



Resource Provisioning su Cloud in CMS

Massimo Sgaravatto - INFN Padova



CMS HLT



- CMS High Level Trigger farm
 - Piu` di 13000 core
 - ~ 200 KHS06
 - Paragonabile a tutti i Tier1 di CMS messi insieme
- Risorsa da sfruttare anche per l'offline
 - Durante i Long Shutdown
 - Ma anche durante i run
 - Se la macchina e` in maintenance, se per l'online non servono tutte le risorse, etc.
 - Senza comunque interferire con le attivita` DAQ
- → Cloud per uso opportunistico di questa risorsa



Cloud CMS HLT



- Basata su OpenStack Grizzly
- Include al momento piu` di 4000 core
 - Link a 60 Gbps tra P5 (dove sta questa Cloud) e Meyrin
- Una VM per macchina fisica
- Macchine possono essere spostate da DAQ a OpenStack e viceversa
 - Attraverso un apposito tool



CERN Agile Infrastructure



- Private cloud del CERN
 - OpenStack based Cloud (Icehouse)
 - Multi-esperimento
 - Attualmente ~ 120K core in tutto (risorse fisicamente a Meyrin e Wigner)
- Attualmente CMS ha 3 progetti su questa cloud
 - CMS Tier0 on AI
 - Include al momento 6000 core (saranno 12000)
 - Risorse sia a Meyrn che a Wigner
 - CMS Evolution
 - 200 core
 - Primo progetto creato. Attualmente usato per test
 - CMS Wigner test
 - 400 core
 - Usato soprattutto per testare l'accesso allo storage (EOS) da Wigner
- Le VMs istanziate per CMS su AI hanno 4 core, 8 GB di RAM
 - Problemi con flavor piu` grandi (problemi a trovare hypervisor in grado di ospitarle)



CMS-HLT vs CERN-AI



- CERN-AI
 - Gestita da IT, multi-esperimento, CMS ha una quota su queste risorse
- Cloud CMS-HLT
 - Pienamente sotto il controllo di CMS, e CMS ne è l'unico utente
 - Qui sono possibili tuning/interventi che non sono possibili su CERN-AI
 - Es. pre-stage dell'immagine su tutti gli hypervisor, scratch del DB tra un run e un altro, ...
- Deployment di OpenStack fatto in modo diverso nelle 2 Cloud, e non sempre la cosa è trasparente per l'utente
 - Es. limite di 8 KB per user-data in CERN-AI, che ha richiesto una modifica in Condor

