# PIK Computing Centre

Andrey Kiryanov

# CACS with HACS

# Power and cooling infrastructure

- Both main power and cooling can sustain up to 300 kVA of load
  - With maximum load the computing equipment can run on batteries for about 15 minutes, which is enough for a graceful shutdown
  - Cooling system has its own UPS and runs on batteries for about 30 minutes
  - Diesel powers the cooling system pumps if everything else fails
  - No guaranteed power in case of a major failure, but power line is redundant
- Currently the Computing Centre equipment is worth 120 kVA of load
  - Roughly 40 minutes of runtime on batteries
  - Over 50% of rack space is unused (12 racks populated out of 28)
  - Network infrastructure is designed for full capacity (except InfiniBand)
  - We can move in new servers without any modifications to the infrastructure

Grid'2018, 10-14 September 2018, Dubna, Russia

# Chillers and refrigerant tanks

# Computing equipment

- Peak theoretical performance is ~362 Tflops

- Real LINPACK results:
  - ~200 Tflops on Xeon CPUs (no AVX-512), effectiveness ~80%
  - ~68 Tflops on Xeon Phi (KNL) CPUs (AVX-512), effectiveness ~50%

- Computing equipment:
  - 160 nodes with Xeon CPUs: 2.4 GHz, 28 cores, 128 GB RAM per node (4.5 GB RAM per core) – **4 480 cores**
  - 40 nodes with Xeon Phi (KNL) CPUs: 1.4 GHz, 68 cores (272 virtual), 96 GB RAM per node – **10 880 virtual cores**
  - 16 nodes with Xeon CPUs: 2.4 GHz, 28 cores, 1 TB RAM + 1.6 TB NVMe SSD – **448 cores**
  - 2 nodes with Xeon CPUs: 2.4 GHz, 28 cores, 1.5 TB RAM – **56 cores**

СУПЕРКОМПЬЮТЕРЫ
tOP 50

http://top50.supercomputers.ru/

| 5 | Санкт-Петербург Суперкомпьютерный | 1468/20552 | узлов: 623 (2xXeon E5-2697v3 2.6 GHz 64 GB RAM) узлов: 56 (2xXeon E5-2697v3 [Acc: 2xTesla K40] 2.6 GHz 64 GB RAM) | 715.94 | 1,015.10 | Группа компаний РСК |

Grid'2018, 10-14 September 2018, Dubna, Russia

# Storage

- Lustre with 2.9 PB of raw disks (~2.3 PB of visible storage + 29 TB of metadata)

  - 2.10.4 LTS release + Mellanox OFED

  - Dense DELL disk shelves connected via SAS

  - Connection through 100 Gbps InfiniBand

- Ceph with 2.5 PB of raw disks (two full racks)

