# JINR Tier1 Center Status and Plans

**Andrey Baginyan,   Anton Balandin,  Sergey Belov,**

**Andrey Dolbilov,   Alexey Golunov,  Natalia Gromova,**

**Ivan Kadochnikov, Ivan Kashunin, Vladimir Korenkov, Valery Mitsyn,**

**Igor Pelevanyuk, Sergei Shmatov, Tatiana Strizh, Vladimir Trofimov,**

**<u>Nikolay Voytishin</u>, Victor Zhiltsov**
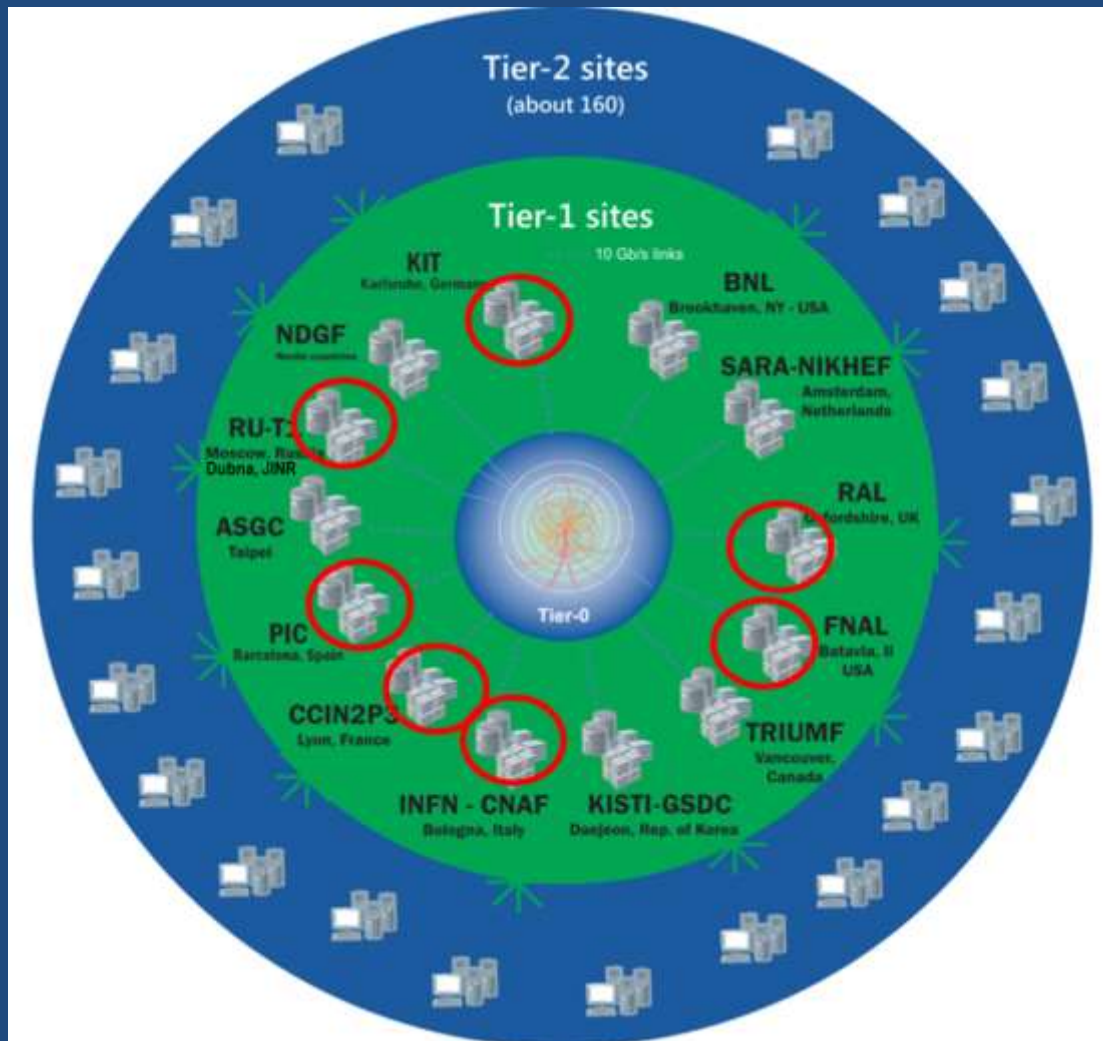
**ROLCG-2018**
**Cluj**
**18-10-18**