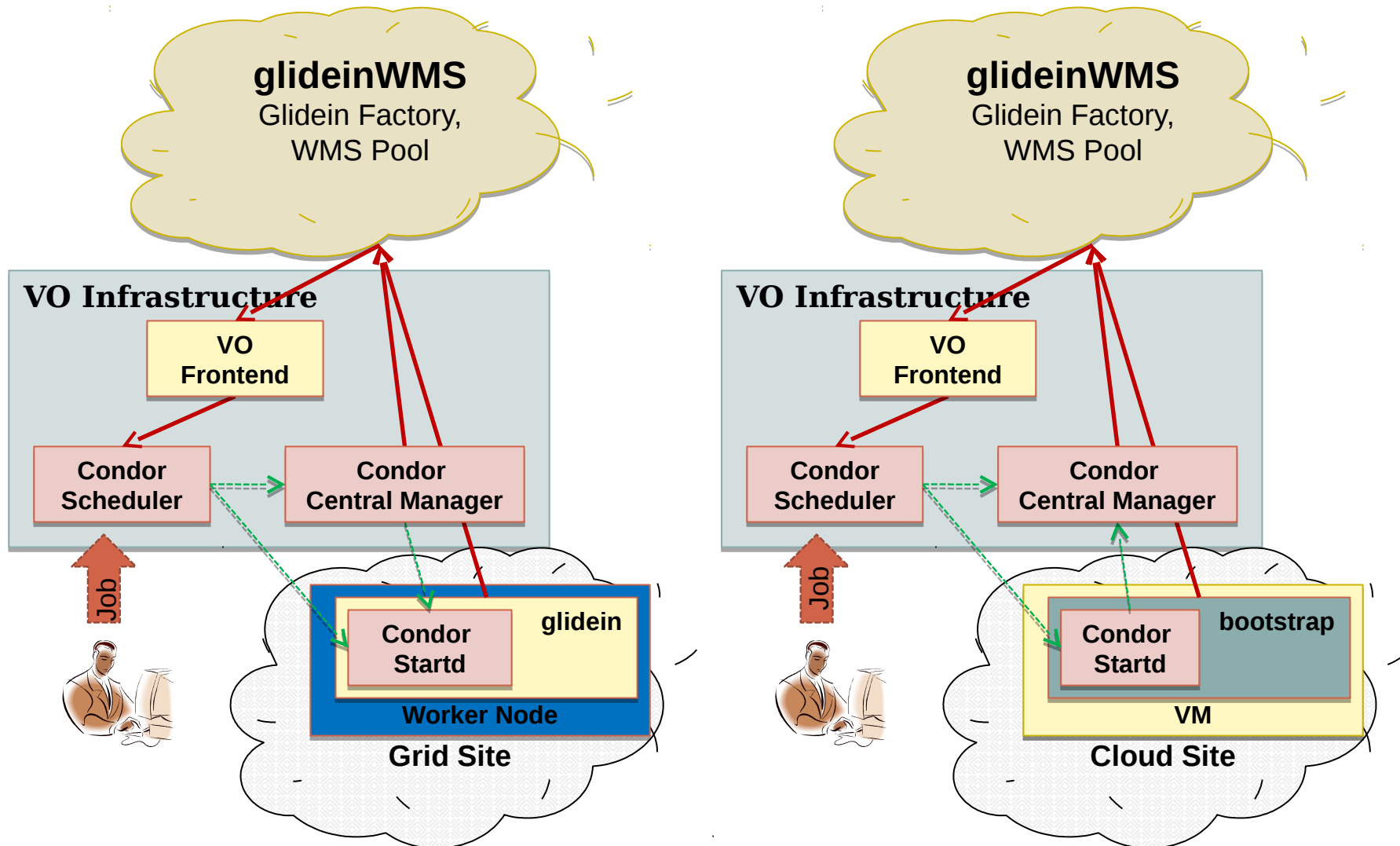


CMS-HLT vs CERN-AI (cont.ed)



