

Developments for ATLAS Supercomputing and Production system

R. Mashinistov et al



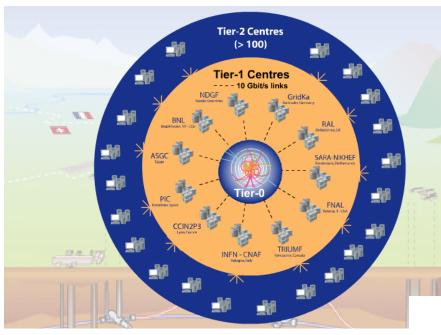


Main topics

- Introduction
 - WLCG
 - PanDA
 - Production System
- Supercomputing
 - ATLAS Pledged Resources
 - Supercomputing in HEP
 - HPC@NRC KI
- Production System
 - •New Task interface
- Plans



WLCG - Worldwide LHC Computing Grid



- >ATLAS Computing Model:
- >11 Clouds: 10 Tis + 1 TO (CERN)
 Cloud = T1 + T2s + T2Ds
 T2D = multi-cloud T2 sites
- > 2-16 T2s in each Cloud

PanDA

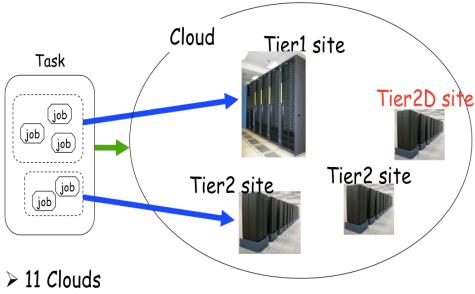
Workload Management System

Task → Cloud : Task brokerage

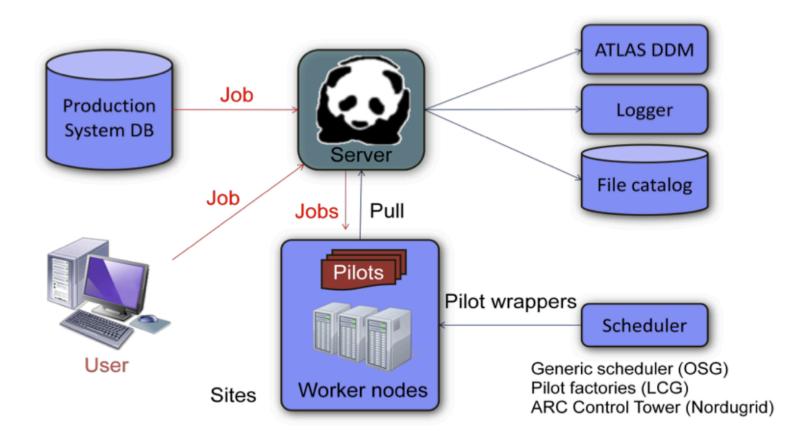
Jobs → Sites: Job brokerage

- Basic unit of work is a job:
 - Executed on a CPU resource/slot
 - May have inputs
 - Produces outputs
- JEDI layer above PanDA to create jobs from ATLAS physics and analysis 'tasks'

Current scale – one million jobs per day



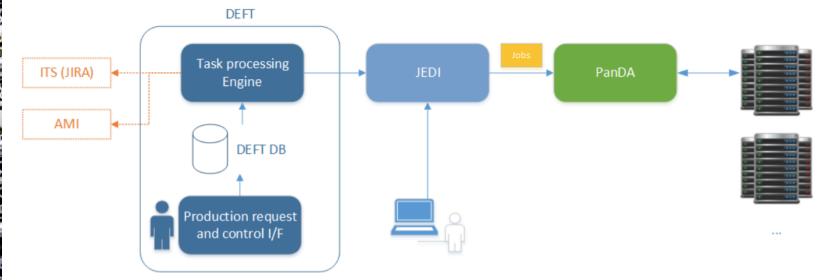
PanDA Workload Management





ATLAS Production System

 Prodsys2 - DEFT (Database Engine for Tasks) and JEDI (Job Execution and Definition Interface)





ATLAS pledged resources

Totally Tiers provide

Physical CPU	Logical CPU	HEPSPEC06	CPU Pledge					
179,625 532,910		5,716,182	3,083,512					
			Tape Pledge, Gb					
Total Online Storage,Gb	Disk Pledge, Gb	Total Nearline Storage, Gb	Tape Pledge, Gb					

Dedicated to ATLAS:

	CPU, HEP- SPEC06	Disk, Tbytes	Tape, Tbytes				
Required	1175000	103000	98000				
Provided	1275226	110293	103190				
% from total	22	34	43				

