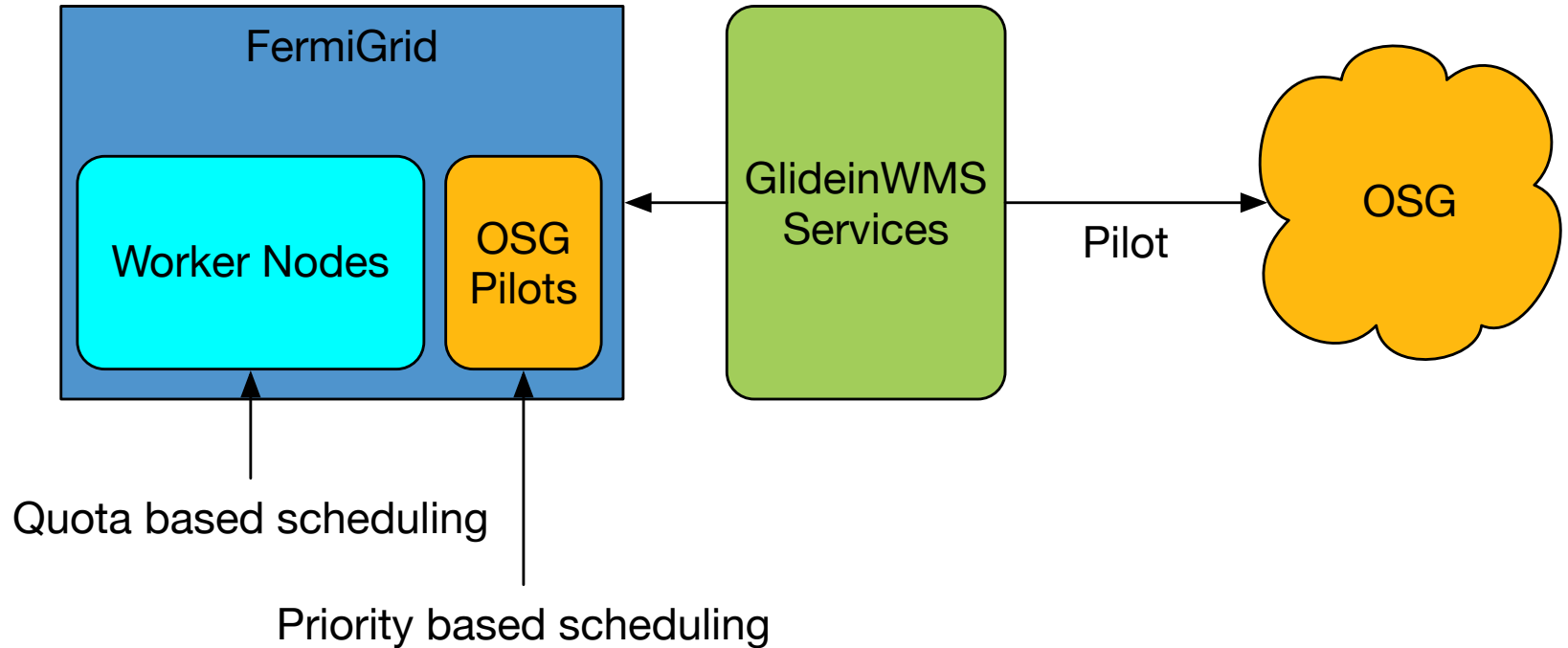
- FIFEBatch was a GlideinWMS pool
  - All slots are similar – controlled by pilot (glidein)
  - Used the glideinWMS Frontend to implement policies
  - Used the OSG Factory for pilot submission
  - Pilot “shape” defined by Factory
  - All of the benefits of glideinWMS and OSG
- All FNAL experiment jobs ran within the FifeBatch pool
- FIFEBatch managed by experimental support team
- GPGrid Managed by Grid Computing team

# SC-PMT - GP Grid Processing requests: Large memory or multi-core as single slot



- We began to see increased demand for large memory or multi-core slots
- For context:
  - A “standard” slot was defined as 1 core, 2GB RAM
- Partitionable slots limited by the pilot size
- Unable to use extra worker resources beyond what is claimed by the pilot