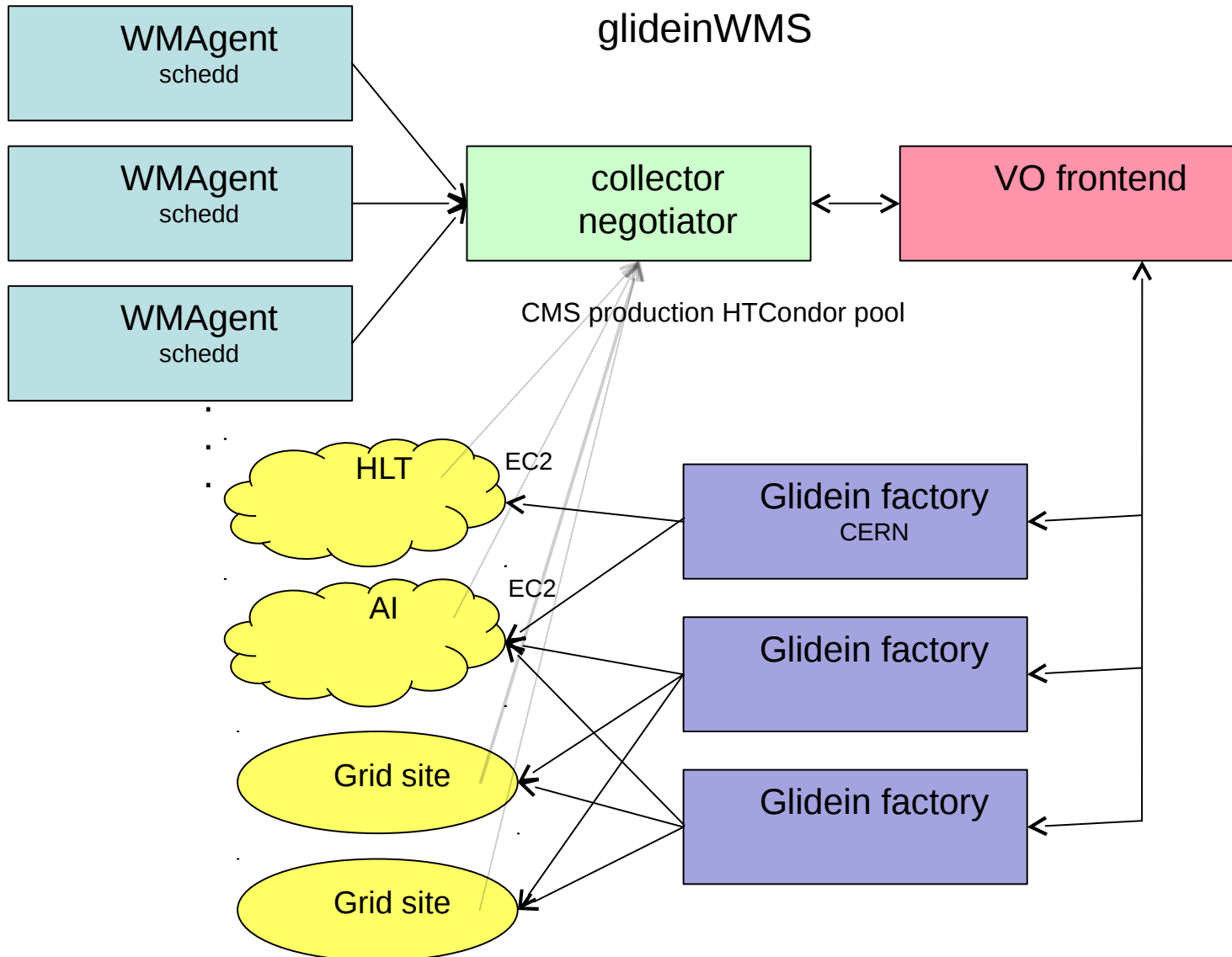
- Alcuni aspetti possono essere rilevanti per una cloud, meno per l'altra
 - Es. Start-up
 - Importante riuscire a “riempire” ASAP le risorse HLT non appena sono disponibili, visto che non saranno disponibili per sempre
 - Meno critico per CERN-AI dove le risorse sono dedicate e lo slow start puo` essere compensato da una lunga lifetime della VM (attualmente 1 mese)
 - Es. VMs in errore
 - Fondamentale fare ASAP il cleanup di queste VMs in CERN-AI, visto che sono “contate” come risorse in uso
 - Meno critico per CMS-HLT, dove la quota OpenStack e` infinita (non si riesce a istanziare nuove VM semplicemente quando non ci sono piu` hypervisor disponibili)

Job submission





Job submission (cont.ed)





In produzione !



