# PanDA @ NRC KI

R. Mashinistov
National Research Centre "Kurchatov Institute", Moscow, Russia
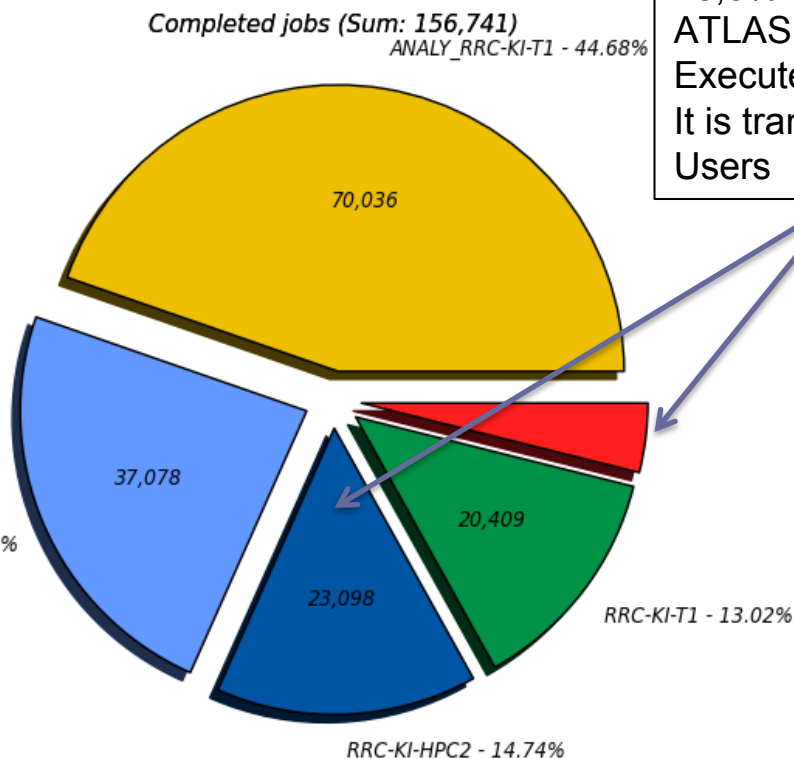
# HPC2 supercomputer at NRC KI



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 (ANALY_RRC-KI grid site)

# First steps

- PanDA@NRC KI
  - server, auto pilots factory, monitor and database server (MySQL)
- After APF was installed in 2014 we immediately set up Analysis queue
  - Condor-SLURM connector
- In 2015 we set up Production queue



Completed jobs (Sum: 156,741)

ANALY_RRC-KI-T1 - 44.68%

70,036

RRC-KI-T1_MCORE - 23.66%

37,078

23,098

20,409

RRC-KI-T1 - 13.02%

RRC-KI-HPC2 - 14.74%

18,6%
ATLAS MC Production and Analysis
Executed on SC@NRC-KI
It is transparent for Prodsys and
Users

- ANALY_RRC-KI-T1 - 44.68% (70,036)
- RRC-KI-T1_MCORE - 23.66% (37,078)
- RRC-KI-HPC2 - 14.74% (23,099)
- RRC-KI-T1 - 13.02% (20,409)
- ANALY_RRC-KI-HPC - 3.90% (6,119)

# Computing portal @ NRC KI

- Web-user interface (FLUSK)
- File catalog & data transfer system
  - Plugin based (ftp, http, grid and etc.)
- Authorization OAuth 2.0
- API for external applications
  - Token authentication
  - Upload/fetch files
  - Send a job
  - Get job status/statistics
- Private FTP storage for each user

# User friendly jobs running and monitoring interface
# Submit jobs and obtain results



- Job setup interface
  - Easy setup distributive, input files, parameters and output file names
  - Local authentication (users don't need a certificate)

- Jobs monitor

- After submitting a job you can monitor its status in the "Job list" tab. When the job finishes, you can get your results from the detailed info page
  - Links to the input/output files

Ruslan Mashinistov

# Portal operation

- •OAuth 2.0 authentication
- •Asynchronous operations – Celery
- •API for external applications

