



# IHEP Site Report

---

Hepix Autumn 2019 -- Amsterdam

Jingyan Shi (shijy@ihep.ac.cn)  
On behalf of Computing Center, IHEP

# Outline

---



1

## Brief Introduction

2

## Operating Status

- Local Cluster
- Grid Site
- Network

3

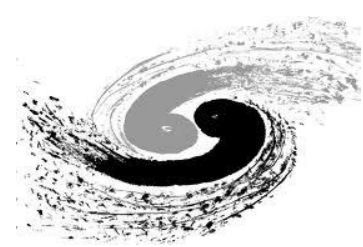
## Activities in progress

- HTCondor Cluster
- Storage
- Grid site

4

## Summary

# Brief Introduction to IHEP



**BESIII** (Beijing Spectrometer III at BEPCII)

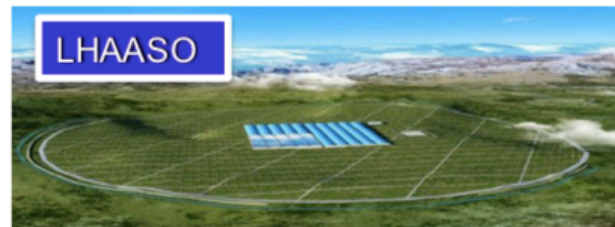


**DYB** (Daya Bay Reactor Neutrino Experiment)



**JUNO** (Jiangmen Underground Neutrino Observatory)

**YBJ** (Tibet-ASgamma ARGO-YBJ Experiments)



Large High Altitude Air Shower Observatory



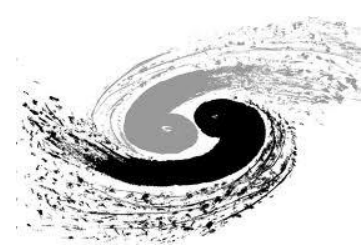
Hard X-Ray Moderate Telescope



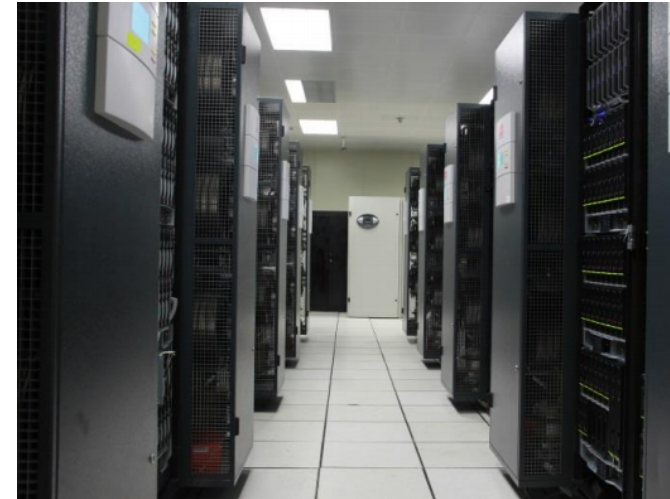
Circular Electron Positron Collider

# Computing Resources

---



- 20,000 cpu cores, 100 GPU cards to for more than 10 experiments
  - HTCondor cluster runs for HTC jobs
  - Slurm cluster runs for HPC jobs
  - WLCG tier 2 site
- About 30PB storage
  - Luster and Eos are two main file systems
  - Caster for tape storage
- Network
  - IP V4/ IP V6 dual stack
  - Ether net(100Gb) / IB (100Gb) supported
  - LHCOne joint