# Outline

- How we started
- Main functions
- Infrastructure
- Network and telecommunication channels
- Resources
- Monitoring
- How well does it work?
- Plans for 2019

# LHC Computing Model + WLCG

WLCG computing enabled physicists to announce the discovery of the Higgs Boson on 4 July 2012



42 countries
170 computing centers
2 million tasks run every day
800,000 computer cores
500 petabytes on disk and
400 petabytes on tape

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

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

**Tier-2 :**
Simulation
End-user analysis

# Joint NRC "Kurchatov Institute" – JINR Tier1 Computing Centre

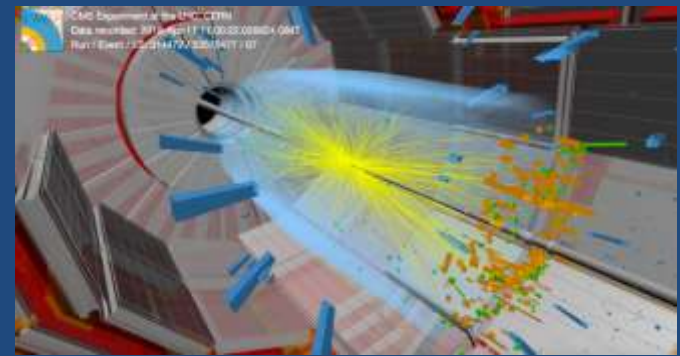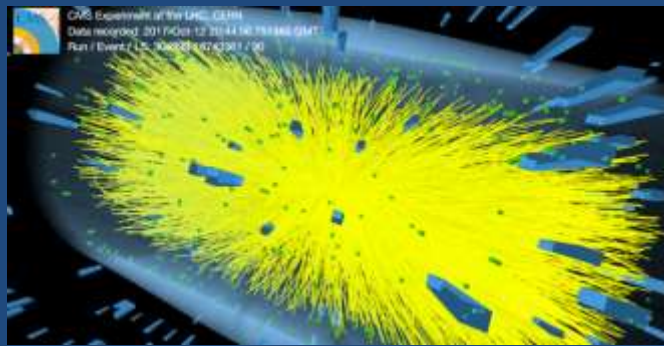➢ **Proposal to create the WLCG Tier1 center in Russia: March 2011, accepted in October 2012**

➢ *The Federal Target Programme Project:* **«Creation of the automated system of data processing for experiments at the LHC of Tier1 level and maintenance of Grid services for a distributed analysis of these data»**
*Duration:* **2011 – 2013**

**Russia Tier1 full scope start-up in WLCG in 2014 NRC "Kurchatov Institute" supports ATLAS, ALICE and LHCb, JINR supports CMS (Compact Muon Solenoid)**

**Systematic increase of computing capacity and data storage is needed in accordance with the experiment requirements**

ORGANISATION EUROPEENNE POUR LA RECHERCHE NUCLEAIRE
EUROPEAN ORGANIZATION FOR NUCLEAR RESEARCH
Laboratoire Européen pour la Physique des Particules
European Laboratory for Particle Physics

GENÈVE, SUISSE
GENEVA, SWITZERLAND

Mail address:
Dr. Ian Bird
CERN, IT Department
CH-1211 GENEVE 23
Switzerland
Tel: +41 22 767 5888

E-mail : Ian.Bird@cern.ch

Prof. Mikhail Kovalchuk
Director of National Research Centre "Kurchatov Institute"
1, Akademia Kurchatova pl.,
Moscow 123182, Russia

Prof. Victor Matveev
Director of Joint Institute for Nuclear Research
Joliot-Curie 6
141980 Dubna, Moscow Region, Russia

Votre référence/Your reference:
Notre référence/Our reference:

Subject: Acceptance of the proposal to build Tier 1 centres in Russia

Geneva, October 12, 2012

Dear Directors,

As you know, the proposals from the National Research Centre – "Kurchatov Institute" and the Joint Institute for Nuclear Research, Dubna, to build Tier 1 centres for LHC data analysis were discussed in the recent WLCG Overview Board held on September 28. I am very happy to report that the proposals were well received by the members of the board, and that the decision was made to accept the Russian sites as a new "Associate Tier 1". This decision will be noted in the formal minutes of the meeting.