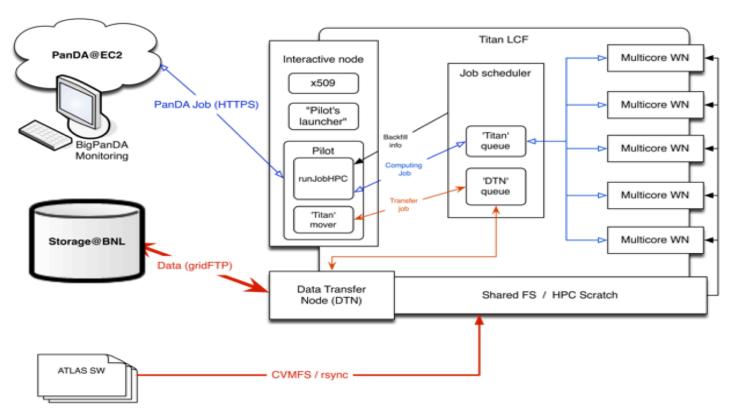
http://wlcg-rebus.cern.ch/apps/pledges/resources/



Supercomputing in High Energy Physics. Leadership Computing Facilities

- ➤ Modern High Performance Computing involves a very large number of cores connected through a high speed network
- ➤ It is not a requirement that HPCs are able to run any possible task
- ➤ What matters is the total number of cycles that can be offloaded from the Grid
- Expanding PanDA from Grid to Leadership Class Facilities (LCF) required changes in the system
 - ➤ Each LCF is unique
 - Unique architecture and hardware
 - ➤ Specialized OS, "weak" worker nodes, limited memory per WN
 - Code cross-compilation is typically required
 - Unique job submission systems
 - ➤ Unique security environment
 - ➤ Pilot submission to a worker node is typically not feasible
 - Pilot/agent per supercomputer or queue model

Extending PanDA to Oak Ridge Leadership Computing Facilities



- ATLAS (BNL, UTA), OLCF, ALICE (CERN, LBNL, UTK) :
- SAGA (a Simple API for Grid Applications) framework as a local batch interface.
- Pilot (payload submission) is running on HPC interactive node and communicating with local batch scheduler to manage jobs on Titan.
- MPI wrapper/overlay scripts that allow to run multiple single threaded workload instances in parallel
- "Backfill" functionality in pilot
- Outputs are transferred to BNL T1 or to local storage



HPC2 supercomputer at NRC KI





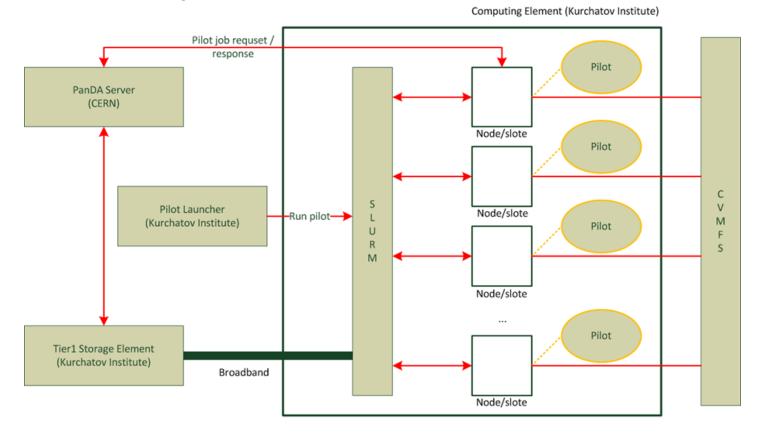
High Performance Cluster - HPC2 second generation HPC with peak performance 122,9 TFLOPS (commissioned 2011). #2 in 15-th issue of Russian top50 Supercomputers

- ◆10240 CPU cores = 1280 nodes 2x Intel Xeon E5450 3,00GHz 4 core 16 Gb RAM;
- UI node only allows to run jobs in batch system (SLURM) or to compile the code
- Shared FS Lustre for WN's and UI
- WN's has an access to WAN
- CVMFS connected to WN's
- Broadband to Tier-1 Storage Element



Integration of PanDA with HPC2 supercomputer at NRC KI

- Pilots are running on each node
- •Number of Pilots = number of available cores (and number of virtual CondorLocal slots)
- Remote running via interactive node (UI)