# Outline

---



1

**Brief Introduction**

2

**Operating Status**

- **Local Cluster**
- **Grid Site**
- **Network**

3

**Activities in progress**

- **HTCondor Cluster**
- **Storage**
- **Grid site**

4

**Summary**

# Updates to Infrastructure

---

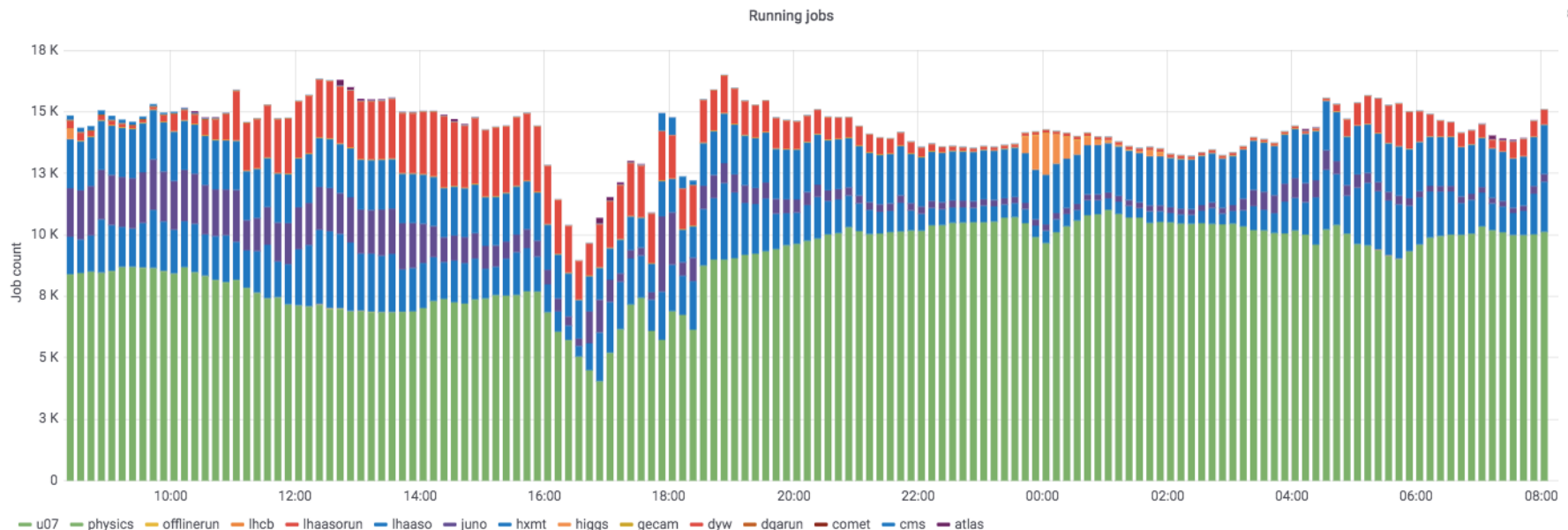
- New work node -- 6160 cpu cores
  - H3C B5700 and Lenovo SN550
    - CPU Intel Xeon Gold 6248 20 cores 2.50GHz
    - Memory 128GB
    - Disk 960GB SSD
- Kernel Upgrade (CVE-2019-11477, CVE-2019-11478, CVE-2019-11479)
  - This is an Intel CPU vulnerability
  - Upgrade kernel to version 2.6.32-754.17.1.el6.x86\_64



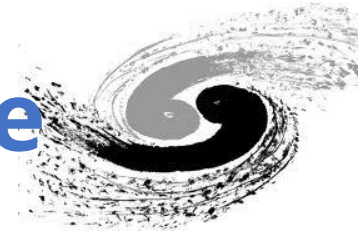
# HTCondor Cluster Status



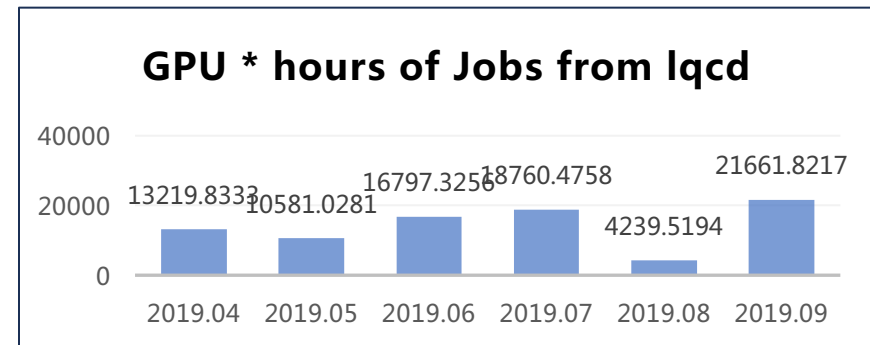
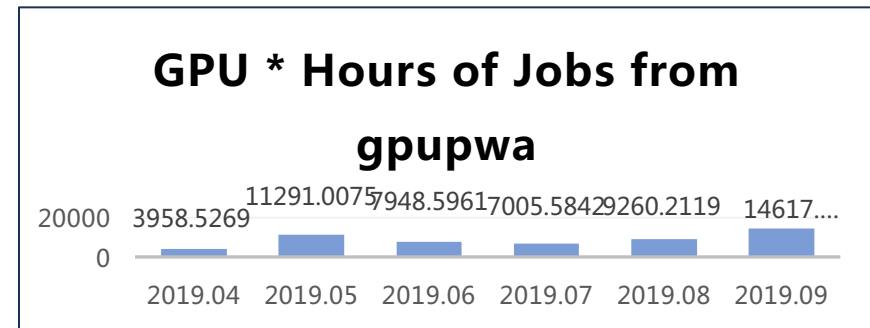
- Upgraded HTCondor to the 8.8.4
  - More stable
- Job memory limitation added
  - 2GB~4GB/job, depending on the memory the work node owned
  - Switch off swap of work node
- Totally 17,860,655 jobs and 4,802,874 hours last half year
- Job slot utilization is over 87%