The next step is now to proceed to signing the WLCG Memorandum of Understanding. The WLCG project office will assist in drafting this MoU, which should be signed by the relevant funding agencies for the two Russian Institutes, or their designated agents.

I am at your disposal for any assistance or to provide further details of the process.

Yours Sincerely,

Dr. Ian Bird
LHC Computing Grid Project Leader
IT Department
CERN

Cc: Prof. Sergio Bertolucci, Dr. Viacheslav Ilyin

4

**In agreement with the CMS Computing model, the JINR Tier1 site provides:**

- acceptance of an agreed share of raw data and Monte Carlo data;
- provision of access to the stored data by other CMS Tier1/Tier2/Tier3 sites of the WLCG;
- service of FTS-channels for Russian and Dubna Member States Tier2 storage elements including monitoring of data transfers

**USER-VISIBLE SERVICES**

- Data Archiving Service
- Disk Storage Services
- Data Access Services
- Reconstruction Services
- Analysis Services
- User Services

**SOME SPECIALIZED SYSTEM-LEVEL SERVICES**

- Mass storage system
- Site security
- Prioritization and accounting
- Database Services

5

Tier1 prototype

Free space

**2014**

| CPU (HEPSpec06) | 14 400 |
|---|---|
| Disk (Terabytes) | 660 |
| Tape (Terabytes) | 72 |

**2018**

| CPU (HEPSpec06) | 72 310 |
|---|---|
| Disk (Terabytes) | 8 319 |
| Tape (Terabytes) | 10 825 |

# Tier1 Infrastructure

- Close-coupled, chilled water cooling InRow
- Hot and cold air containment system
- MGE Galaxy 7000 – 2x300 kW energy efficient solutions 3Ph power protection with high adaptability
- Installation of two new transformers (2.5 MW)
- Guaranteed power supply using two diesel generators

# Network and telecommunication channels



The network infrastructure is meant to provide a 100% availability and reliability of the storage and computing resources of the JINR Tier-1 center.



**Local Area Network** – 10 Gbps, planned upgrade to 100 Gbps
**Wide Area Network** – 100Gbps,
**LHCOPN** - 2x10Gbps
**LHCONE** – 10 Gbps
Upgrade WAN to 2x100Gbps planned
IPv6/IPv4 - enabled

# Tier1 resources



T1_RU_JINR Logical CPU, T1_RU_JINR HEPSPEC06, T1_RU_JINR Disk (TB), T1_RU_JINR Tape (TB)

# Tier1 resources 2018

## Computing Elements (CE)

* *Worker Node (WN)*
  **Typically SuperMicro Blade**

**100 64-bit machines: 2 x CPU (Xeon X5675 @ 3.07GHz, 6 cores per processor); 48GB RAM, 2x1000GB SATA-II; 2x1GbE.**

**175 64-bit machines: 2 x CPU (Xeon E5-2680 v2 @ 2.80GHz, 10 cores per processor), 64GB RAM; 2x1000GB SATA-II; 2x1GbE.**

**Total: 4720 core/slots for batch.**

* *Software*
  **OS: Scientific Linux release 6 x86_64.**
  **BATCH : Torque 4.2.10 (home made)**
  **Maui 3.3.2 (home made)**
  **CMS Phedex**