WEB-based user interface → WWW

External tools endpoint → API

Private ftp storage for user's input files → FTP

WEBPANDA interface → New job → PanDA server

files

VWN VWN VWN
NFS
Cloud

WN WN WN
LUSTRE
HPC

Ruslan Mashinistov

# HPC Pilot for biology jobs

- Pilots are running on auxilliar node
- Pilot runs job on the WN's via SLURM

# Biology application

- Next Generation Genome Sequencing (NGS)
- Analysis of ancient genomes sequencing data (Mammoths DNA) using popular software pipeline PALEOMIX can take a month even running it on the powerful computer resource. PALEOMIX include typical set of software used to process NGS data
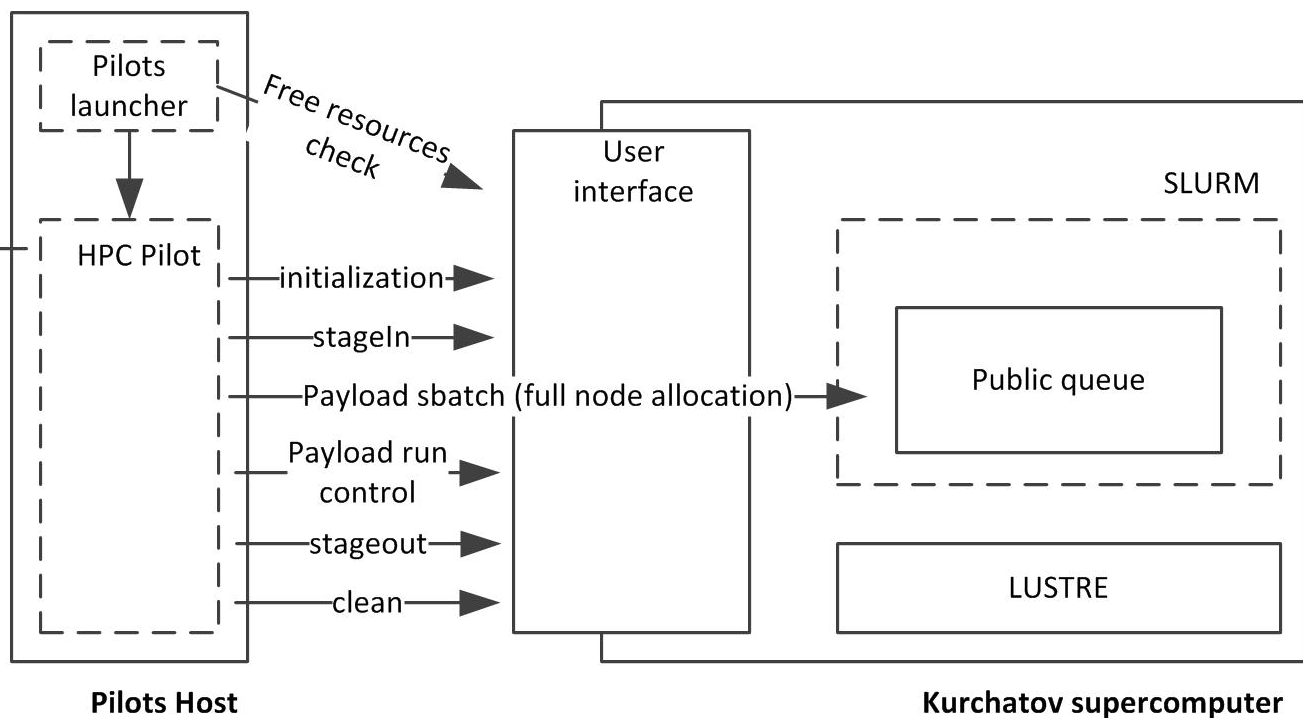- We adapted the PALEOMIX pipeline to run it on a distributed computing environment powered by PanDA.
- To run pipeline we split input files into chunks which are run separately on different nodes as separate inputs for PALEOMIX and finally merge output file, it is very similar to what it done by ATLAS to process and to simulate data.
- Using software tools developed initially for HEP and Grid can reduce payload execution time for Mammoths DNA samples from weeks to days.