# Slurm GPU Cluster - Infrastructure

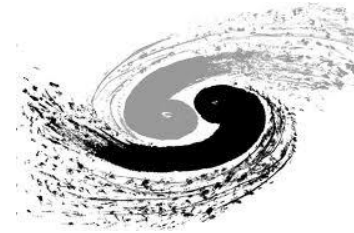


- Slurm GPU Cluster : OS SL7.5 + Slurm 18.08
  - Resource
    - 1 control node
    - 2 login nodes: 14 NVIDIA v100 gpu cards
    - 10 worker nodes : 80 NVIDIA v100 GPU cards
      - 256GB memory
      - 10Gb Ethernet and 100Gb IB connection
  - Aim at lqcd, BES partial wave analysis, machine learning etc.
  - 800TB Lustre storage

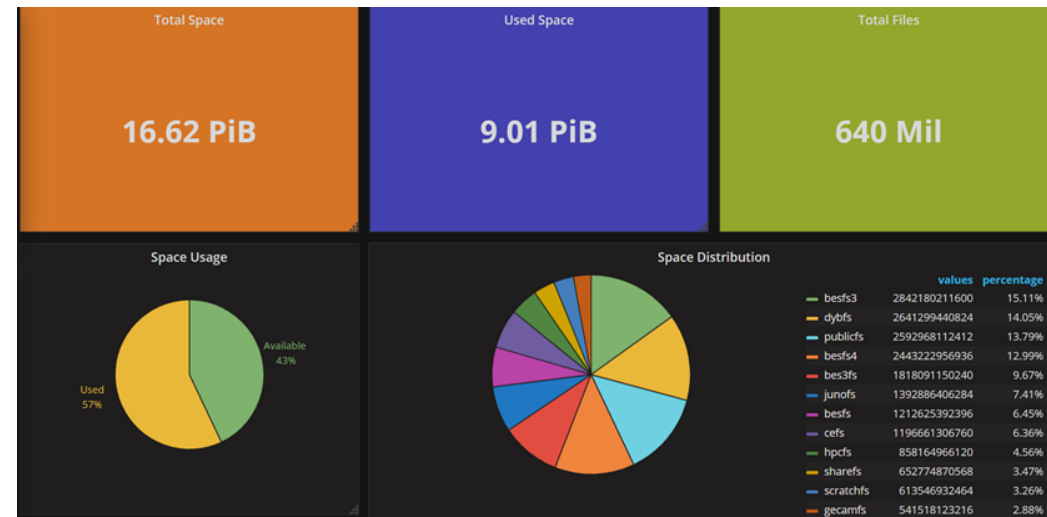
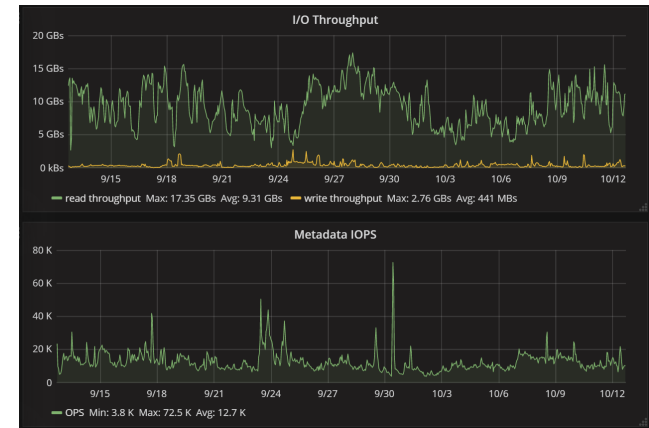