## Storage Elements(SE)
## Storage System:  dCache

* *Hardware*
*Typically Supermicro and DELL*
*1st - Disk Only:*
  31 disk servers: 2 x CPU (Xeon E5-2650 @ 2.00GHz); 128GB RAM; 112TB h/w ZFS (24x6000GB NL SAS); 2x10G.
  12 disk servers: 2 x CPU (Xeon E5-2660 @ 2.60GHz); 128GB RAM; 70TB ZFS (16x6000GB NL SAS); 2x10G.
  8 disk servers: 2 x CPU (Xeon E5-2650 @ 2.29GHz)  128GB RAM; 150TB ZFS (24x8000GB NLSAS), 2x10G
  12 disk servers: 2 x CPU (Xeon E5-2660 @ 2.60GHz)  128GB RAM; 150TB ZFS (24x8000GB NLSAS), 2x10G

**Total space: 7.3PB**
  3 head node machines: 2 x CPU (Xeon E5-2683 @ 2.00GHz); 128GB RAM; 4x1000GB SAS h/w RAID10; 2x10G.
  8 KVM for access protocols support.

*2nd - support Mass Storage System:*
  8 disk servers: 2xCPU (Xeon X5650 @2.67GHz); 96GB RAM; 63TB h/w RAID6 (24x3000GB SATAIII); 2x10G; Qlogic Dual 8Gb FC.
  8 disk servers: 2 x CPU (E5-2640 v4 @ 2.40GHz); 128GB RAM; 70TB ZFS (16x6000GB NLSAS); 2x10G; Qlogic Dual 16Gb FC.

**Total disk buffer space: 1.1 PB.**
**1 tape robot: IBM TS3500, 10 PB**
  3440xLTO-6 data cartridges; 12xLTO-6 tape drives FC8.
  3 head node machines: 2 x CPU (Xeon E5-2683 v3 @ 2.00GHz); 128GB RAM; 4x1000GB SAS h/w RAID10; 2x10G.
  6 KVM machines for access protocols support

* *Software*
  dCache-3.2
  Enstore 4.2.2 for tape robot.

# Tier1 data exchange 2017 - 2018







**Buffer and TAPE usage:**

- **157 sites worldwide transfer data FROM us;**
- **140 sites transfer data TO us;**
**Leaders are CERN, KIT(Germany), RAL (UK).**

**Storage Element Disk only usage:**

- **316 sites worldwide transfer data FROM us;**
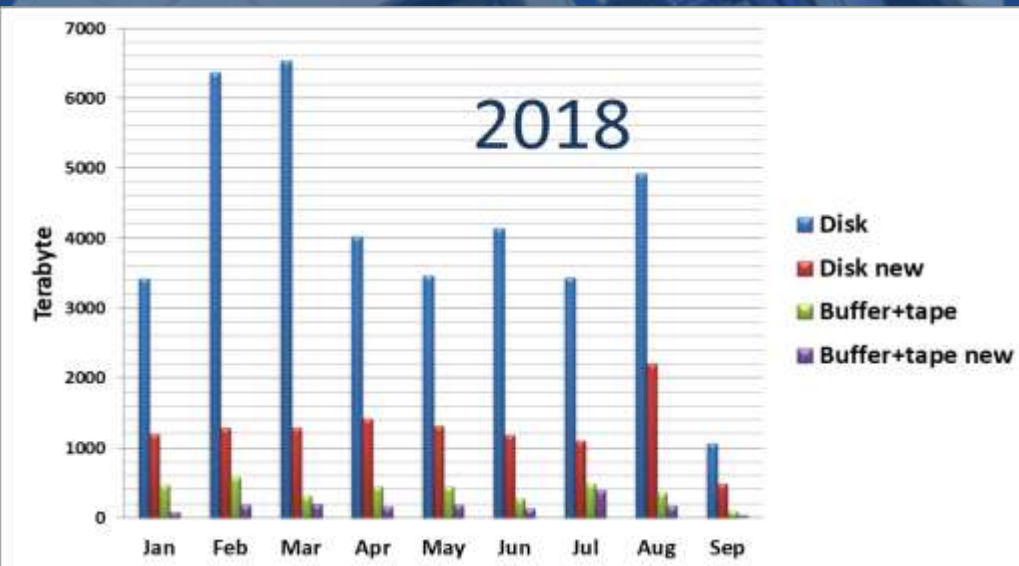- **150 sites transfer data TO us.**