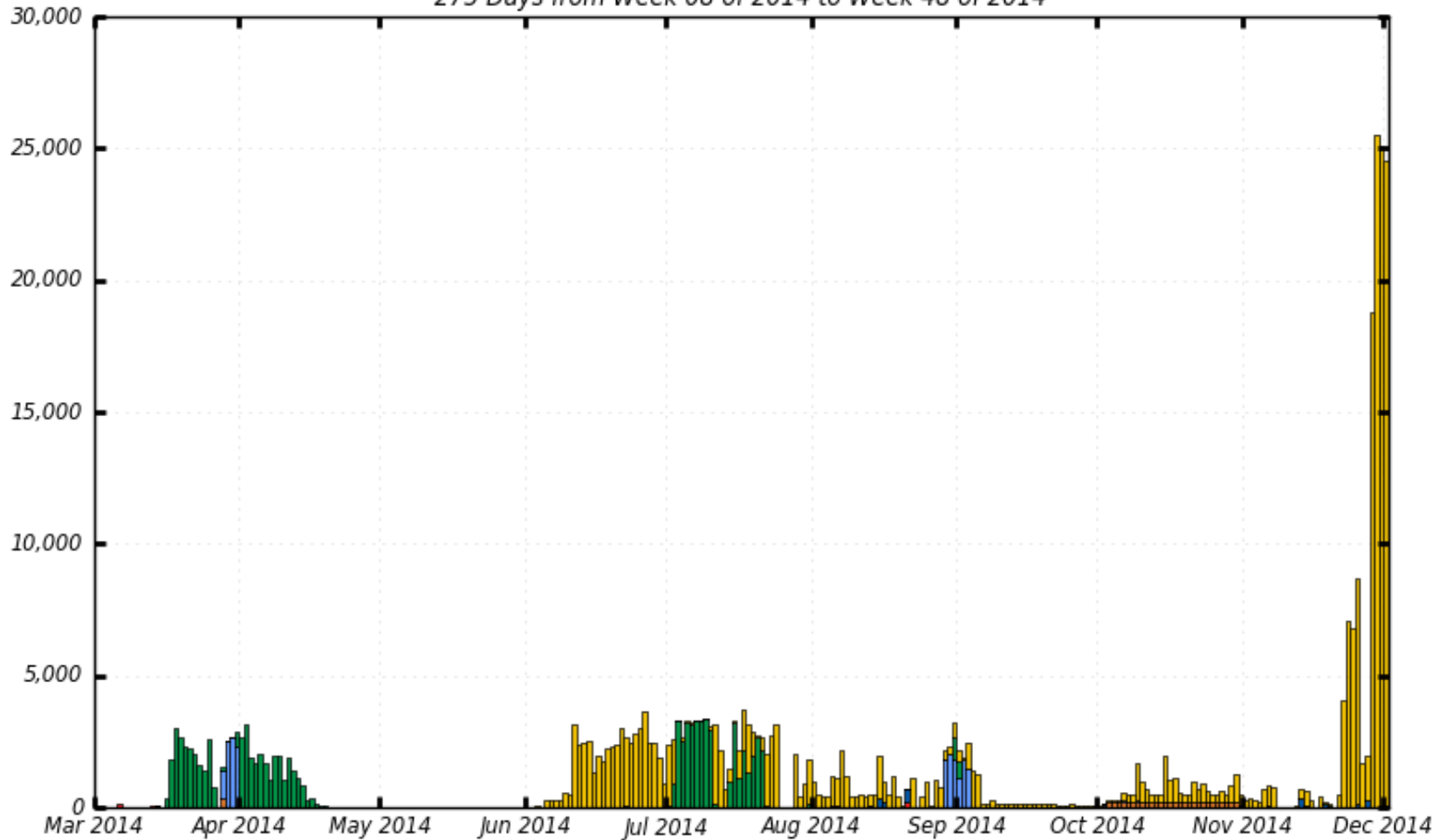
- **E` stato richiesto un intenso troubleshooting ...**
 - Sono stati necessari diversi fix/workaround/tuning in diverse componenti della job submission chain
 - Esempi di problemi visti (e risolti)
 - VMs in Error non cancellate
 - Leak delle ssh-keypairs
 - INFN ha contribuito usando una istanza di prova di GlideinWMS nel T2 di Padova-Legnaro che crea istanze su CERN-AI
- **... ma da Maggio le Cloud HLT e AI sono in produzione**
 - Usate per reprocessing e produzioni MC