# Storage Statistics



- Space capacity
  - Lustre :17 PB total, 9 PB used, 3 PB will be added soon
  - EOS: 4 PB total, 3.3 PB used, 2PB will be added soon
- Performance -- Aggregate bandwidth
  - Read :17.35 GB/s peak, 9.3 GB/s average
  - Write :2.76 GB/s peak, 0.4 GB/s average
- Availability Time
  - >99%
- All the disk servers have been upgraded to Lustre 2.10.6
  - To support newer linux kernel(3.x) and new coming hardware
  - computing nodes are running older Lustre on top of 2.x Linux kernel



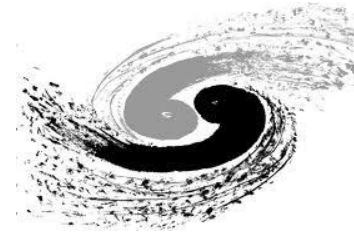
# AFS Authentication Upgrade

---



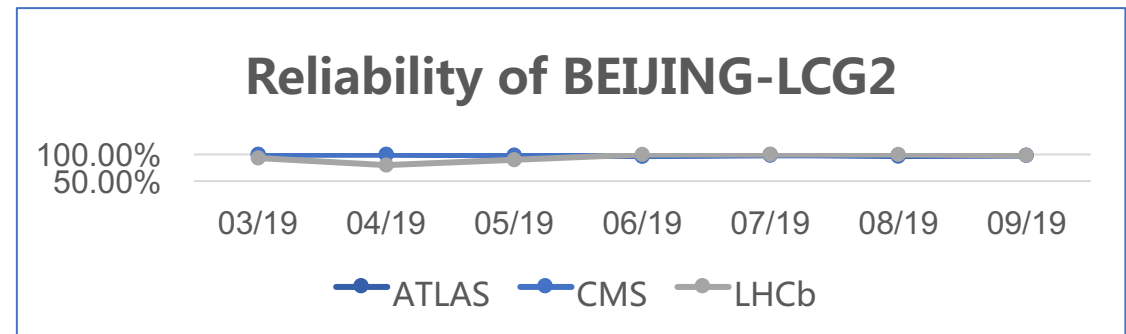
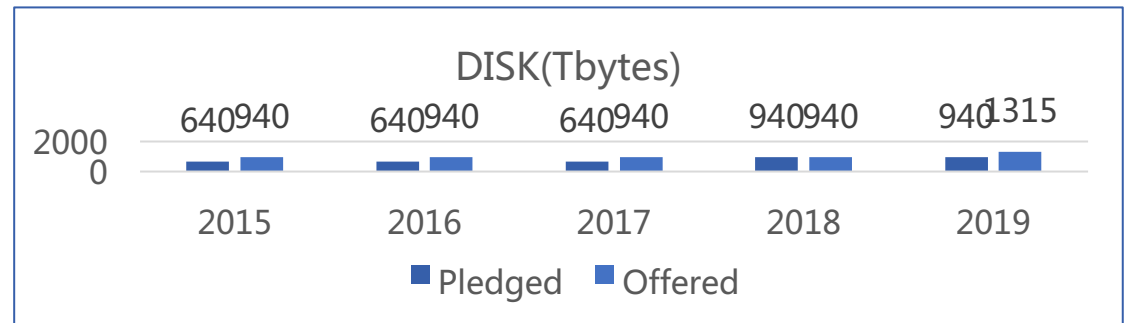
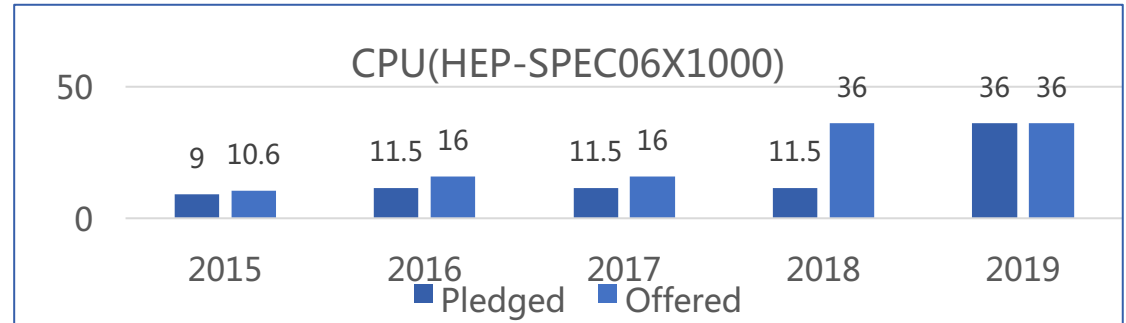
- Upgraded AFS authentication from AFS kaserver to kerberos 5
  - Improve security: AFS kaserver has weak security properties
  - Success to get tokens when login nodes
- Features
  - More flexible: Account authentication is independent from AFS file system
  - Deployed with the master/slave configuration to provide high availability of Kerberos 5 KDC service
  - Support password-free authentication in Login farm