# Combined: GPGrid + FifeBatch = FermiGrid

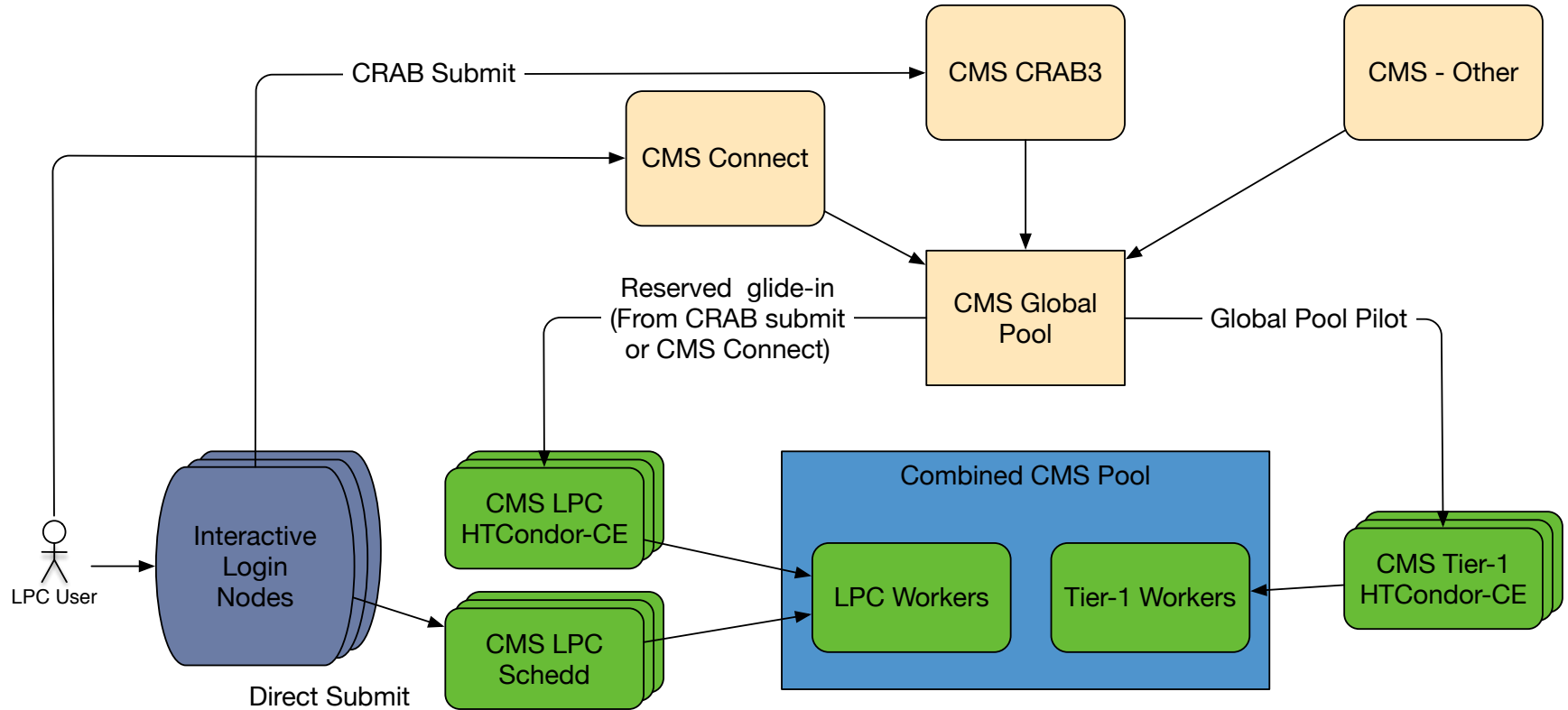




# CMS Tier-1 + LPC

- New requirements:
  - Make LPC available to CMS Connect
  - Make CRAB3 jobs run on LPC resources
- LPC workers reconfigured to remove all extra storage mounts
  - Now LPC workers look identical to the Tier-1 workers
- LPC needed Grid interface for CMS Connect and CRAB3
  - The Tier-1 was already Grid-enabled
- However, 2 competing usage models:
  - Tier-1 wants to be fully utilized
  - LPC wants resources at the time of need