Use Cloud HLT



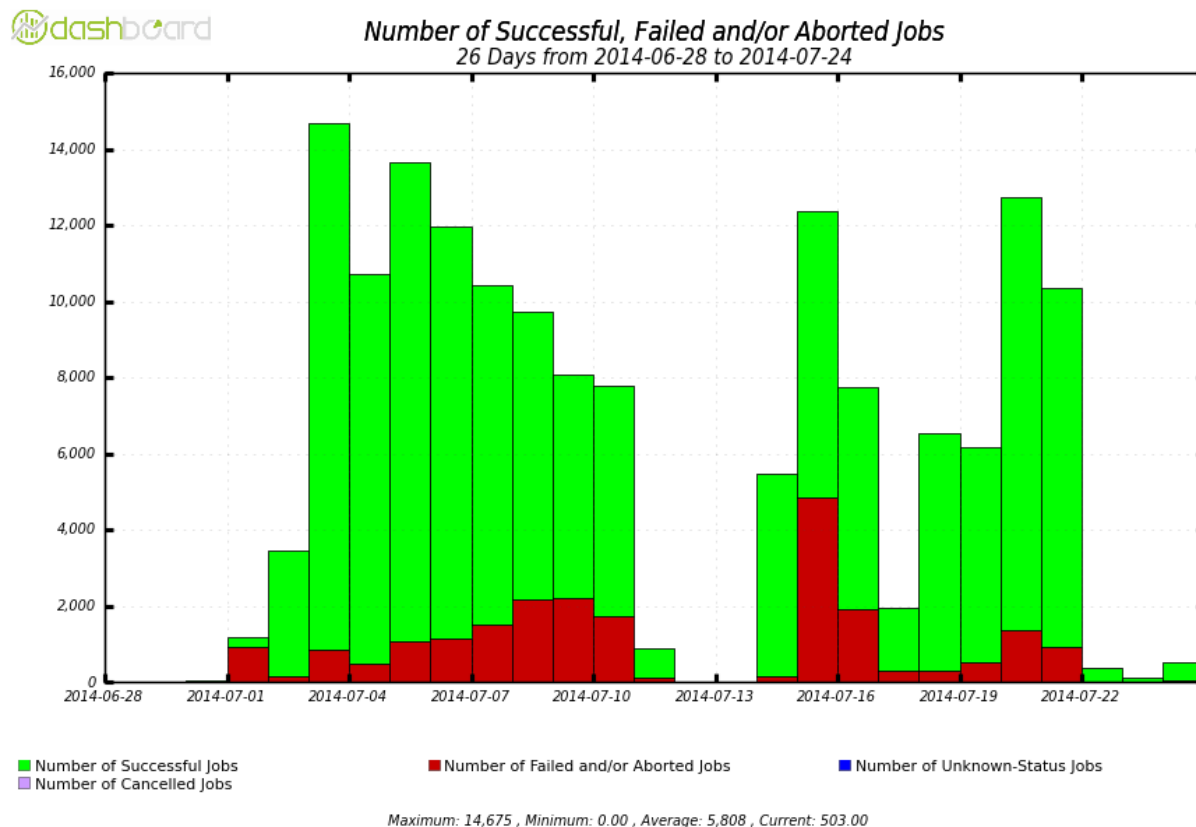
Running Jobs by Activities
275 Days from Week 08 of 2014 to Week 48 of 2014

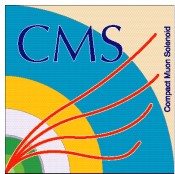


production reprocessing test unknown relval
cloud-testing integration

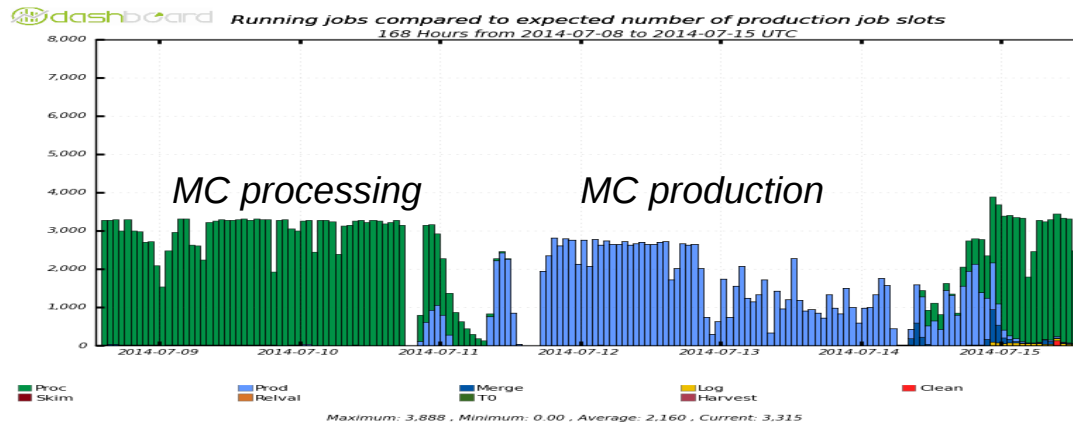
Maximum: 25,479 , Minimum: 0.00 , Average: 1,343 , Current: 24,557