# BEIJING-LCG2 Tier2 Resources



- CPU: 1896 cores
  - Intel Golden 6140 1008 Cores
  - Intel E5-2680V3: 696 Cores
  - Intel X5650 192 Cores
- Batch: Torque 4.2.10
- VO: ATLAS, CMS, LHCb

- DPM: 775TB
  - 4TB \* 24slots with Raid 6, 5 Array boxes
  - DELL MD3860 8TB\*60slots
- dCache: 540TB
  - 4TB \* 24slots with Raid 6, 6 Array boxes
  - 3TB \* 24slots with Raid 6, 2 Array boxes



The Site keeps a good reliability at most of the time

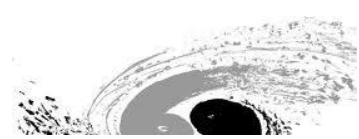
# BEIJING-LCG2 Tier2 Operations

---

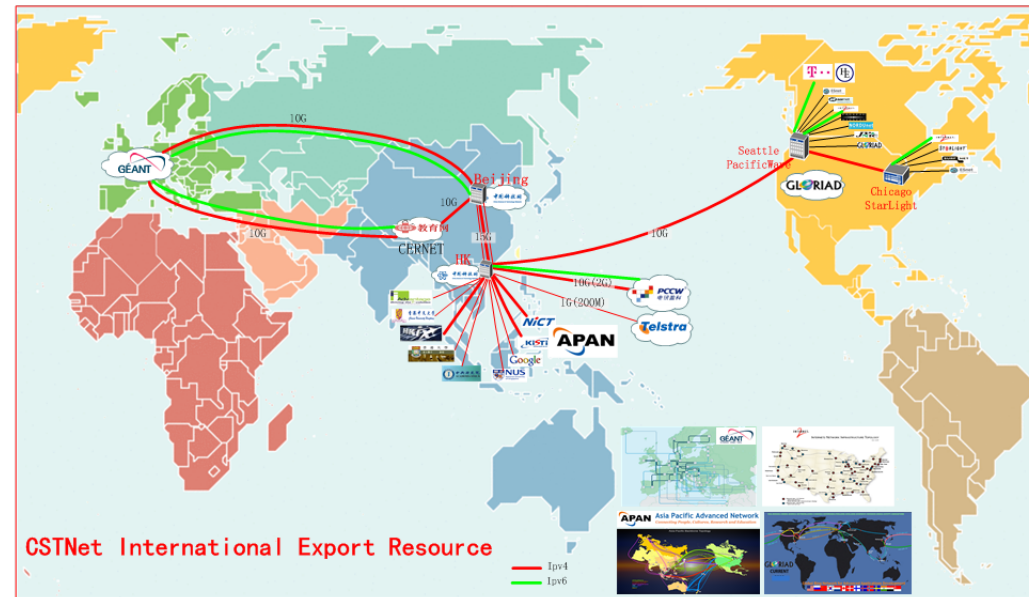
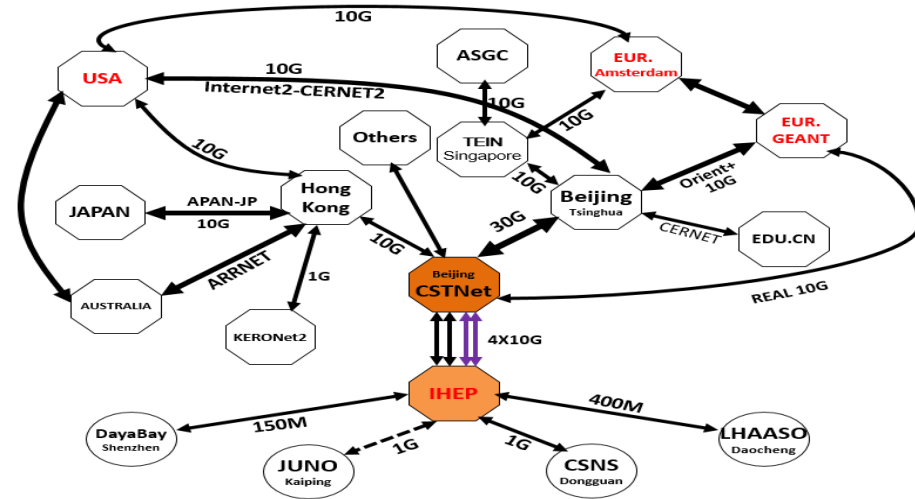