# CMS Tier-1 + LPC



# CMS - Docker

## HTCondor-CE

### Job Router

Sets WantDocker = MachineAttrFERMIHTC\_DOCKER\_CAPABLE0  
Sets DockerImage = image expression

## HTCondor Worker

### Advertises:

FERMIHTC\_DOCKER\_CAPABLE=True  
FERMIHTC\_DOCKER\_TRUSTED\_IMAGES= <comma separated list>

## LPC Schedd

### Job Transform

Sets WantDocker = MachineAttrFERMIHTC\_DOCKER\_CAPABLE0  
Sets DockerImage = image expression

## GlideinWMS Pilot

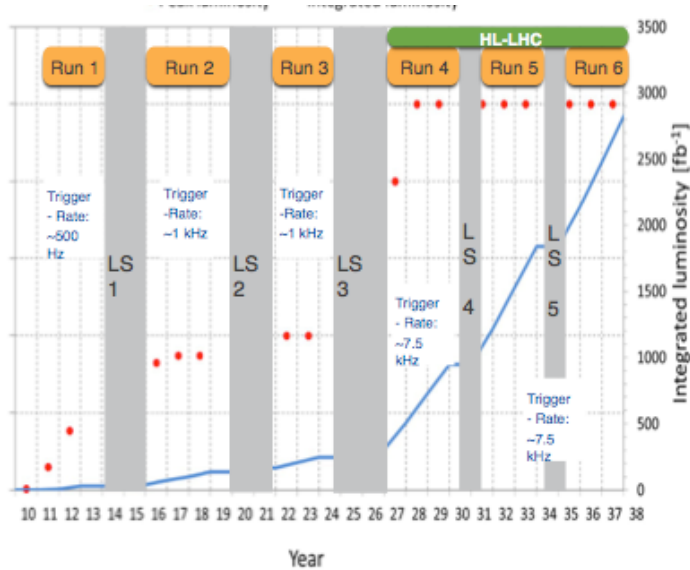
### Advertises:

FERMIHTC\_DOCKER\_CAPABLE=False

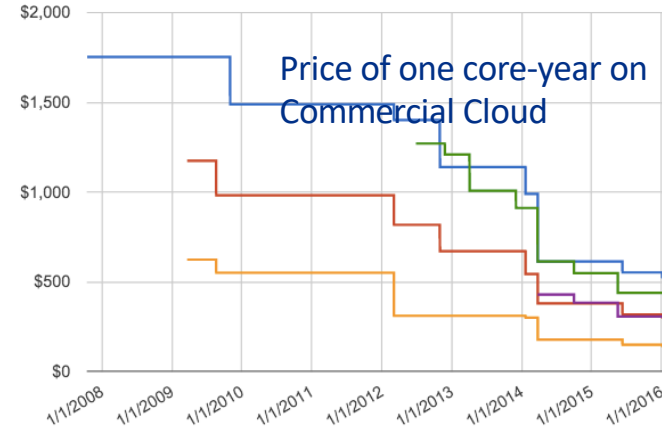
# HEPCloud - Drivers for Evolving the Facility

- HEP computing needs will be 10-100x current capacity

Two new programs coming online (DUNE, High-Luminosity LHC), while new physics search programs (Mu2e) will be operating

