# TIER-1 CMS at JINR: Past, Present and Future

N.S. Astakhov, A.S. Baginyan, S.D. Belov, A.G. Dolbilov, A.O. Golunov, I.N. Gorbunov, N.I. Gromova, I.S. Kadochnikov, I.A. Kashunin, V.V. Korenkov, V.V. Mitsyn, I.S.Pelevenyuk, S.V. Shmatov, **T.A. Strizh**, E.A. Tikhonenko, V.V. Trofimov, N.N. Voitishin, V.E. Zhiltsov

Joint Institute for Nuclear Research

GRID 2016

July 4 - 9 2016

# The project of launching Tier-1 at JINR

In November, 2011 at the suggestion of A.A. Fursenko at a session of the Committee on Russia-CERN cooperation, a decision was accepted on the creation in Russia of a Tier-1 center for LHC experiments on the base of NRC "Kurchatov institute" and JINR.

On September 28, 2012, a session of the Supervisory Council of WLCG project approved a work schedule on the creation of Tier-1 in Russia.

**First stage (December, 2012) -** creation of a prototype of Tier-1 at NRC KI and at JINR.
**Second stage (November, 2013) -** installation of equipment for the base Tier-1 center, its testing and finishing up to required functional characteristics.
**Third stage (March, 2015) –** finalization of this complex and commissioning a full-scale Tier-1 center for CMS in JINR.

**In the sequel, a systematic increase of computing capacity and data storage is needed in accordance with the experiment requirements.**

# LHC Computing Model



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

**Tier-1 (11→14 centres):**
- Permanent storage
- Re-processing
- Analysis
- Simulation

**Tier-2 (>200 centres):**
- Simulation
- End-user analysis

# JINR Computing Centre for Data Storage, Processing and Analysis

In accordance with the CMS computing model, the centers Tier-1 for CMS provide a wide range of reliable services for the whole CMS collaboration using standard grid-interfaces of the WLCG project and additional CMS services. Also required is a high-level availability of all Tier-1 services and a constant, in 7x24 mode, technical support for troubleshooting.

The main functions of the Tier-1 centre include:

- organization of sequential data processing (data acquisition, fast browsing (skimming), reprocessing),
- storage of bulk raw experimental and simulated data
- providing data to other Tier-1-2-3 sites for copying and physical analysis,
- receiving "raw" data from Tier-0 centre for a long-time (several years) storage, replication of such data onto other Tier-1 centres.
- data analysis using algorithms of the CMS collaboration.

# Creation of CMS Tier-1 in JINR

- Engineering infrastructure (a system of uninterrupted power supply, climate - control);
- High-speed reliable network infrastructure with a dedicated reserved data link to CERN (LHCOPN);
- Computing system and storage system on the basis of disk arrays and tape libraries of high capacity;
- 100% reliability and availability.

Visible to the user services of the Tier-1 centre include:

1) services of archival storage,

2) services of disk storage,

3) services of access to data and provision of data exchange,

4) services of control and recovery of the system,

5) services for data processing and analysis,

6) services to organize a user access to the WLCG resources.

# CMS Tier-1

**March 2015 – CMS Tier1 Inauguration**

LHCOPN – 10Gbps, 2400 cores (~ 30 kHS06),
5 PB tapes (IBM TS3500), 2.4 PB disk
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

Tape Robot

Computing elements

Uninterrupted power supply

NETWORK

Cooling system

Russia:
NRC KI

US-BNL

Amsterdam/NIKHEF-SARA

Bologna/CNAF

Taipei/ASGC

Ca-
TRIUMF

JINR

NDGF

CERN

US-FNAL

UK-RAL

Barcelona/PIC

Lyon/CCIN2P3

De-FZK

26 June 2009

# Architecture and hardware resources

The computing cluster: 220 physical machines, 3400 cores /54,4 kHS06
To support the batch processing system, a special server with a cluster resource allocation system and a scheduler has been installed.

Storage systems - dCache software.
1. dCache installation - only disk servers and used for operational data storage with fast access to them.
2. dCache installation - disk servers and a tape robot. The disks serve as a buffer zone, tape robot is intended for a long-time storage of data.

Totally, 2 installations have now 3.4 PB of effective disk space, and the tape robot has a 5.4 PB of data storage capacity. To support the storage and access to data, 8 physical and 14 virtual machines have been installed.

# The servers to support the grid-WLCG environment

Logical structure of the Tier-1



WLCG services are installed with grid middleware EMI-3. Currently, 21 services WLCG are installed to provide :

➢ user and virtual organizations (VO) authorization;

➢ task run from VO remote services.

➢ the WLCG information system;

➢ different algorithms of remote testing and verification of the service environment on local resources.