- Success e failure durante CSA14 MC processing
- Failure dovuti a problemi di accesso allo storage

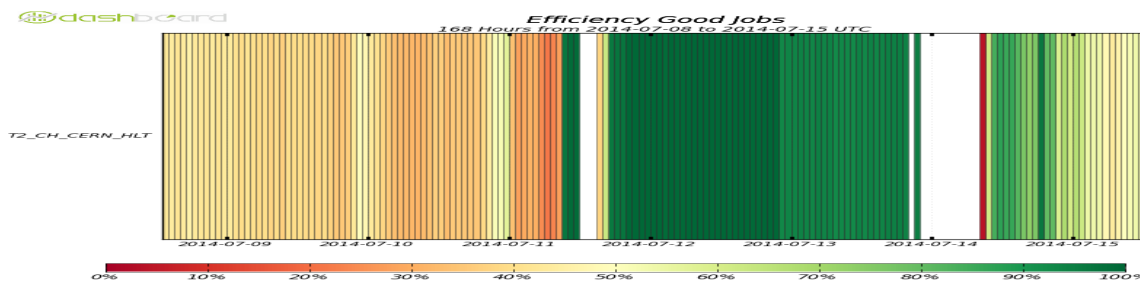




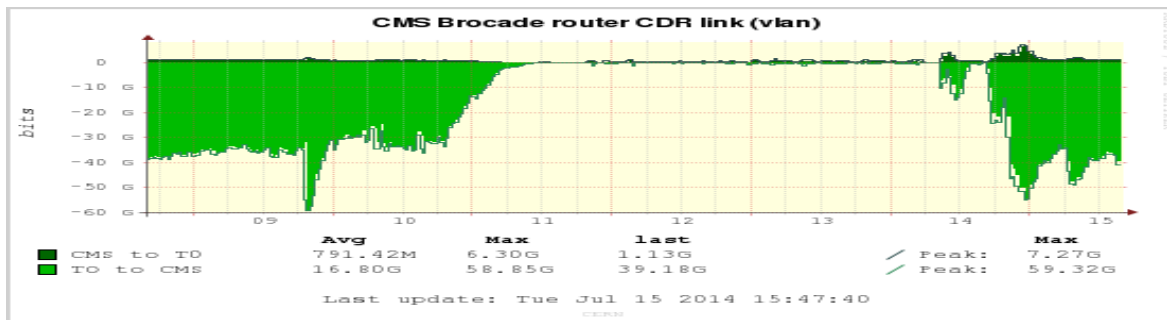
Cloud HLT efficienza (cont.ed)



