# Optimization for Bioinformatics genome sequencing pipelines by means of HEP computing tools for Grid and Supercomputers

A.Novikov, V. Aulov, D. Drizhuk, A. Klimentov, R. Mashinistov, A.Nedoluzhko, A.Poyda, F. Sharko, I. Tertychnyy, A. Teslyuk
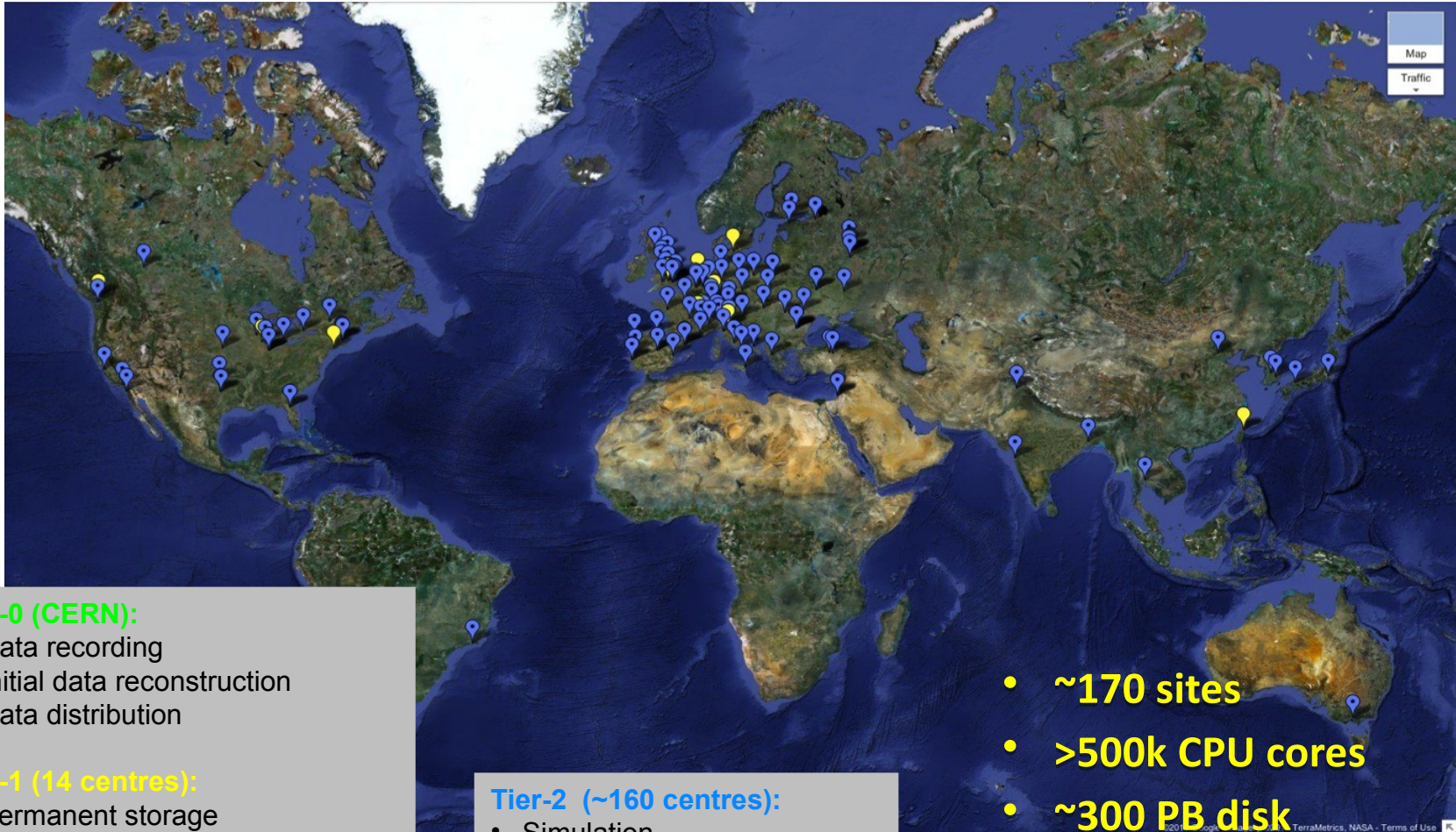
*NRC "Kurchatov Institute"*