  - 13.2.1 Mimic release

  - Standard Supermicro disk servers

  - Connection through 2x10 Gbps Ethernet

  - InfiniBand is also available, currently used for replication

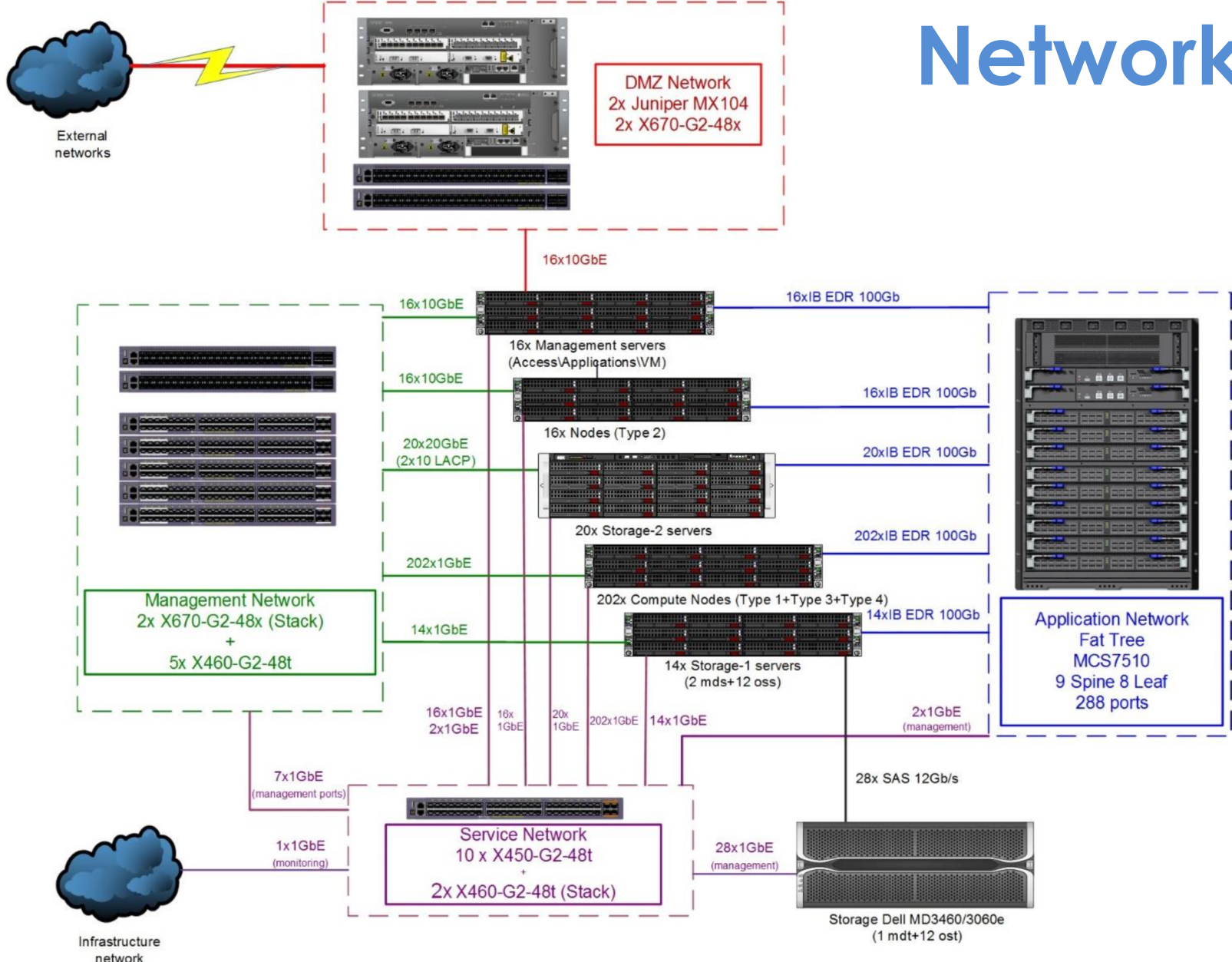Grid'2018, 10-14 September 2018, Dubna, Russia