- Adding support for VO LHCb
  - Resource for LHCb : 1008 CPU cores and 360TB disks.
- Join LHCOne and enable ipv6 for data transfer.
- Upgrade servers and work nodes to Centos7
  - Develop UMD4 auto installation and configuration modules for Centos7.
  - Upgrade DPM storage element to the latest version.
- HTCondor-CE testing is under going

# Internet connection



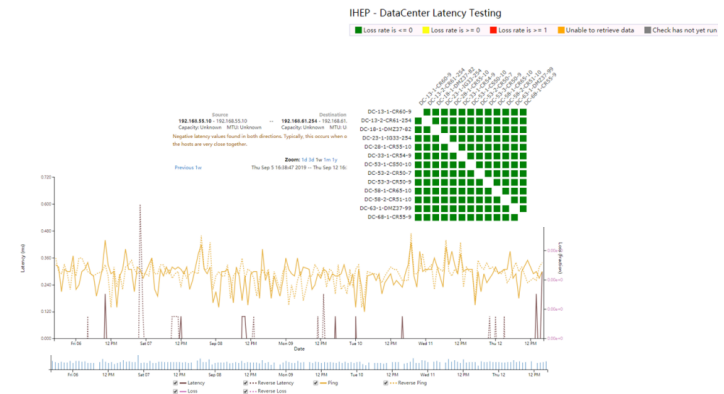
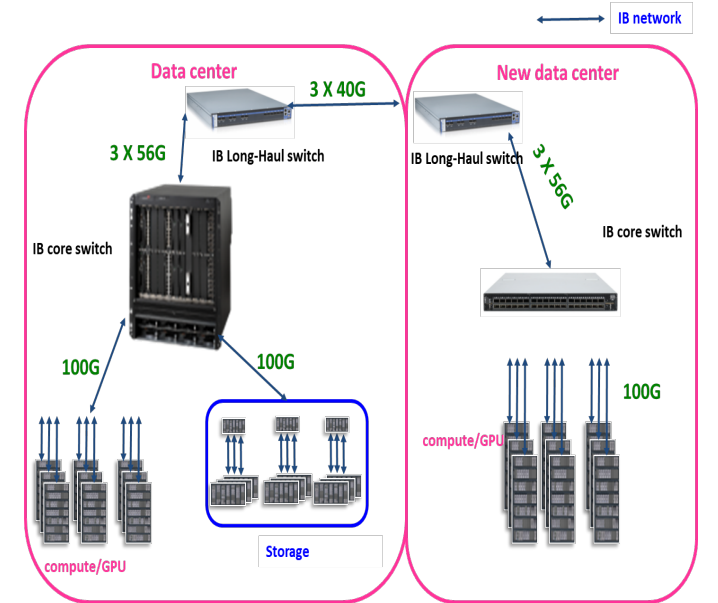
- 4 X 10G links to CSTNet
  - 2X10G for LHCONE
  - 2X10G for normal traffic
- LHCONE update at IHEP
  - New Peer to GEANT by CSTNET was finished last month
  - Route between IHEP and Europe has been changed from Orient+(CERNet - GEANT) to REAL link(CSTNet - GEANT).
  - Peer points with Internet2/ESNet by CSTNET was ready two weeks ago
  - More peers for LHCONE will be ready
    - GEANT ( Amsterdam)
    - APAN (Asia Pacific Area Network)



# Data Center Network Updates



- InfiniBand network for HPC is ready
  - 100Gbps backbone, in production
  - RDMA\_Write bandwidth:  $\approx 11675\text{MB/s}$
  - RDMA\_Write latency:  $< 0.95\ \mu\text{s}$
  - 15 IB nodes for HPC now
- 100G Ethernet for DCN
  - Upgrade the data center core switch
    - Add a 100Gb/s blade module, provide 6 100Gb/s ports
  - A new 25GE TOR is online, provide 4 X 100Gbps uplink for storage servers, whose Ethernet card is 25Gb/s
- Latency monitoring service for computing platform is online
  - The performance of internal network in Data Center is well