4-9 July 2016, GRID 2016, Dubna

# Content

1. Introduction

2. PanDA WMS for ATLAS experiment at Large Hadron Collider, CERN.

3. PanDA instance at NRC KI, adaptation to non HEP science branches. Portal for bioinformatics task.

4. Approbation on ancient Mammoth DNA sequencing tasks with PALEOMIX software. Optimizations and pipelines support.

5. Conclusion

# World GRID Resources



**Tier-0 (CERN):**
- Data recording
- Initial data reconstruction
- Data distribution

**Tier-1 (14 centres):**
- Permanent storage
- Re-processing
- Analysis

**Tier-2  (~160 centres):**
- Simulation
- End-user analysis

- **~170 sites**
- **>500k CPU cores**
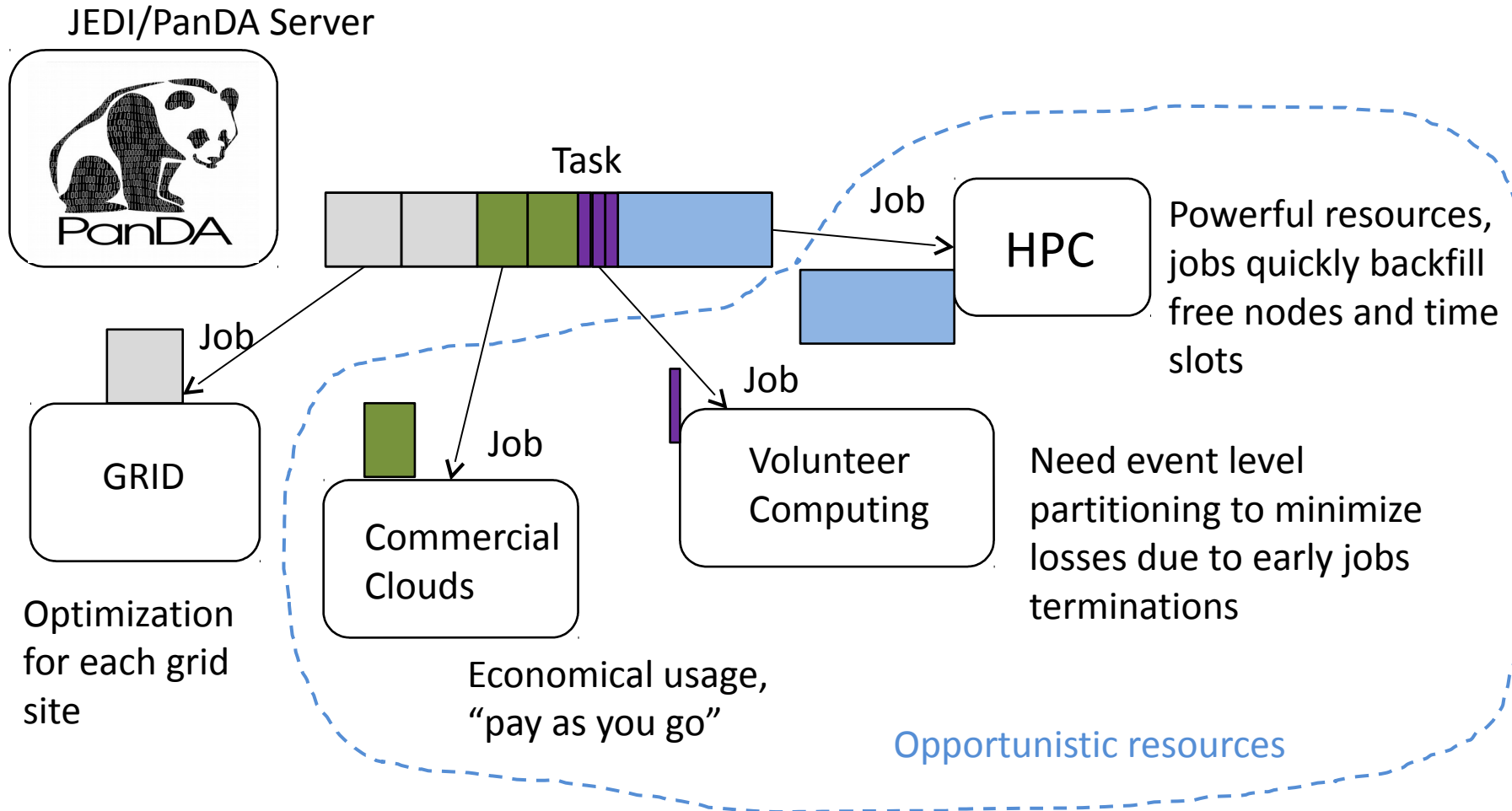- **~300 PB disk**