# Lustre disk shelves and servers
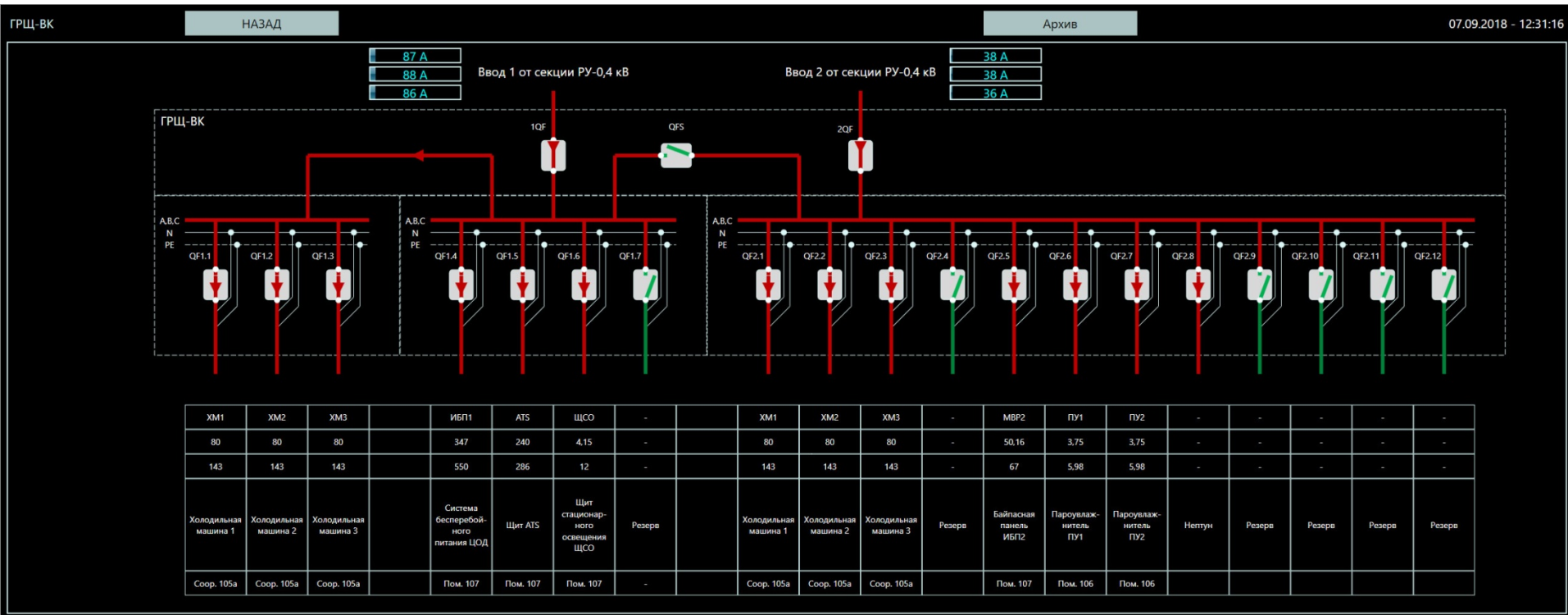
# InfiniBand switch and Ceph servers

# Engineering Systems Monitoring



Integral system covering mains, UPS, cooling, climate, leak sensors, chillers and diesel

Main screen with a bird's-eye view

Grid'2018, 10-14 September 2018, Dubna, Russia

# Engineering Systems Monitoring



## Detailed mains monitoring

Grid'2018, 10-14 September 2018, Dubna, Russia

# Engineering Systems Monitoring



Even more detailed monitoring of a single mains switchboard

Grid'2018, 10-14 September 2018, Dubna, Russia

# Software

- OS: CentOS 7, 64-bit

- High availability: pacemaker

- Management: xCAT

- Auth: FreeIPA

- Batch system: Slurm

- Cluster monitoring: Ganglia, Nagios, Zabbix

- MPI: Open MPI, Intel MPI, Platform MPI

- Compilers: GCC, Intel Compiler

- Intrusion prevention: fail2ban

- Virtualization: KVM

# User perspective

- Heterogeneous resources

- Different types of nodes are organized in distinct queues
  - "Standard" CPUs
  - Big/Huge memory
  - KNL

- User data reside on Lustre
  - Home directory is accessible from all nodes
  - Common software and shared data in dedicated areas
  - Environment modules for switching between compilers and MPI flavours

Grid'2018, 10-14 September 2018, Dubna, Russia

# High Availability

- Most critical services are running in HA mode
    - xCAT: two servers in master-slave mode
    - Lustre: every disk shelf is connected to two servers
        - Resources can be reallocated in real time, allowing for almost transparent Lustre server maintenance
    - Slurm: two servers in master-slave mode
    - MariaDB: Galera cluster with two servers
    - FreeIPA: two servers with replication
    - Ceph: multiple MON, MGR, MDS
- Automatic failover for all services except xCAT