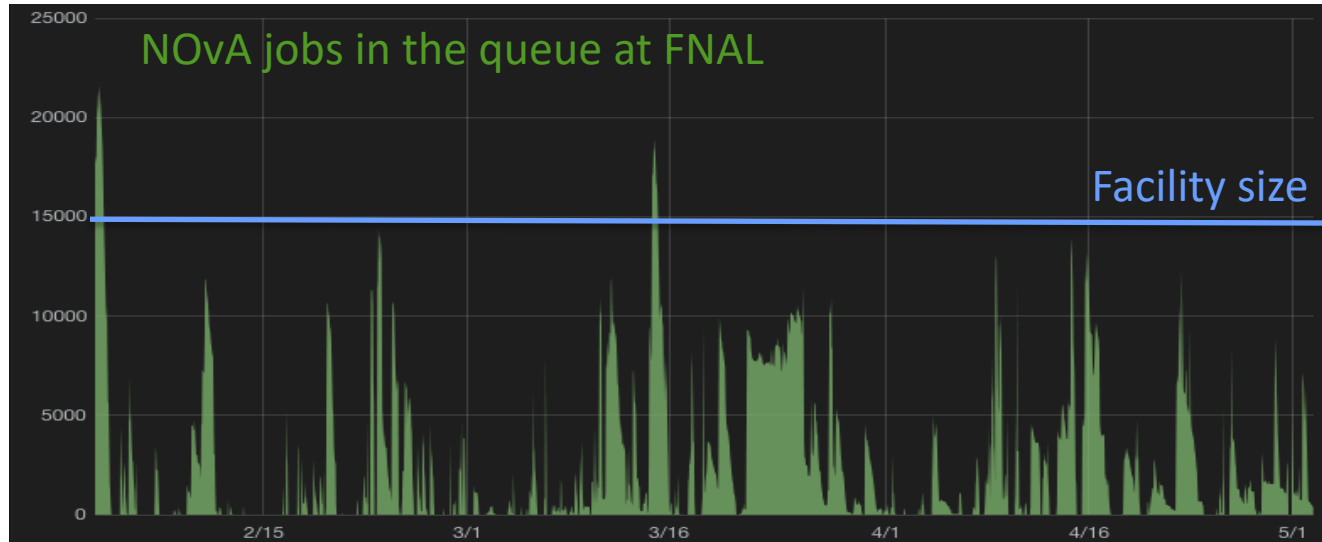


- Scale of industry at or above R&D
  - Commercial clouds offering increased **value** for decreased **cost** compared to the past



# HEPCloud - Drivers for Evolving the Facility: Elasticity

- Usage is not steady-state
- Computing schedules driven by real-world considerations (detector, accelerator, ...) but also ingenuity – this is research and development of cutting-edge science



# HEPCloud - Classes of Resource Providers

## Grid

- Virtual Organizations (VOs) of users trusted by Grid sites
- VOs get allocations → **Pledges**
  - Unused allocations: opportunistic resources

“Things you borrow”

Trust Federation

## Cloud

- Community Clouds - Similar trust federation to Grids
- Commercial Clouds - **Pay-As-You-Go** model
  - Strongly accounted
  - Near-infinite capacity → **Elasticity**
  - Spot price market

“Things you rent”

Economic Model

## HPC

- Researchers granted access to HPC installations
- Peer review committees award **Allocations**
  - Awards model designed for individual PIs rather than large collaborations