http://wlcg.web.cern.ch/

# PanDA for HEP

**PanDA** - Production and Distributed Analysis, is a workload management system (WMS) and a project, developed for ATLAS experiment at Large Hadron Collider, CERN.
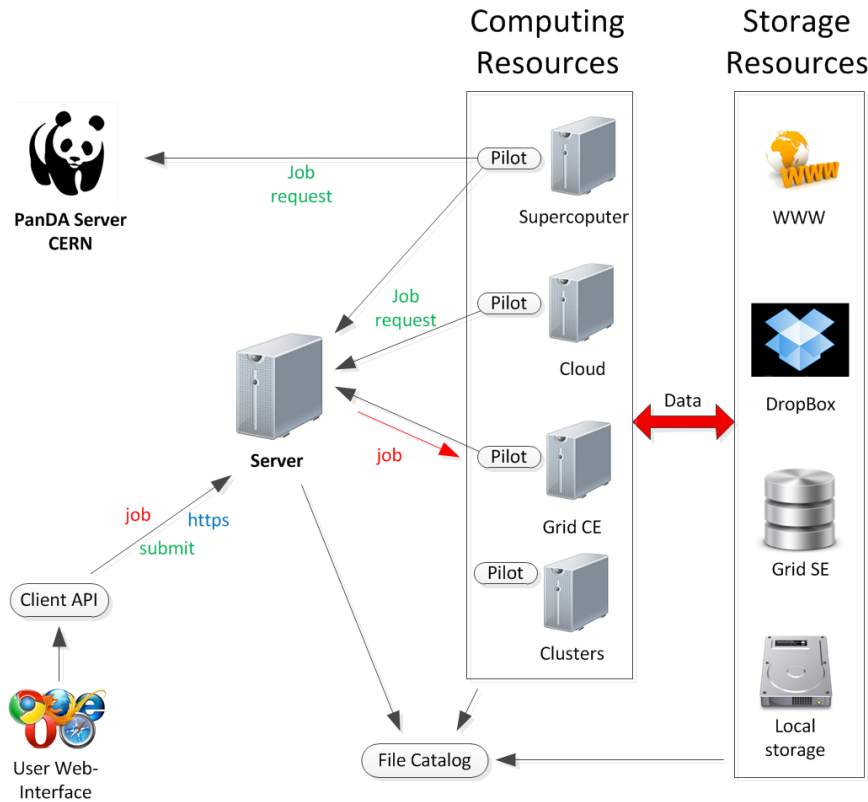
Main features:

1. Single central queue from hundreds of sites for users.

2. Reducing error rate and operational costs at site support level by a system launching pilots jobs.

3. Support for different grid middlewares (and their versions) with unified high level view of workflow for users.

4. Hide from users complexity of low level workflow control automatization. For example, user can submit single task with many input files, and the system in auto mode will define how to split task on subtasks, when and where execute them and re-run in case of failures.

5. Utilize single PanDA WMS for processing real (experiment) and simulated (event generation tasks) data, users analysis and for all other operations with data.

6. Support possibility of integration with cloud infrastructure, HPC, etc.

# Dynamic job definition and workload partitioning in PanDA

JEDI/PanDA Server



Task

Job → HPC — Powerful resources, jobs quickly backfill free nodes and time slots

Job → GRID

Optimization for each grid site

Job → Commercial Clouds

Economical usage, "pay as you go"

Job → Volunteer Computing

Need event level partitioning to minimize losses due to early jobs terminations

Opportunistic resources

# PanDA at NRC KI (MegaPanDA)



High Performance second generation cluster HPC2 with peak performance 122,9 TFLOPS (commissioned 2011) (10240 cores = 1280 nodes 2x Intel Xeon E5450 3,00ГГц 4 cores 16 Гб RAM). #2 in 15-th issue of Russian top50 Supercomputers

# Adaptation of bioinformatics tasks

Users, in particular, bioinformatics need some adaptations:

- Need more then default resources (memory >1-2Gb per job).

- Possibility to run multiple-core or MPI jobs.

- Will work with FTP rather then grid catalog.

- Cannot submit tasks without grid X.509 certificate.

- Has peaks and declines in overall workload.

- Give support for special software (Bowtie2, Abyss, PALEOMIX).

# Web user interface

**NEW JOB**

SOFTWARE: bowtie2: 1.5.2

INPUT FILES: drag & drop — Обзор... Файлы не выбраны.

1 files ready for upload!