NATIONAL RESEARCH CENTRE
«KURCHATOV INSTITUTE»

PETERSBURG NUCLEAR PHYSICS INSTITUTE

# Intrusion prevention

- We use a pool of login nodes for user access

- Key-based SSH, no password authentication

- Primary firewall on a main router enforces static access policy

- Individual firewalls on all hosts with global addresses

- Custom-compiled version of fail2ban with IPv6 support

- Ban log on a shared filesystem
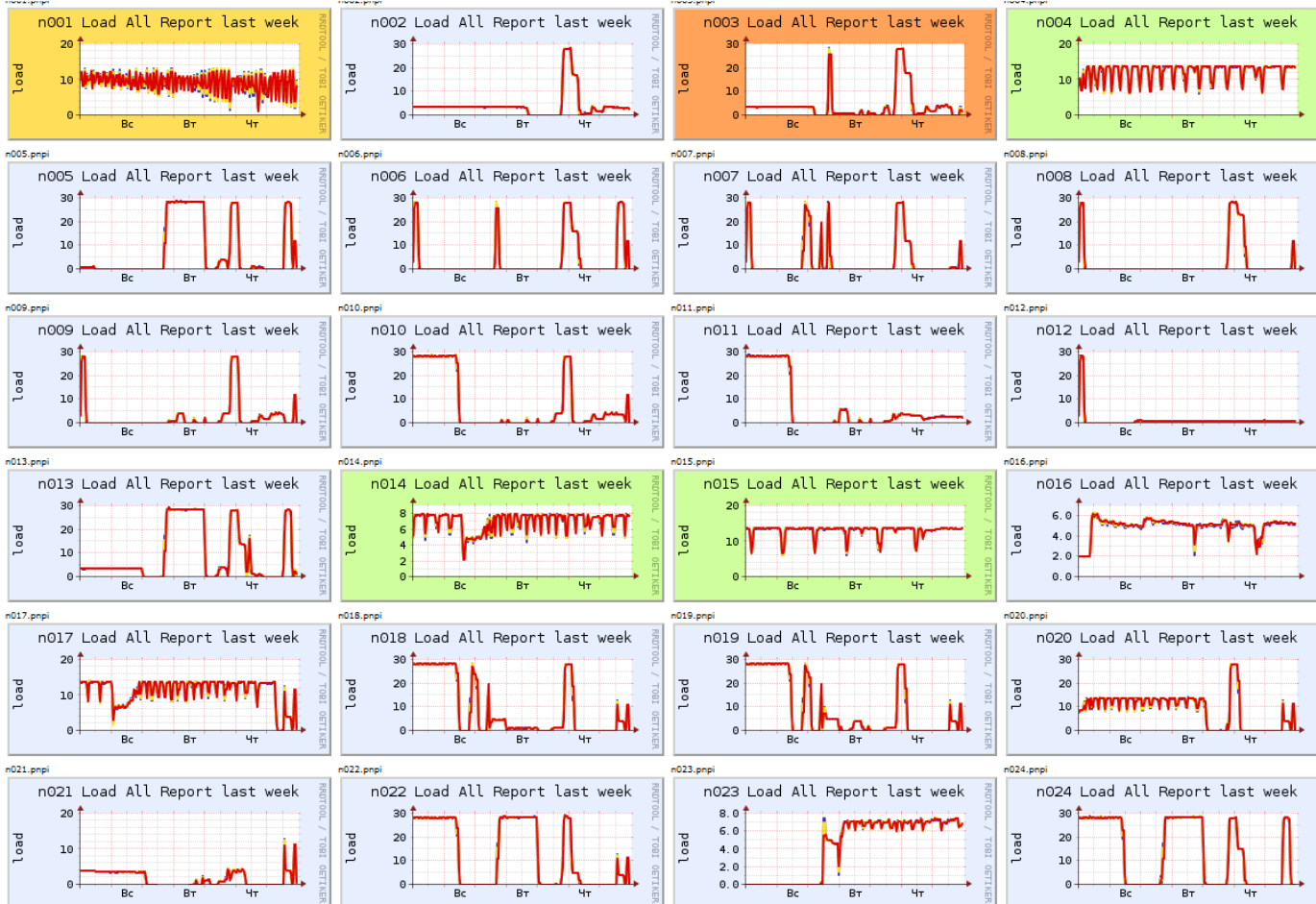
  - If two nodes ban the same IP it will propagate to the whole pool of nodes