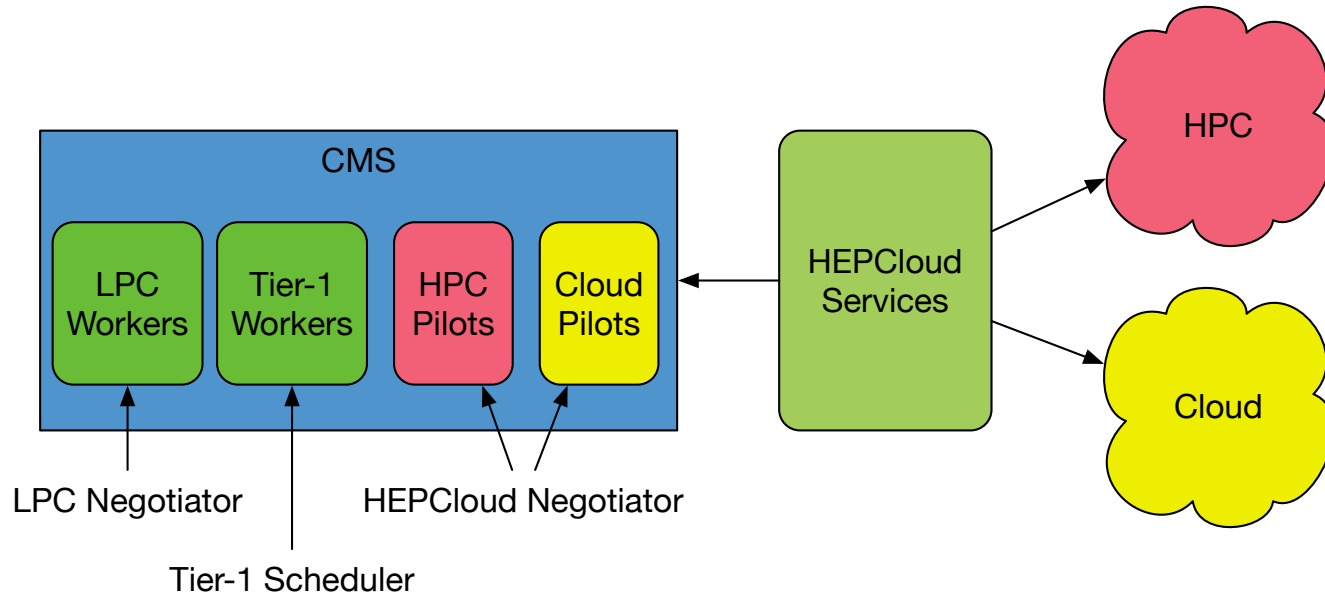
“Things you are given”

Grant Allocation

# HEPCloud

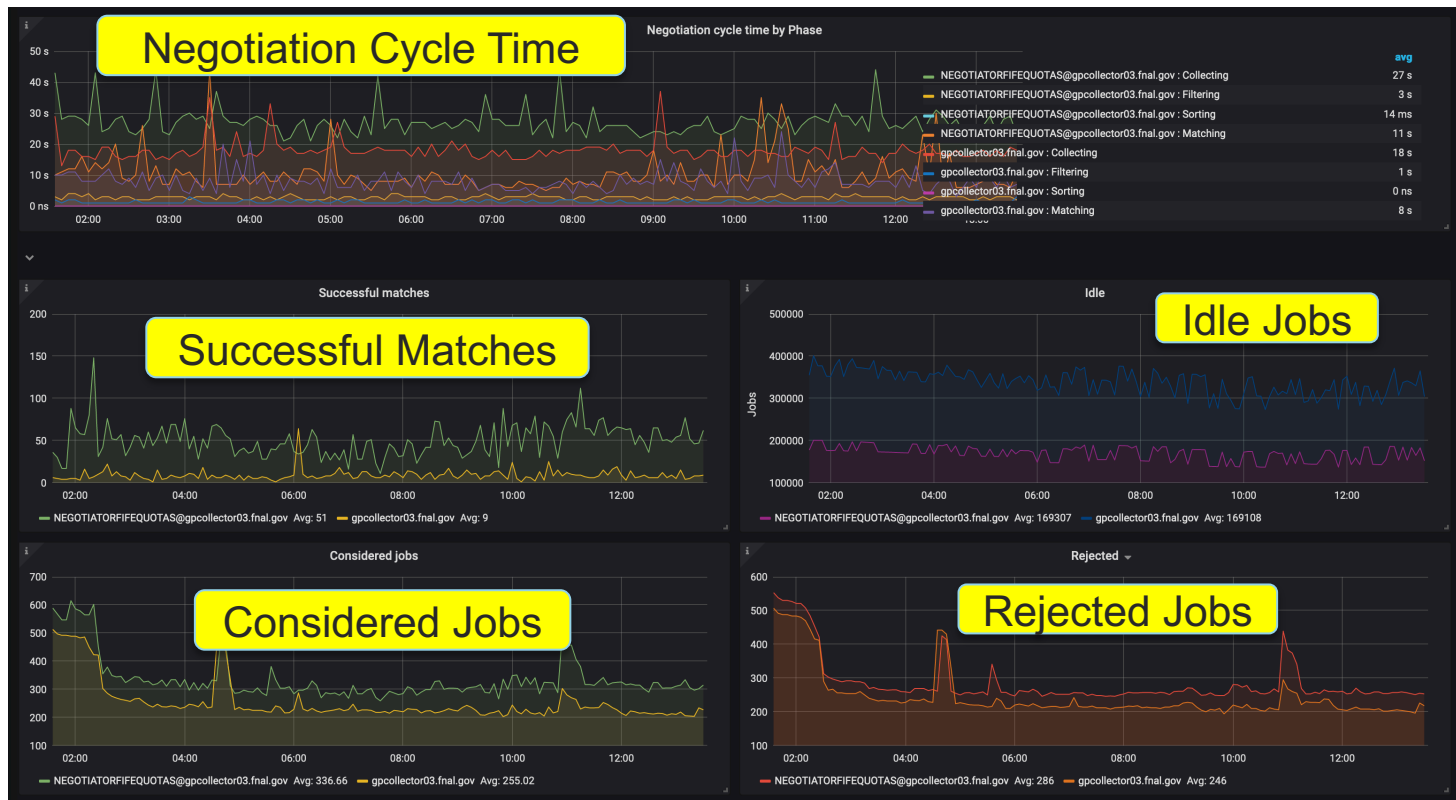
- New DOE requirements: Use LCF Facilities
- HEPCloud adds Cloud and HPC resources to the pool
- Cloud and HPC resource requests are carefully curated for specific classes of jobs
  - Only want appropriate jobs to land on Cloud and HPC resources
  - Additional negotiator also gives more flexibility in handling new resource types