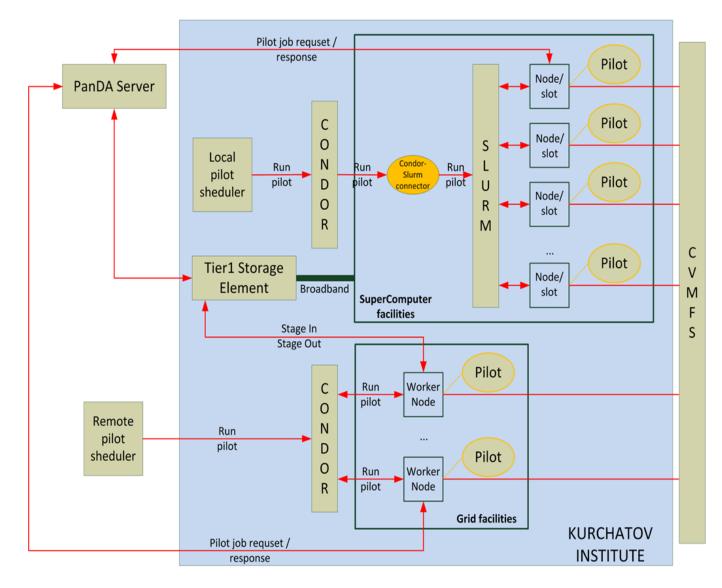




Integration of Russian Tier-1 Grid Center with High Performance Computers at NRC-KI (CHEP15)

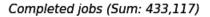
- In 2014 the pioneering work on implementation of the portal, which will combine the Kurchatov Institute's Tier-1 site and the supercomputer HPC2 into Kurchatov Institute's Tier-1 center was done.
 - This center allows starting tasks optionally not only on Grid but on the supercomputer as well and to collect the results in the common storage.
- For ATLAS production and analysis tasks new PanDA sites within ATLAS Tier-1 site of National Research Center "Kurchatov Institute" were defined for supercomputer.
 - ANALY_RRC-KI-HPC site for user analysis 2 nodes (16 cores, 2GB RAM per core)
 - User analysis tasks for MEPhI. RAW to D3PD reconstruction of high mu pp events (up to 70 interactions) for TRT performance at high occupancy study. (ATLAS TRT SW group)
 - Results presented at <u>Nuclear Electronics & Computing</u>, <u>NEC-2015</u>
 (28 Sept- 3 Oct 2015, Montenegro)
 - https://cds.cern.ch/record/2007680
 - RRC-KI-HPC2 site for production 32 nodes (256 cores, 2 GB RAM + 1 GB SWAP per core)

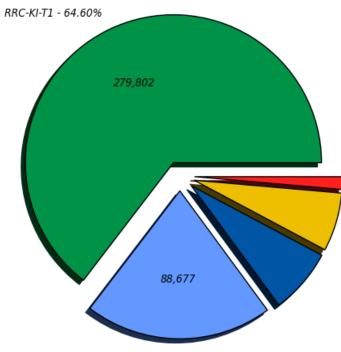
Integration scheme of Russian Tier-1 Grid Center with High Performance Computers at NRC-KI (CHEP15)



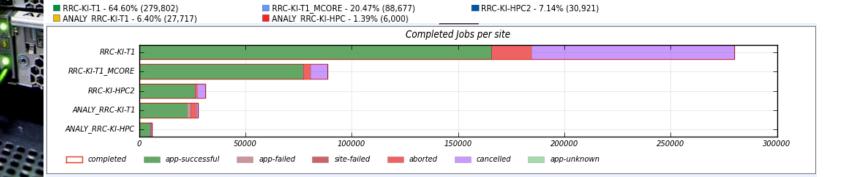
ATLAS site validation







RRC-KI-T1_MCORE - 20.47%



ATLAS site validation

IRZ-LMU_CZPAP_MCORE

AVALY_RRC-KI-HPC

AVALY_HPCZN

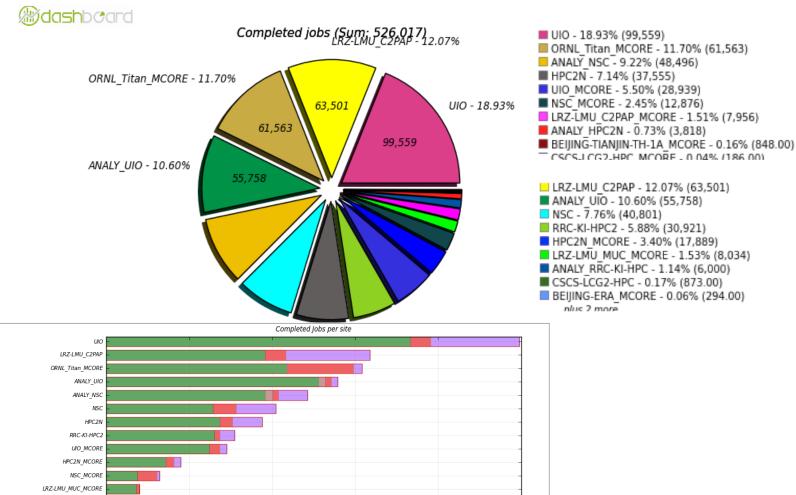
CSCS-LCG2-HPC

BEIJING-TIANJIN-TH-1A_MCO..