# Virtual Infrastructure

- 14 servers with 256 GB RAM as hypervisors

- KVM-based VMs managed by xCAT

- Disk volumes on Ceph

- SR-IOV for InfiniBand and 10 Gbps Ethernet

  - Zero VM I/O overhead

  - Lustre over IB works seamlessly

  - Needed some xCAT modifications

  - Libvirtd allocates ethernet VFs automatically, but fails miserably with IB because of a longer MAC
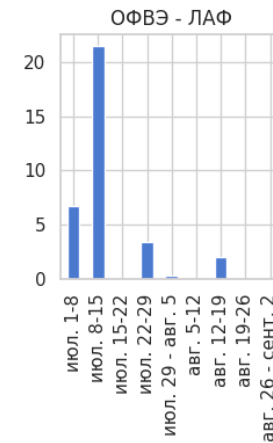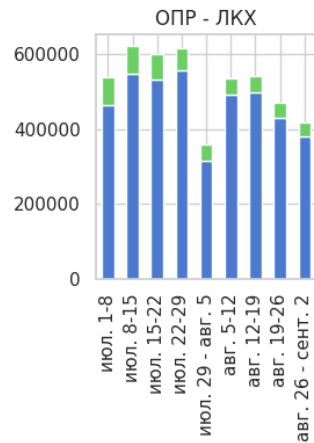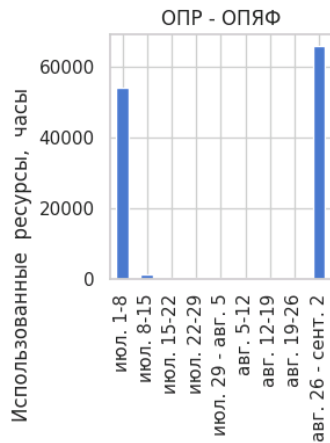
# Ganglia Monitoring



Grid'2018, 10-14 September 2018, Dubna, Russia

# Slurm Monitoring



Per-core job distribution

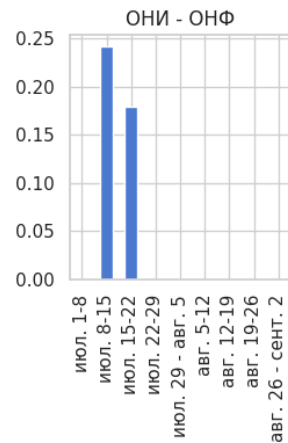Grid'2018, 10-14 September 2018, Dubna, Russia

