

**CHEP** 2016  
Declination

22nd International Conference on Computing in High Energy and Nuclear Physics, Hosted by SLAC and LBNL, Fall 2016



# Stability and scalability of the CMS Global Pool: Pushing HTCondor and glideinWMS to new limits

**Antonio Pérez-Calero Yzquierdo**

on behalf of the **CMS Collaboration, Computing and Offline, Submission Infrastructure Group**

CHEP 2016 San Francisco, USA  
11th October, 2016

**Ciemat**  
Centro de Investigaciones  
Energéticas, Medioambientales  
y Tecnológicas



**PIC**  
port d'informació  
científica



# Highlights (I)



**CMS Global Pool:** single HTCondor pool covering all Grid computing processing resources pledged to CMS, plus significant Cloud and opportunistic CPUs. Resource allocation is handled by GlideinWMS while job to resource matchmaking is managed by HTCondor

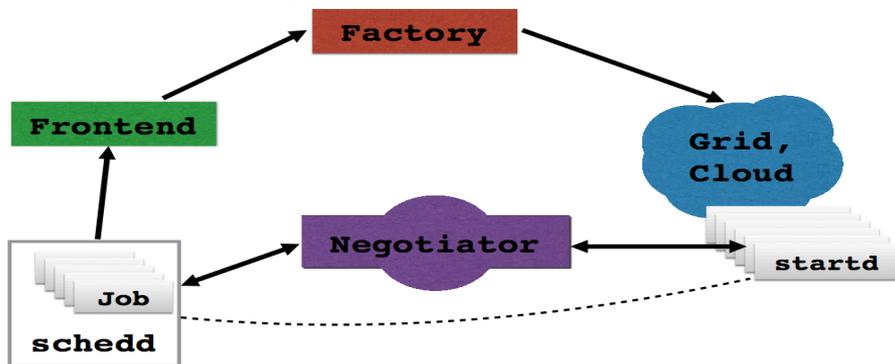
**Scalability:** The size of the pool has been increasing year by year (+40% 2015 into 2016), currently reaching peaks of 160k CPU cores. Main limits to scalability are the I/O between its components (schedds, startds, collector and negotiator), the combinatorics at the negotiator (jobs x pilots) and the speed of individual components

**Scalability tests with OSG** were conducted in 2015. The main recommendation was to put the HTCondor Communication Broker (CCB) on separate hardware from the rest of the Central Manager. Solving this I/O limitation allowed to push the scale up to 200,000 CPUs.

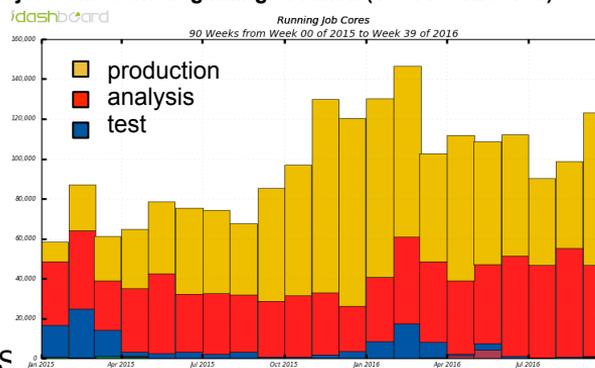
**Negotiator scalability** bottlenecks (I/O and combinatorics) were mitigated by running multiple parallel negotiators

**Transition to multi-core pilots and jobs** allows job memory requirement per CPU being reduced, as it also reduces the scale of the job/pilot matchmaking by the number of CPUs per job

Elements of the CMS global pool



Monthly avg. concurrently used CPU cores running CMS jobs since the beginning of Run II (T0+T1s+T2s+T3s)



scalability of the CMS



# Highlights (II)



**Further scalability limitations** observed during 2016 on the schedds and negotiator agents were mitigated with a number of interventions improving I/O between them and also making them faster.

**Next scalability tests:** a new round of scale tests with the OSG is planned for the fourth quarter of 2016. Interested in examining even higher scales but also include the effects of multi-core pilots and a more diverse job mix to model actual and future CMS usage. Look for scaling limitations and stability issues well above the current 150k CPU cores, in anticipation to 2017 and beyond

**Stability and High Availability** of the CMS global pool are achieved by means of the main agents (CM and FE) being deployed in HA mode while redundant infrastructure in different availability zones (EU and USA) is used for the other elements (pilot factories and job schedds)

Key to the success of the **Global Pool stably reaching ever higher scales** has been CMS's close coordination with the **HTCondor developers, the glideinWMS developers, and the OSG**

HA setup of the CMS global pool