INPUT FILES: ftp dir — ftpdir1

INPUT FILES: http url — http://storage.com/myfile.txt — Add

INPUT FILES: guid — web.it_f2e1920f-9b22-4878-94ab-82b1eaef2983 — Add

INPUT FILES: container — Add

☐ One file one job

CORES: 8

RUN SCRIPT: mkdir tmp; mkdir out; bam_pipeline run --max-threads=2 --jar-root=$JAR_ROOT --temp-root=tmp --destination=out ../Mammoth.aaaaaaaaaa.yaml

TAGS: run2  x

Reset   Send job

Update period: 5 min

Show 10 entries      Search:

| ID | Owner | PandaID | Distributive | Created | Modified | Attempt | Status |
|----|-------|---------|--------------|---------|----------|---------|--------|
| 2384 | ruslan | 2455 | paleomix_bam [1.1.0] | 21.03.2016 5:53 | 22.03.2016 3:46 | 1 | finished |
| 2383 | ruslan | 2454 | paleomix_bam [1.1.0] | 18.03.2016 8:23 | 20.03.2016 7:16 | 0 | finished |
| 2382 | ruslan | 2453 | paleomix_bam [1.1.0] | 18.03.2016 8:23 | 20.03.2016 7:16 | 0 | finished |
| 2381 | ruslan | 2452 | paleomix_bam [1.1.0] | 18.03.2016 8:23 | 20.03.2016 6:40 | 0 | finished |
| 2380 | ruslan | 2451 | paleomix_bam [1.1.0] | 18.03.2016 8:23 | 20.03.2016 5:55 | 1 | finished |
| 2379 | ruslan | 2450 | paleomix_bam [1.1.0] | 18.03.2016 8:22 | 20.03.2016 5:31 | 4 | finished |
| 2378 | ruslan | 2449 | paleomix_bam [1.1.0] | 18.03.2016 8:22 | 19.03.2016 21:16 | 3 | failed |

| GUID | TYPE | LFN | LINK |
|------|------|-----|------|
| web.it_78385e8b-f44f-4442-b61f-47328c6d32d2 | input | loxAfr3.fasta | [http] |
| web.ruslan_707a6ebf-27c8-4ec9-af7d-1d87f408c6aa | input | Mammoth.aaaaaaaacs.yaml | [http] |
| web.it_m_113bde47-7fd7-4442-9ab2-3e9737d001b4 | input | Mammoth.1.aaaaaaaacs.fastq | [http] |
| web.it_m_4beb6833-24e7-490e-9e9b-a88ff8358be1 | input | Mammoth.2.aaaaaaaacs.fastq | [http] |
| web.ruslan_86f9a9a5-32c2-4f85-8248-4c4902e1dbdc | output | results.tgz | [http] [ftp] |
| web.ruslan_a9bd8724-1135-43ef-a3c0-a6cb2498afc2 | log | job.c3748a94-3add-4af6-b74b-7a6ba9dcd87e.log.tgz | [http] [ftp] |

Tasks setup interface: easy definition for distributive, input files, parameters and output file names.

Local authentication with OAuth 2.0 (users don't need a certificate).

FTP (or other) storage for user I/O and flexible data  transfer system in backend.

Jobs monitor.

Results from the detailed info page (available only for finished tasks) are links to the I/O files.

# Ancient DNA analysis

Ancient DNA (aDNA) is DNA isolated from ancient specimens such as archaeological an paleontological remains.

Ancient DNA is analyzed from:

- Mummies;
- Organisms preserved in amber;
- Plant materials found in ancient tombs;
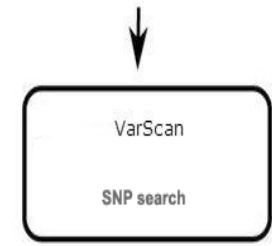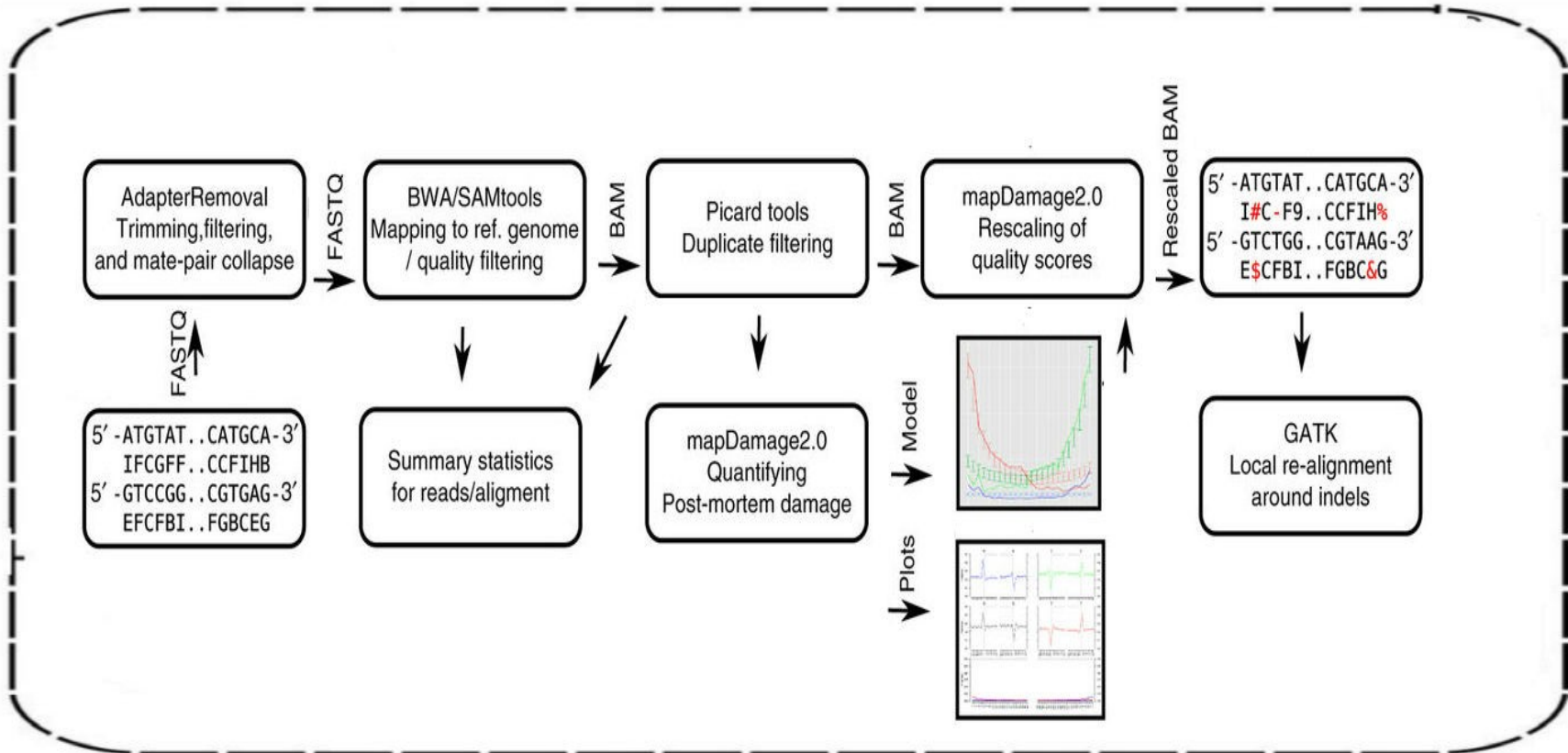- Bacteria;
- Bones;
- Permafrost
- Etc.

**Difficulties of DNA analysis**:

- DNA degradation;
- DNA contamination
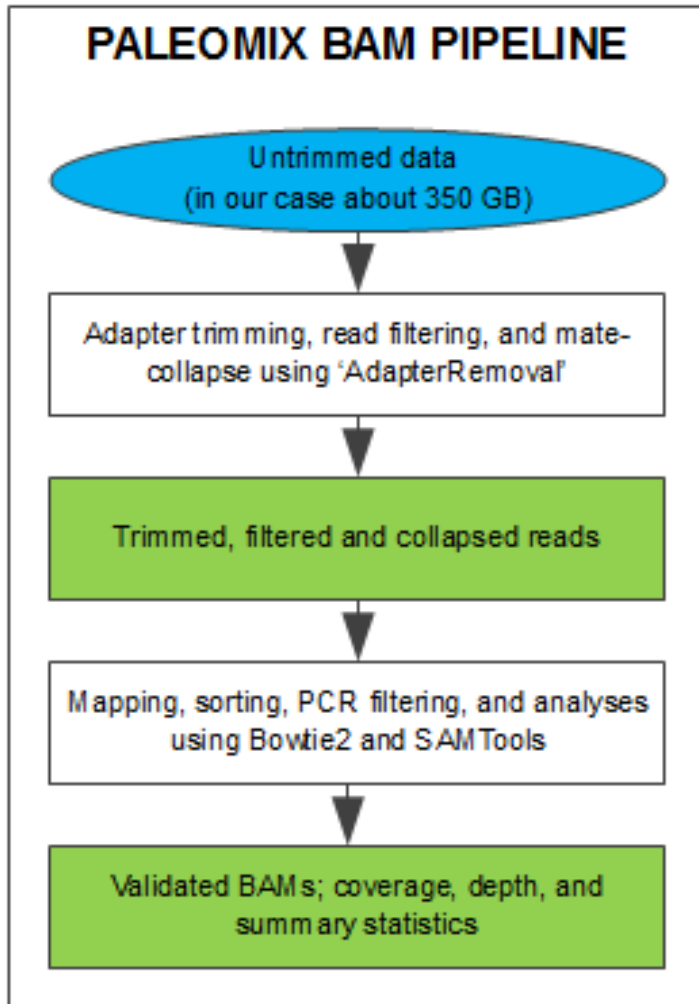- DNA postmortem mutation
- and etc