# Outline

---



1

## Brief Introduction

2

## Operating Status

- Local Cluster
- Grid Site
- Network

3

## Activities in progress

- HTCondor Cluster
- Storage
- Grid site

4

## Summary

# Activities in progress– HTCondor cluster

---



- Migrate from SL7 to CentOS
  - Tests are undergoing
- Plan to run all jobs in container
  - Motivation
    - SL6 doesn't not support new hardware
    - Experiments do not want to upgrade to SL7
    - Easy to dispatch job to remote site
  - Singularity container job test: SL6 and SL7 images with dedicate experiment file directories mounted
  - Plan to start with the new coming work node



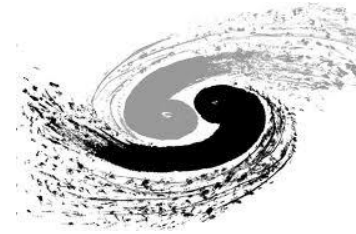
# GPU cards for LQCD Performance Evaluation

---



- Procurement for 80 GPU cards this year
  - Support LQCD
- Performance Evaluation for LQCD to run on GPU cards
  - nvswitch vs. nvlink
    - 20% performance promoted
  - IB network performance
    - 4 GPU cards/100GB IB card
  - Memory: 384 GB

# EOS + JBOD Evaluation



## ● Motivation

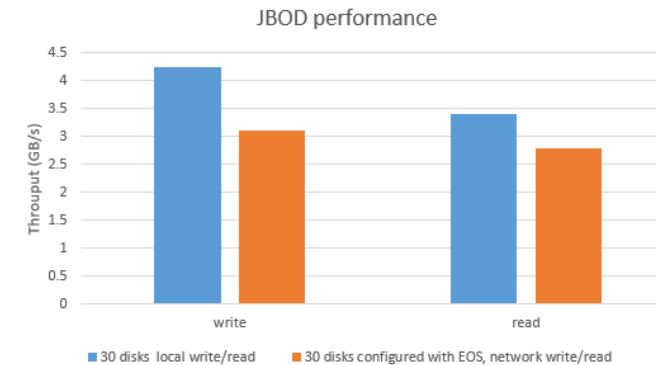
- Current RAID disk arrays can not run full bandwidth, which will become a performance bottleneck
- Use JBOD instead of RAID to provide better performance

## ● Tests on JBOD

- Preliminary tests showed that the speed was increased by about two times compared with RAID.
- A single SATA disk is basically 130MB/s, the aggregation of 30 SATA disks can reach 4 GB/sec.
- 60+ disks JBOD should be configured with two servers.

## ● Next Step

- Purchased 4 DELL ME484 JBOD arrays and 8 servers, totally raw capacity 4PB.
- Will be extended to the LHAASO EOS instance in Q4, configured with replica layout in EOS.



## DELL ME484 JBOD

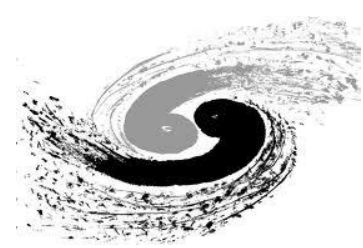
5U84 drive expansion

Direct attach for 13G and 14G PowerEdge servers

Support for direct attach SAS using 12Gb SAS HBA

# CTA at IHEP

---



- Motivation from CASTOR 1 to CTA
  - EOS is already adopted by IHEP for the disk storage.
  - EOS+CTA provides a unified interface to access disk and tape.
- Development of DB backend
  - MySQL is used widely at IHEP.
  - With the help from Steven Murray at CERN, MySQL is supported.
- Deployment at IHEP with virtual tape library
  - 3 dedicated servers are purchased.
  - The full software stack is deployed into one server using kubernetes.
- Prototype at IHEP is under preparation
  - IBM TS2900 Tape Autoloader with LTO 7 driver
- Next step:
  - setup the prototype and measure the performance.



# Summary

---



- Both computing and storage scale expanded
- Software upgrades has been done and the IHEP site keeps running smoothly
- Taking efforts to meet the requirements from the experiments
  - Container job
  - JBOD storage
  - LQCD performance



---

Thank you!

Question?