~ 4-5 weeks

135 chunks

...

~ 4 days

megaPanDA

# JEDI Oracle->MySQL migration

- Originally all ID's were maintained using Oracle sequences
- While migration PanDA to MySQL this concept was replaced by auto increment columns.
- While migrate JEDI to MySQL we met the inconsistency issue

Oracle: sequence

MySQL: Auto increment

**TaskBuffer**

PanDA

JEDI

**JobGenerator**

| jobsdefined4 | |
|---|---|
| **PANDAID** | ... |
| 1 | ... |
| 2 | ... |
| ... | |

| jobswaiting4 | |
|---|---|
| **PANDAID** | ... |
| 1 | ... |
| 2 | ... |
| ... | |

| jobsactive4 | |
|---|---|
| **PANDAID** | ... |
| 1 | ... |
| 2 | ... |
| ... | |

# JEDI Oracle->MySQL migration

- Back to the sequences concept
- But as MySQL don't have Sequences we implemented it
- We add new table "sequence" and implemented 2 functions (nextval() and curval()). This allowed to handle sequences in a very similar way how it's done in Oracle.
- All the changes in the code a localized in the WrappedCursor
  - 1st WrappedCursor realisation:

    *Schema_name.SEQUENCE_NAME.nextval -> NULL*

  - WrappedCursor update:

    *Schema_name.SEQUENCE_NAME.nextval -> Schema_name.nextval(SEQUENCE_NAME)*

| sequence | name | increment | min_value | max_value | cur_value | cycle |
|---|---|---|---|---|---|---|
| servicelist | cloudtasks_id_seq | 1 | 1 | 9223372036854775807 | 1 | 0 |
| site | filestable4_row_id_seq | 1 | 1 | 9223372036854775807 | 47 | 0 |
| siteaccess | jedi_dataset_cont_fileid_seq | 1 | 1 | 9223372036854775807 | 192 | 0 |
| sitedata | jedi_datasets_id_seq | 1 | 1 | 9223372036854775807 | 162 | 0 |
| siteddm | jedi_output_template_id_seq | 1 | 1 | 9223372036854775807 | 83 | 0 |
| sitehistory | jedi_work_queue_id_seq | 1 | 1 | 9223372036854775807 | 1 | 0 |
| sites_matrix_data | subcounter_subid_seq | 1 | 1 | 9223372036854775807 | 1 | 0 |
| sitesinfo | jobsdefined4_pandaid_seq | 1 | 1 | 9223372036854775807 | 24 | 0 |
| sitestats | group_jobid_seq | 1 | 1 | 9223372036854775807 | 1 | 0 |
| subcounter_subid_seq | prodsys2_task_id_seq | 1 | 1 | 9223372036854775807 | 1 | 0 |

# Summary

- ❑ Migration to MySQL
  - ❑ panda-server  **DONE**
  - ❑ jedi-server **PROGRESS**
- ❑ Update the JEDI installation to make it standard PROGRESS
- ❑ Twiki for PanDA@NRC-KI
  - ❑ https://twiki.cern.ch/twiki/bin/view/PanDA/BigPanDAforNRCKI
  - ❑ Needed to be reviewed and updated
- ❑ Setupper & Adder plugins
  - ❑ Implemented for NRC KI but not documented
  - ❑ Class Setupper – called while receiving a job
    - ❑ Checks input files, it weights and checksum.
    - ❑ Gets srm-links to the files in given datasets
  - ❑ Class Adder – called while pilot sends job 'finished'('failed') status
    - ❑ Files already in grid and we have srm-links
    - ❑ Registers output data in Rucio
  - ❑ Plugins could be set in panda_server.cfg
- ❑ PanDA queue definition in schedconfig table & "cached" json files for pilot
  - ❑ AGIS handle it for central ATLAS server
  - ❑ Not documented for others

# Backup slides

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