BEIJING-ERA_MCORE

CSCS-LCG2-HPC_MCORE

MPPMU-HYDRA_MCORE





Prodsys2: Production task monitor

- Task Processing engine: translates production requests to tasks
- ☐ Task monitoring: filtering, sorting and etc.
- Task manipulation: abort, clone, finish, change priority, reassign and etc.



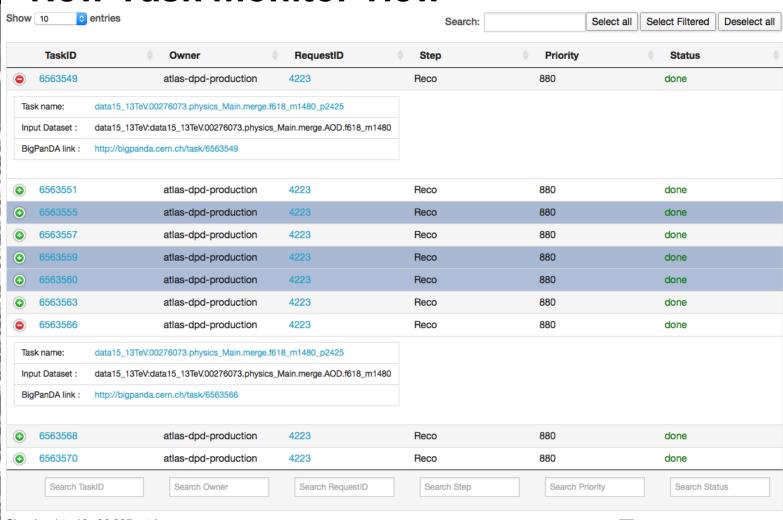
sk Name C		Step 0	 Task ID 	⇒ Priority ⇒	Total Jobs	Done Jobs ≎	Failure %		Status 0	Submit time	Timestamp	Provenance	٥ ٥.
mc15_13TeV.363075.Sherpa_CT10_llvv_renorm025.merge.e4681_p2482	5580	Reco	7409843	560	0	0	None	6000000	waiting	Jan 11 13:23	Jan 11 13:23	GP	587
mc15_13TeV.301171.PowhegPythia8EvtGen_AZNLOCTEQ6L1_Wmintaunu_2250M2500.merge.e3663_s2608_s2183_r7326_r6282_p2495	5577	Reco	7408485	560	3	4	0	18770	done	Jan 11 10:03	Jan 11 13:3	3 GP	58
mc15_13TeV.301157.PowhegPythia8EvtGen_AZNLOCTEQ6L1_Wplustaunu_4500M5000.merge.e3663_s2608_s2183_r7326_r6282_p2498	5577	Reco	7408482	560	3	4	0	18584	done	Jan 11 10:03	Jan 11 12:3	3 GP	5
mc15_13TeV.301147.PowhegPythia8EvtGen_AZNLOCTEQ6L1_Wplustaunu_1250M1500.merge.e3663_s2608_s2183_r7326_r6282_p2498	5577	Reco	7408479	560	3	0	0	0	running	Jan 11 10:02	Jan 11 11:3	3 GP	5
mc15_13TeV.301145.PowhegPythia8EvtGen_AZNLOCTEQ6L1_Wplustaunu_800M1000.merge.e3663_s2608_s2183_r7326_r6282_p2495	5577	Reco	7408475	560	3	4	0	15746	done	Jan 11 10:02	Jan 11 12:3	3 GP	5
mc15_13TeV.301142.PowhegPythia8EvtGen_AZNLOCTEQ6L1_Wplustaunu_250M400.merge.e3663_s2608_s2183_r7326_r6282_p2495	5577	Reco	7408471	560	13	10	0	13805	running	Jan 11 10:02	Jan 11 12:3	3 GP	5
mc15_13TeV.301163.PowhegPythia8EvtGen_AZNLOCTEQ6L1_Wmintaunu_400M600.merge.e3663_s2608_s2183_r7326_r6282_p2495	5577	Reco	7408468	560	5	0	0	0		Jan 11 10:01	Jan 11 11:3	3 GP	5
mc15_13TeV.301167.PowhegPythia8EvtGen_AZNLOCTEQ6L1_Wmintaunu_1250M1500.merge.e3663_s2608_s2183_r7326_r6282_p2495	5577	Reco	7408464	560	3	4	0	17438	done	Jan 11 10:01	Jan 11 12:3	3 GP	5
mc15_13TeV.301160.PowhegPythia8EvtGen_AZNLOCTEQ6L1_Wmintaunu_120M180.merge.e3663_s2608_s2183_r7326_r6282_p2495	5577	Reco	7408461	560	30	10	0	2078		Jan 11 10:01	Jan 11 12:3	3 GP	5
mc15_13TeV.301144.PowhegPythia8EvtGen_AZNLOCTEQ6L1_Wplustaunu_600M800.merge.e3663_s2608_s2183_r7326_r6282_p2495	5577	Reco	7408457	560	3	4	0	14241	done	Jan 11 10:01	Jan 11 13:3	3 GP	5
mc15_13TeV.301178.PowhegPythia8EvtGen_AZNLOCTEQ6L1_Wmintaunu_5000M.merge.e3663_s2608_s2183_r7326_r6282_p2495	5577	Reco	7408453	560	3	4	0	17631	done	Jan 11 10:00	Jan 11 13:3	3 GP	5
mc15_13TeV.301168.PowhegPythia8EvtGen_AZNLOCTEQ6L1_Wmintaunu_1500M1750.merge.e3663_s2608_s2183_r7326_r6282_p2495	5577	Reco	7408449	560	3	4	0	17861	done	Jan 11 10:00	Jan 11 13:3:	3 GP	5