activity



CPU efficiency

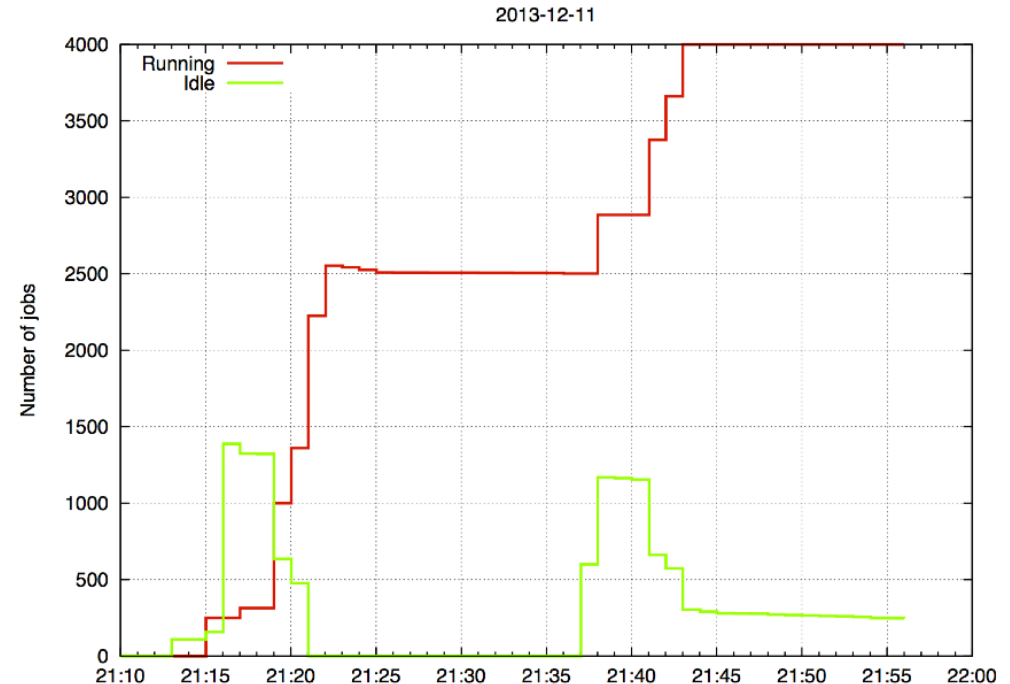


Network between P5 and CERN



Start-up

- HLT: partendo da scratch
~ 25' per avere 4000 job running
 - Istanziamento della VM, start di Condor startd (che deve contattare il collector), get e run dei job
- **Molto** piu` lento in CERN-AI
 - Attualmente viene creata una VM ogni 4'
 - Con rate alte molte VM vanno in Error
 - Si compensa considerando VM che durano 1 mese





Gestione immagini



- Non si usa uCernVM
- Immagini SL6.x create ad-hoc
 - CVMFS, middleware gLite per WN, pilot launcher, xrootd client, Ganglia gmond
- Per T0
 - Costruite via OZ
 - Integrazione con automatic build system implementato da CERN-IT
 - Sorgenti in GIT CERN: quando si committa qualcosa, viene triggerato un build dell'immagine che viene caricata in Glance



Cloud per 2011 HI rereco



- Processing necessario per la conferenza Quark Matter di Maggio 2014
 - 2 workflow: uno assegnato a HLT, l'altro a Vanderbilt
- A fine Marzo 2014 realizzato che non si sarebbe riusciti a finire in tempo a causa di problemi con il workflow assegnato a Vanderbilt
- Deciso allora di runnare
 - un clone del workflow su piu` siti con accesso a dati via xrootd
 - Abortito a causa di troppi failure
 - un clone del workflow su Cloud AI+HLT
 - Completato con successo in tempo per la conferenza



Alcuni next step



- Cloud usate anche per l'analisi
- Aggiunta di altri siti Cloud, in un'ottica di migrazione dei siti Grid **in maniera efficiente**
 - CMS interessato a testare il fairshare scheduler
- Non far istanziare le VM dalla GlideinWMS factory ma via Heat, in particolare per HLT
 - Per facilitare l'aggiunta/rimozione di risorse, e per risolvere alcuni problemi nel job submission framework attuale (es. API EC2 fragile, necessita` di una ssh-keypair per VM)
 - L'immagine installerà e configurerà Condor in modo da connettersi al GlideinWMS user collector
- Application checkpointing e restart
 - Per lo use case HLT, dove è necessario fare il preempting delle risorse quando servono per DAQ



Conclusioni



- Cloud HLT e CERN-AI utilizzate per attività di vera produzione
 - Sia per processing che per produzioni MC
- Contributo INFN
 - C. Grandi
 - Co-coordinamento (con D. Collings)
 - M. Sgaravatto
 - Contributo nel troubleshooting della job submission chain, creazione e update immagini per T0, operations

Grazie a Andrew Lahiff (RAL) a cui ho rubato molto del materiale che ho usato in questa presentazione