# JINR Tier-1 Connectivity Scheme

# LHCOPN JINR-T1 Traffic (last year)

# Tier-1 Network Structure





Three-layer architecture of the Tier-1 at JINR

Network segment of Tier-1 module at JINR built of Brocade equipment

A diagram of a network architecture implemented Tier-1 module as of 2016. The module consists of twenty-four (24) racks, sixteen (16) of which are filled with a server equipment, and eight (8) allocated for the module cooling and to create the desired climate. Ten (10) racks are filled with disk servers. Three (3) racks are allocated by the computing blade servers. Servers that provide Grid infrastructure occupy three (3) racks.

The operation of switches assumes creation of a virtual factory, which combines up to 32 devices, capable of self-balancing traffic across all pathes.

# JINR Tier1- Tier2 Monitoring

1. **Monitoring the states of all nodes and services- from the supply system to the robotized tape library**
2. **Global real time survey of the state of the whole computing complex**
3. **In case of emergency, alerts are sent to habilitated persons via e-mail, SMS, etc.**
4. **~690 elements are under observation**
5. **~ 3500 checks in real time**

# JINR Tier1 Dashboard

# Tier-1 Brocade Factory Traffic



Output traffic

Input traffic

# Default Metrics for the site T1_RU_JINR

# Last month jobs

# Last month pending and running jobs

# Efficiency good&all jobs

# events processed



NEvents Processed for good jobs in MEvents (Million Events) (Sum: 33,598)

- T1_US_FNAL - 40.15% (13,490)
- T1_RU_JINR - 17.25% (5,797)
- T1_FR_CCIN2P3 - 15.71% (5,279)
- T1_IT_CNAF - 12.67% (4,259)
- T1_UK_RAL - 6.97% (2,343)
- T1_ES_PIC - 4.57% (1,534)
- T1_DE_KIT - 2.62% (881.00)
- T0_CH_CERN - 0.05% (16.00)

Average CPU time spent on one Good Event in secs.

Average WC time spent on one Good Event in secs.

# WallClock Efficiency based on success/all accomplished jobs

# Efficiency based on success/all accomplished jobs

# Resource Utilization

# Data transfer from the JINR T1 via Production and Debug Phedex instances

Pie chart: requests to JINR-T1

- DE 18,14%
- FR 15,44%
- IN 11,99%
- HU 10,78%
- KR 9,01%
- PL 7,37%
- BE 6,78%
- BR 6,02%
- US 4,92%
- IT 2,88%
- GR 2,82%
- TW 1,38%
- HR 0,95%
- BG 0,89%
- PK 0,29%
- UK 0,29%
- NL 0,04%
- CA+SK+FI+UA 0,01%

Raw data transfers to JINR-T1:
250-300 MBps
>1 TB/hour
~30TB/day

# Job processing  by activities

# Happyface installation at JINR

# Development plans for the Tier-1 centre

Planned yearly growth of Tier-1 resources – absolute values and percentage growth over the previous year.

|  | 2016 | 2017 | 2018 | 2019 |
|---|---|---|---|---|
| **Processor capacity of the core/kHS06** | 3400/54,4 | 4200/67,2 (24%) | 5200/83,2 (23%) | 10000/160 (52%) |
| **Disk storage (TB)** | 3390 | 5070 (49%) | 6100 (20%) | 8000 (80%) |
| **Tape storage (TB)** | 10000 | 20000 | 20000 | 20000 |

**CPU needs per event**

**LHC Upgrade 2019-2021. Computing Needs**

Run4
**ATLAS + CMS**

Run3
2020-2022
**ALICE + LHCb**

Run2 :
2015 - 2018

Run1 :
2009 - 2013

- CPU needs (per event) will grow with track multiplicity (pileup) and energy
- Storage needs are proportional to accumulated luminosity
- Grid resources are limited by funding and fully utilized

# Goals and tasks of the MICC engineering infrastructure:

The engineering infrastructure should provide a reliable functioning of the Complex 24 hours a day, 7 days a week round-the-year



A scheme of arranging the equipment in the computer hall on the 2-nd and 4-th floors

For the NICA project the data stream has the following parameters:

- high speed of the event set (up to 6 kHz),
- in central Au-Au collision at the NICA energies, about 1000 charged particles are generated,
- predicted event quantity - 19 billion;
- the total amount of initial data can be valued as 30 PB annually or 8.4 PB after compression.

(*MPD Conceptual Design Report v. 1.4* )

**Simulation of the distributed computer infrastructure**



**A model for studying processes has been created:**
✓**Tape robot,**
✓**Disk array,**
✓**CPU Cluster**.

# Importance of the Tier-1 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 Tier-1 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!