Production tasks manage interface

- Current implementation of Request/Task monitor
 - Server-side
 - Initial data load takes time
 - Slow re-filtering, sorting etc.
 - Need to implement caching
- New implementation of Tasks monitor
 - Currently implemented for Tasks at given Request
 - Client-side
 - Still initial data load takes time
 - Fast filtering, sorting, search on the page
 - New features
 - Filter per column
 - Expandable fields
 - □ Common search including expandable fields
 - Smart row selection
 - Still under developmen

New Task monitor view

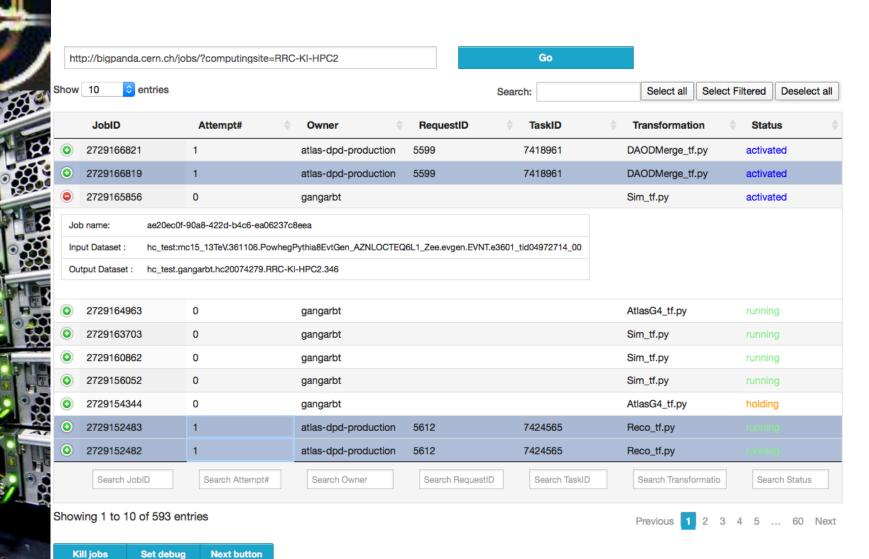


Showing 1 to 10 of 2,335 entries

Previous 1 2 3 4 5 ... 234 Next

Abort Finish Retry Reassign Parameters Obsolete Kill jobs Login

Implementation for Job manage





Plans

- HPC: Efficiency and stability increase of the PanDA sites for ATLAS
 - ☐ Increase "maxrun" jobs number
 - New python implementation of Condor-SLURM connector with new features
 - Learning failed jobs logs
- Prodsys2: Continue to support and develop the interface



Acknowlegement

This work was funded in part by the U.S. Department of Energy, Office of Science, High Energy Physics and Advanced Scientific Computing Research under Contracts No. DE- SC0012704, No. DE-AC02-98CH10886 and Contract No. DE- AC02-06CH11357. NRC-KI team work was funded by the Russian Ministry of Science and Education under Contract No 14.Z50.31.0024. Supercomputing resources at NRC-KI are supported as a part of the center for collective usage (project RFMEFI62114X0006, funded by Russian Ministry of Science and Education).