Древние жители Кавказа из раскопок археолога В.Р.Эрлиха
(Государственный музей искусства народов Востока, Москва). Источник: www.mk.ru
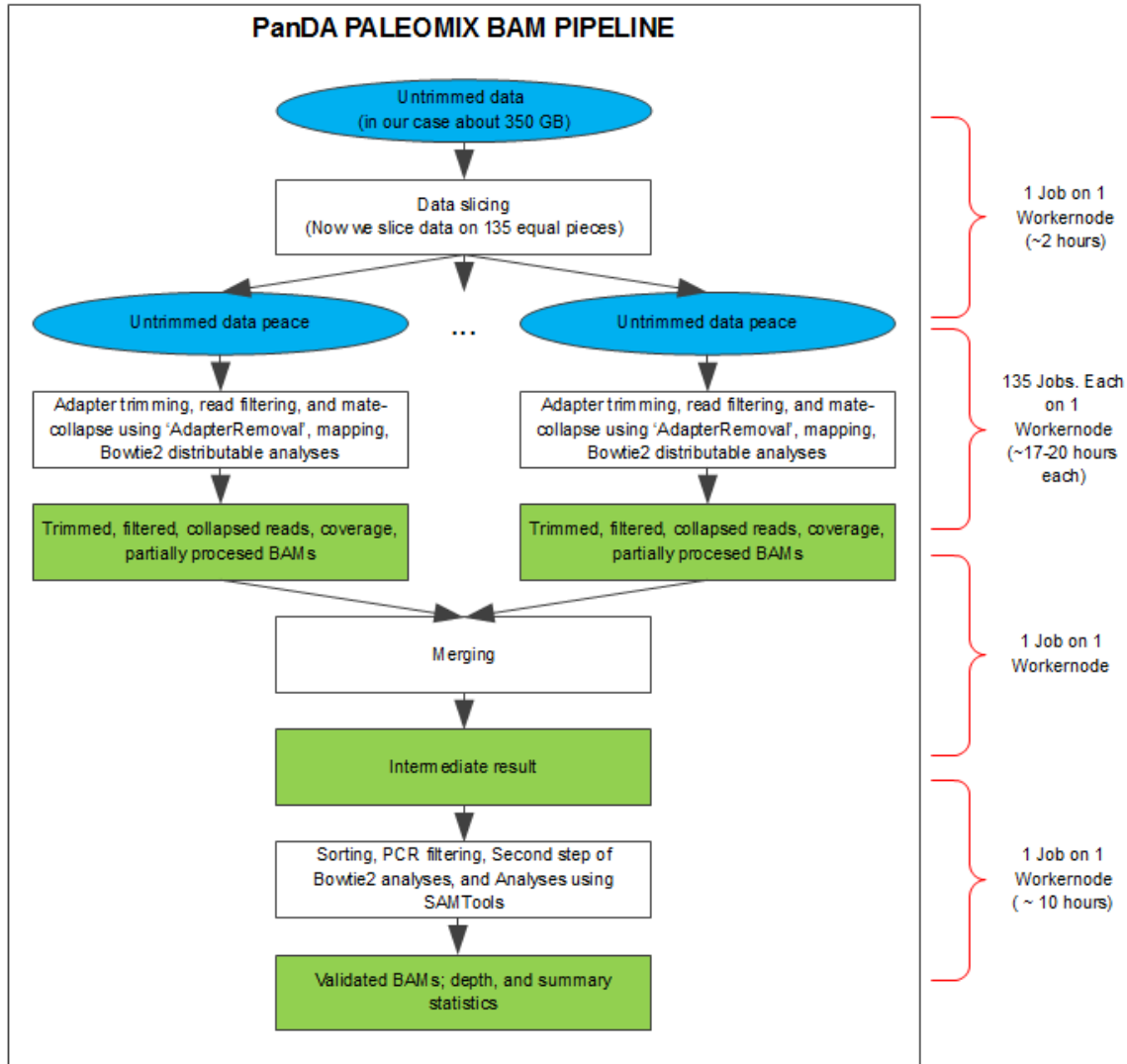
# Data analysis: PALEOMIX



