HTCondor at HEPiX, WLCG and CERN – Status and Outlook

Helge Meinhard / CERN HTCondor week 2018 Madison (WI) 22 May 2018

CERN material courtesy by Ben Jones



HEPiX

• From our Web site https://www.hepix.org:

"The HEPiX forum brings together worldwide Information Technology staff, including system administrators, system engineers, and managers from the High Energy Physics and Nuclear Physics laboratories and institutes, to **foster a learning and sharing experience** between sites facing scientific computing and data challenges."

- Workshops: Twice per year, one week each
 - Open attendance, everybody (including non-HEP!) welcome
 - Plenaries only, no programme committee, no proceedings
 - Honest exchanges about experience, status and plans
 - Workshop last week at Physics Department, UW Madison
 - Next workshops: 08 12 October 2018 Barcelona (Spain); spring 2019 San Diego (CA); autumn/fall 2019 Amsterdam (The Netherlands)
- Working groups, board, co-chairs (HM, Tony Wong/BNL)



HTCondor at HEPiX (and WLCG)

- HTCondor often mentioned at HEPiX in site reports and dedicated presentations (computing track)
- Clear consolidation: Previously plethora of solutions (PBS/Torque, *GridEngine, LSF, ...), most sites now on (or moving to) HTCondor or (HPC only) Slurm
- Similarly: CEs for grid submission: Consolidating on HTCondor CE (with HTCondor) and ARC-CE (with HTCondor and Slurm)
- Big topic recently: analysis job submission from Jupyter notebooks
- WLCG in December 2017 at pledging sites: 211M HS06 days (30% over pledges), equivalent to 700k average today's cores
 - Significant contributions from non-pledging sites, volunteers, ... ("opportunistic usage")



HTCondor in WLCG

Site	Batch scheduler	Site	Batch scheduler
CERN	See later	US T2	Mostly HTCondor
BNL	HTCondor	LBNL	Slurm
FNAL	HTCondor	IHEP	HTCondor, (Slurm)
KIT	HTCondor	DESY	HTCondor, (Slurm)
Nordic T1	Slurm	FZU	Migration to HTCondor
CC-IN2P3	UGE, considering HTC		ongoing
	LITC and an	U Tokyo	LSF
RAL	HICondor	CSCS	Slurm
Nikhef	PBS	0000	Oldini
PIC	Migration to HTC 60% done	GRIF	HTCondor
		CoEPP	HTCondor
CNAF	Migration to HTC started		



CERN: Previously at HTCondor week...

- At the 2016 HTCondor week, we had a production setup
- Since then we have increased in size, and also the scope of what we're asking the batch system to do
- The rest of this talk will cover where we are with our deployment, the evolution of our use cases, and some future work



Batch Capacity





Last 2 years (on fifemon since 2016 Condor Week)





Migration status

- Grid workload migrated entirely
- No technical issues preventing rest of capacity moving to HTCondor
- Remaining use cases are some Tier-0 reconstruction & calibration that will move at end of Run 2 (end 2018)





CERN Data Centre: Private Openstack Cloud



Two submission use cases

	Grid	Local
Authentication	X509 Proxy	Kerberos
Submitters	LHC experiments, COMPASS, NA62, ILC, DUNE	Local users of experiments, Beams, Theorists, AMS, ATLAS Tier-0
Submission method	Submission frameworks: GlideinWMS, Dirac, PanDA, AliEn	From condor_submit by hand, to complicated DAGs, to Tier-0 submit frameworks.
Storage	Grid protocols. SRM, XRootD…	AFS, EOS



Compute Growth Outlook

- Resources looking very tight for Run 3
- No new datacenter & exiting Wigner
- Requirement to maximize the use of any compute we can, wherever it is acquired.





HTCondor Infra in numbers

- 2 pools
 - Share + extras: 155k cores
 - Tier-0 (CMS and ATLAS): 30k cores
- 13 + 2 production htcondor-ce
- 10 + 1 production "local" schedds
- Main shared pool:
 - 3 negotiators (2 normal + 1 for external cloud resources)
 - 15 sub collectors
- Max 10k jobs per schedd



Multiple resource types

- Standard shared batch farm
- Resources dedicated to one group
 - Special requirements, such as Tier-0 activity
 - Experiments that "own" their own resources, but want central IT service to run it
- Opportunistic resources internally
 - Using spare CPU slots on Disk servers (BEER)
- Opportunistic resources externally
 - XBatch / HNScience Cloud
- Special machines (big memory)



Targeting specific resources

- Beyond specifying just resource characteristics (cpu, memory etc) we have jobs targeting different resources
- Accounting Group matches jobs to dedicated resources
- We use job router / job transforms to provide special routes to special resources like Cloud or BEER
 - Experiments' monitoring is based on concept of "sites" with particular JDL, and for special resources they want extra observability



BEER



- Batch on EOS Extra Resources
- CERN has typically bought same hardware for batch and disk servers
- Disk servers don't use much CPU (or, for physics workload) utilize much filesystem cache
- Familiar to any of you that were at HEPiX last week – see HEPiX talk for performance analysis

https://indico.cern.ch/event/676324/contributions/2981816/



BEER Integration

- Aim: limit HTCondor & jobs to under resource limits disk server can afford
- Minimize config & OS requirement of host disk server
- HTCondor and jobs managed by CGroup with max memory, limit CPUs and I/O
- Jobs in Docker universe to abstract disk server environment
- Drain / evacuate procedures for storage admins!











- Procurement so far has been for flat capacity rather than burst
- HTCondor integration 1.0:
 - Configuration Management to create VMs with certificates to log into pool
 - Experiments again want to observe / monitor as a separate site
 - Separate negotiator, specific htcondor-ce route to match jobs requesting cloud with cloud workers



Future: kubernetes



- Kubernetes to manage HTCondor has a number of potential wins
 - kubefed federation means we can span kubernetes pods across clouds
 - At kubecon demoed federation from CERN to T-Systems, have integrated GCE, Azure, AWS, CERN... https://kccnceu18.sched.com/event/Duoa
 - Simplify requirements for cloud: just need a container engine or just IAAS
 - Internally can use bare metal managed by cloud team, container layer batch team
 - No "virtualization overhead", no hypervisor tax
 - Potential to "hyperconverge" data, services, batch





kubefed init fed --host-cluster-context=condor-host



22 May 2018

...

HTCondor at HEPiX, WLCG and CERN 21



--host-cluster-context condor-host \

--cluster-context tsystems

22 May 2018

Conclusions

- Demands on compute not getting easier
- We need to be able to deploy real workload on any resources we can get our hands on
- HTCondor continues to help us expand and meet these demands
- More technical detail available at European HTCondor workshop 04-07 September 2018 at RAL (and next HTCondor week hopefully)



Questions?