**SE disk only 2018**

| | |
|---|---|
| ***Total*** | 36.50 PB |
| ***new files*** | 11.19 PB |

**SE Buffer + Tape 2018**

| | |
|---|---|
| ***Total*** | 6.78 PB |
| ***new files*** | 3.38 PB |

**SE Buffer + Tape 2017**

| | |
|---|---|
| ***Total*** | 3.64 PB |
| ***new files*** | 1.52 PB |

# Tier1 hardware monitoring +

**For a robust performance of the complex it is necessary to monitor the state of all nodes and services - from the supply system to the robotized tape library.**



**Monitoring data are collected from the wide range of hardware and software related to Tier1**

- cooling systems,
- temperature sensors,
- uninterruptable power supplies (UPS),
- computing servers,
- disk arrays,
- managing services,
- L2 and L3 switches/routers
- tape robot.

**~ 850 elements are under observation**

**~ 8000 checks in real time**

**~ 100 scripts**

**The system allows one, in a real time mode, to observe the whole computing complex state and send the system alerts to administrators and users via e-mail, sms, etc.**

# Tier1 services

Apart from hardware metrics, service metrics are scattered among many internal and external systems.

This information relates to

– data transfers,

– data storage,

– data processing.

In order to keep track of the services admin should regularly check several dozens of web pages. Interpretation of data is more complex.

Aim is:

– Provide a single source of aggregated monitoring information.

– Perform basic analysis of data and provide status of the system.

# Tier1 services monitoring system (I)

Idea is to collect, aggregate and analyze data from different sources. Then provide in comprehensive form on the web page. In case of critical failures – inform administrators.



Web interface

Data Sources

Filter JINR related

Analyze

Database

Service statuses

Store the data

# Tier1 services monitoring system (II)

# PhEDEx operations: issues

PhEDEx system was designed to operate mostly automatically. But sometimes, due to different reasons it requires intervention to fix errors manually.

Source of information about errors is a corresponding PhEDEx webpage. Every error is a big form with source/destination site, time of assigned/start/done, PFNs to/from, transfer/detail/validate logs.

In order to simplify operation python script was written to list important errors and provide relevant information about them.

PhEDEx – CMS Data Transfers

Info **Activity** Data Requests Components Reports Next-gen website

Rate | Rate Plots | Queue Plots | Quality Plots | Routing | Transfer Details | Deletions | **Recent Errors**