# HEPCloud Era

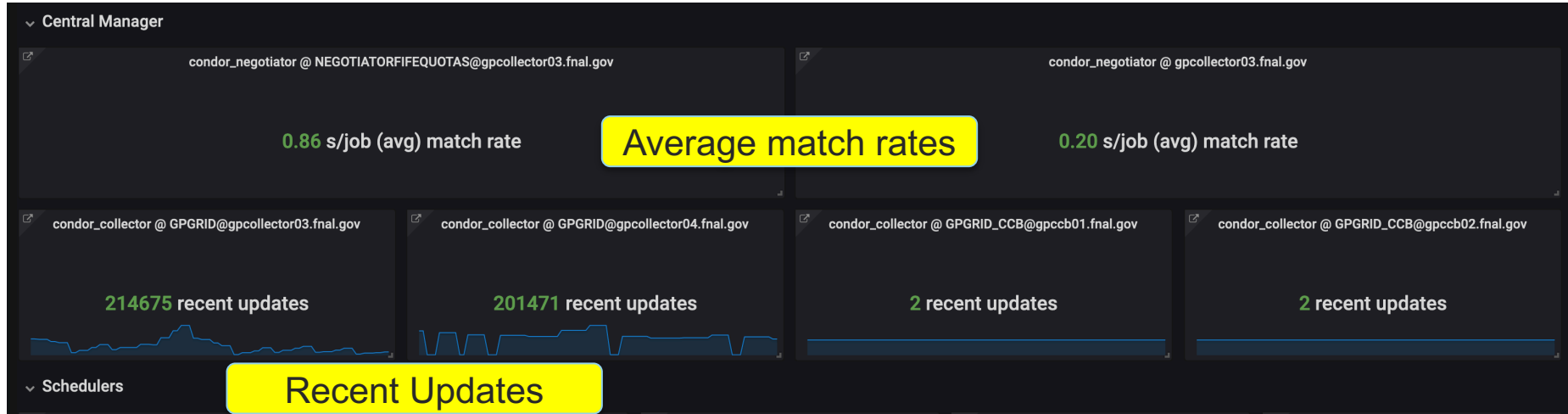




# Monitoring – Negotiation Cycles



# Monitoring – Central Manager



## Next Steps

- CI/CD pipelines for Docker containers
- Containerizing workers? (Kubernetes, DC/OS, etc.)
- HTCondor on HPC facilities with no outbound networking
- Better handling of MPI jobs
  - No dedicated FIFO scheduler
  - No preemption

# Questions, Comments?