# Reports

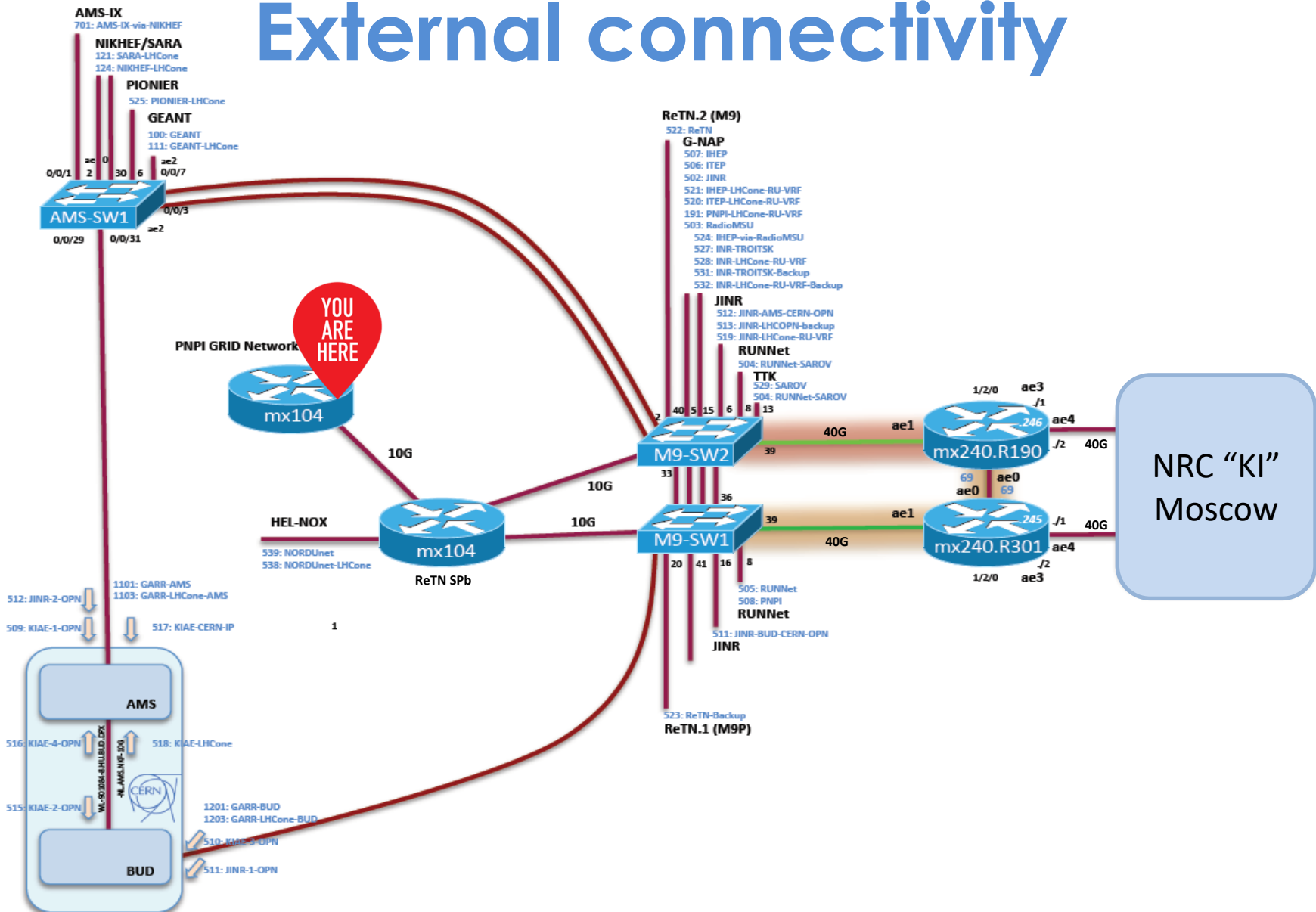Динамика использования ресурсов коллективами

По числу ядер на задание

Grid'2018, 10-14 September 2018, Dubna, Russia

# Challenges

- Different default fan modes on different servers

  - Firmware upgrades and manual tuning

  - IB card may overheat with low fan speed

- KNL performance issues

  - Firmware upgrades and memory mode tuning

- Mellanox OFED memory allocation problems

  - PR is still open with Mellanox

  - Seems to be fixed in 4.4-2 release

- Issues with Ceph EC pools

  - Solved by moving to Mimic and RHEL 7.5 kernel

# Thank you!