PhEDEx

check_phedex_errors.py

Type or errors:
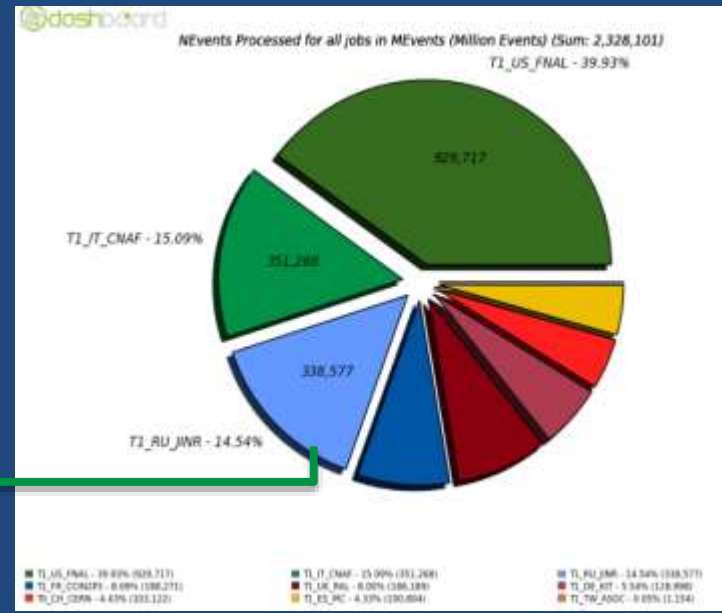nsf - No such file or directory
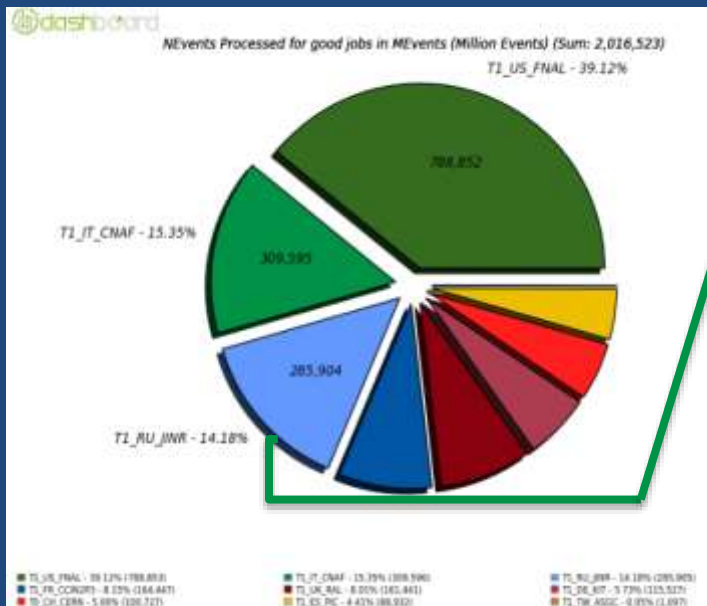csmm - Checksum mismatch
smm - Size mismatch
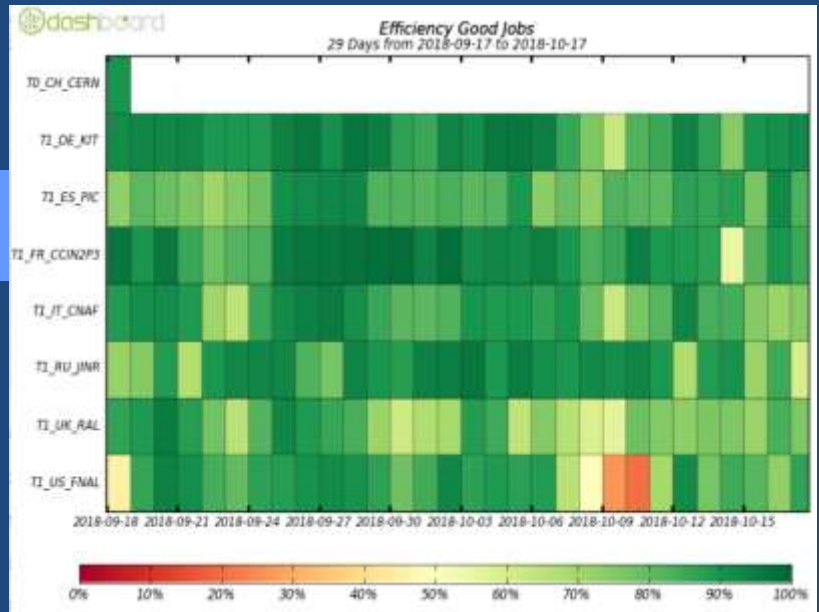uto - User timeout over