*Schubert et al., (Nature Protocols) 2014

# PALEOMIX pipeline task adaptation

## PALEOMIX BAM PIPELINE

Untrimmed data
(in our case about 350 GB)

↓

Adapter trimming, read filtering, and mate-collapse using 'AdapterRemoval'

↓

Trimmed, filtered and collapsed reads

↓

Mapping, sorting, PCR filtering, and analyses using Bowtie2 and SAMTools

↓

Validated BAMs; coverage, depth, and summary statistics

- The PALEOMIX pipeline is a user-friendly package designed for Unix-like systems and largely automates the analyses related to whole genome re-sequencing. It is compatible with a full range of sequence data and performs a series of user-defined analyses, including read trimming, collapsing of overlapping mate-pairs, read mapping, PCR duplicate removal, SNP calling, and metagenomic profiling.

- For ancient DNA sequence data, the PALEOMIX pipeline also supports the quantification of *post-mortem* DNA damage and standard mis-incorporation and fragmentation patterns. Finally, in situations where several genomes are available, the PALEOMIX pipeline can reconstruct Maximum Likelihood phylogenomic trees and reveal the evolutionary phylogenetic relationships among taxa.

- The PALEOMIX pipeline has been developed by researchers from Ludovic Orlando's group at the Centre for GeoGenetics, University of Copenhagen, Denmark. The software and related documentation is available at https://github.com/MikkelSchubert/paleomix

# Optimized PALEOMIX pipeline



**PanDA PALEOMIX BAM PIPELINE**

Untrimmed data
(in our case about 350 GB)

Data slicing
(Now we slice data on 135 equal pieces)

1 Job on 1 Workernode (~2 hours)

Untrimmed data peace ... Untrimmed data peace

Adapter trimming, read filtering, and mate-collapse using 'AdapterRemoval', mapping, Bowtie2 distributable analyses

135 Jobs. Each on 1 Workernode (~17-20 hours each)

Trimmed, filtered, collapsed reads, coverage, partially procesed BAMs

Merging

1 Job on 1 Workernode

Intermediate result

Sorting, PCR filtering, Second step of Bowtie2 analyses, and Analyses using SAMTools

1 Job on 1 Workernode (~ 10 hours)
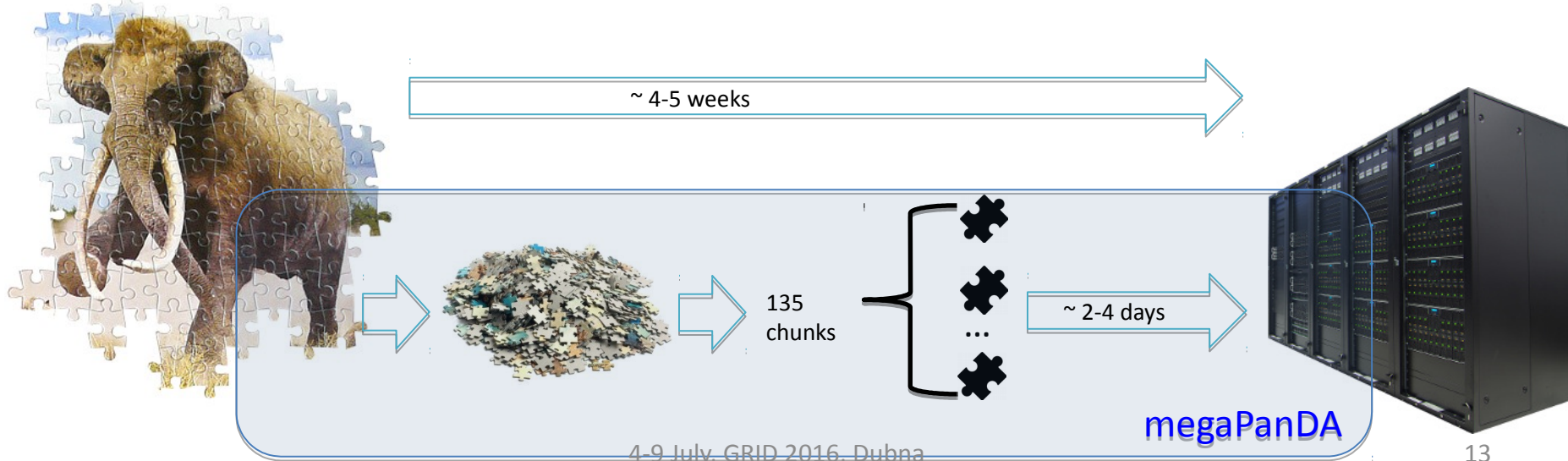
Validated BAMs; depth, and summary statistics