List of files in error state

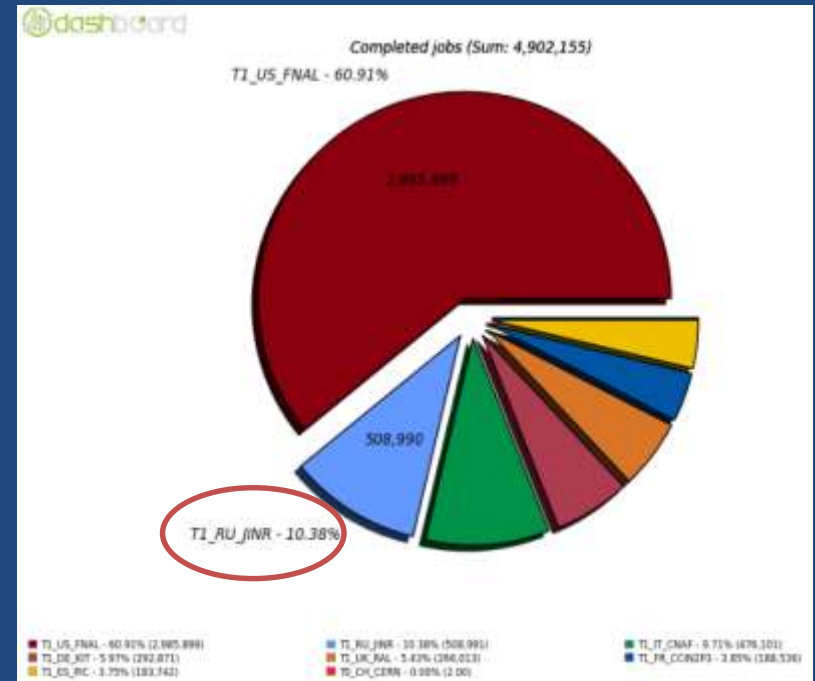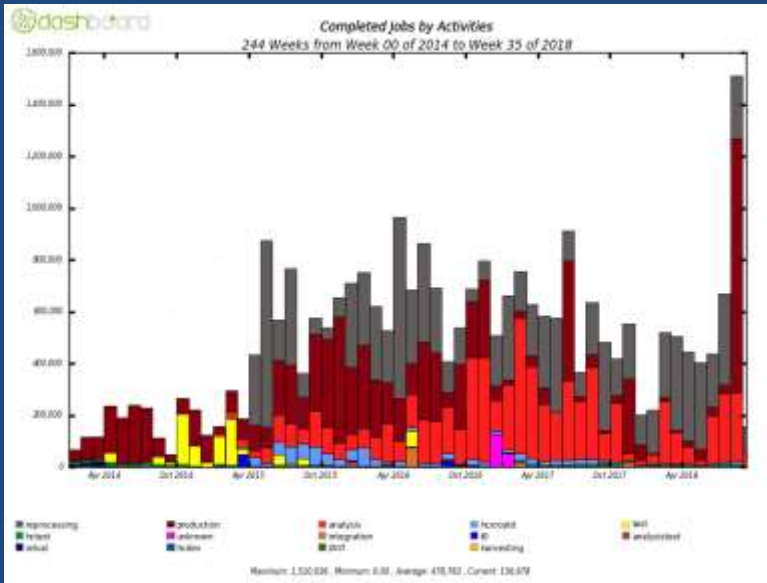# From: 2014-01-01  to: 2018-10-17

**Last Month**

# Jobs Processing by Activities



| Activities | Nevents |
|---|---|
| Analysis | 281 555 814 018 |
| Reprocessing | 10 042 089 961 |
| Production | 5 569 434 071 |
| Test | 428 136 099 |
| Etc. | ….. |

**Total: 308 560 399 956 events**
**Average Rate: 2.086/s**

# Tier1 site reliability

JINR Tier1 (blue) monthly results compared to the average Tier1 reliabilities (orange) as well as with the WLCG target for site reliability (green dotted line) which is set to 97% since 2009, according to the WLCG MoU



**WLCG Sites Reliability**
**CMS**

MOU — RU-JINR-T1 — Average of ALL Tier-0 and Tier-1 sites

# JINR Tier1 site upgrade plans for 2019

CPU (HEPSpec06):

72310 → 90000;

Disk storage volume:

7.3 PB → 8PB;

Tape robot volume:

10 PB → 20 PB.

# Importance of the Tier1 center at JINR

* **Creation of conditions for JINR physicists, JINR Member States, RDMS-CMS collaboration for a full-scale participation in processing and analysis of data of the CMS experiment on the Large Hadron Collider.**

* **The invaluable experience of launching the Tier1 center will be used for creating a system of storage and data processing of megaproject NICA and other scale projects of the JINR-participating countries.**

* **The studies in the field of Big Data analytics assume significance for the development of the perspective directions of science and economy as well as analysis and forecasting of processes in various fields.**

# Thank you for your attention!