The portal fits standard PanDA computational scheme and shows most efficiency for compute-intensive tasks that could be split into many sub-jobs to be computed in parallel. To hide execution complexity and manual routines from end-users we introduced and seamlessly integrated into the portal original pipelines control system, that automatically (and without user prompt) split input data, prepare and run sub-tasks as ordinary PanDA jobs and merge results.

So we assume that every pipeline contains several steps (as tasks), each of which consists of some server side preparation and one or many standard PanDA jobs. Steps with many jobs executed in parallel-mode by PanDA.

# Ancient mammoth genome sequencing pipeline task adaptation

- Next Generation Genome Sequencing (NGS) (DNA sequencing is the process of determining the precise order of nucleotides within a DNA molecule).

- 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 is done by ATLAS to process and to simulate data.

- Using software tools developed initially for HEP and Grid one can reduce payload execution time for Mammoths DNA samples from weeks to days. Performed on data from Genome analysis laboratory at NRC "Kurchatov institute", 350Gb.

~ 4-5 weeks

135 chunks

...

~ 2-4 days

megaPanDA

# About mammoth



- Woolly mammoths (*Mammuthus primigenius* Blum.) were an evolutionary dead end of genus Mammuthus which arose in Africa and migrated to Eurasia almost three million years ago.

- Paleozoologists described as minimum fossils of three species in these genera – *M. meridionalis* (Early Pleistocene), *M. trogontherii* (Middle Pleistocene) and *M. primigenius* (Late Pleistocene). Woolly mammoth appeared 300 – 200 thousands years ago in Siberia and after that colonized Europe and North America.

- The latest mammoth population disappeared on St. Paul Island 6 kyr ago and Vrangel Island 4 kyr ago because of inbreeding and loss of genetic diversity. Novel methods of ancient DNA provide significant abilities for genomic analysis of extinct species such as Pleistocene megafauna species.

# About mammoth (2)

Woolly mammoth calf known as Khroma (who was excavated in October 2008 in the Khroma River, Yana-Indigirka lowland, Yakutia, Eastern Siberia (Russia); the specimen's AMS age was at background levels, that is >50,000 years ).

Фото: Mammothportal.com

# Conclusion and future plans

Problems and needs launching users pipelines tasks for bioinformatics are described. Used developed at NRC KI interfaces and portal for launching task by means of PanDA WMS, initially developed to support HEP experiment ATLAS at LHC, CERN.

We performed the adaptation of the PALEOMIX pipeline to a distributed computing environment powered by PanDA for Ancient Mammoths DNA samples. We used PanDA to manage computational tasks on a multi-node parallel supercomputer.

The approbation was performed on an ancient  mamoth DNA sequencing task, for which the total computational time  dramatically reduced from several weeks to 3-4 days.  This  includes decreasing the total computation time because of jobs brokering, submission and auto resubmission of failed ones by means of PanDA, which also demonstrated it earlier for the HEP applications in the Grid.

Thus using software tools developed initially for HEP and Grid can reduce computation time for bioinformatics tasks such as PALEOMIX pipeline for Ancient Mammoths DNA samples from weeks to days.   This approach allows performing compute-intensive sciences workflows (as HEP, bioinformatics, astrophysics, etc.) using joined compute powers of different infrastructures.

# Acknowledgements

Many thanks to PanDA Core SW team.

# Thank you for attention!
## Questions?
### novikov@wdcb.ru

# Архитектура PanDA WMS



Production managers

PanDA server

Data Management System

Local Replica Catalog

submitter (bamboo/JEDI)

**production job**

https

define

task/job repository (Production DB)

**analysis job**

https

submit

End-user

EGEE/EGI

pull

https

job

Logging System

OSG

pilot

pilot

https

https

NDGF

pilot

arc

ARC Interface (aCT)

pilot

Worker Nodes

condor-g

pilot scheduler (autopyfactory)

NonPanDA components