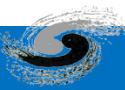# EOS status of IHEP site

Haibo Li

On behalf of Computer Center, IHEP

2018-02-05

# Contents

- IHEP Introduction

- EOS at IHEP

- Issues encountered

- Future Plan

- Summary

# Large science facilities

- IHEP: The largest fundamental research center in China

- IHEP serves as the backbone of China's large science facilities

  - Beijing Electron Positron Collider (BEPCII/BESIII)

  - Yangbajing Cosmic Ray Observatory (ASg & ARGO)

  - Daya Bay Neutrino Experiment

  - China Spallation Neutron Source (CSNS)

  - Hard X-ray Modulation Telescope (HXMT)

  - Accelerator-driven Sub-critical System (ADS)

  - Jiangmen Neutrino Underground Observatory (JUNO)

  - Large High Altitude Air Shower Observatory (LHAASO)

  - High Energy Photon Source (HEPS)

  - Under planning: XTP, HERD,CEPC …

# Major Experiments at IHEP

- BEPCII/BESIII
  - 5PB data in 5 years
- DYB
  - 400TB per year data collected
- JUNO
  - ~2PB raw data per year

lustre

- **LHAASO**
  - ~2PB raw data per year
  - accumulate 20PB+ in 10 years

EOS

- HMXT
- Atlas and CMS Tier2 site
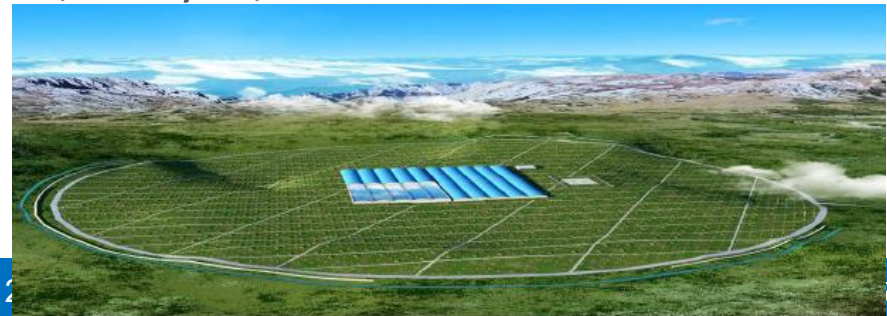  - 940TB disk, 1088 CPU cores

# LHAASO: Large High Altitude Air Shower Observatory

- Mt. Haizi (4410 m a.s.l., 29º21' 27.6" N, 100º08'19.6" E), Sichuan, China

- LHAASO Scientific Goals

  - Origin of GCRs

  - Gamma ray astronomy

  - New physics frontier (dark matter, Lorentz invariance… )

# LHAASO Computing requirements

- **~2 Petabytes** (2 million Gigabytes) of data annually generated by the LHAASO detectors
  - 1.7 PB of raw data, and >200TB of reconstruction data
  - Totally >20PB for ten years
  - Start taking data in 2018, and the data increased gradually
  - Fully completed in 2020
- **>2 Petabytes** of data generated by MC simulation
- To build one distributed computing system containing about 6000 CPU cores to process the data
  - ~ 4500 CPU cores for reconstruction, analysis, …
  - ~ 1500 cores for production

# Contents

- IHEP Introduction

- EOS at IHEP

- Issues encountered

- Future Plan

- Summary

# Storage resources at IHEP

- Local cluster

    - 11 PB+ disk storage

        - Lustre:8.5PB

        - EOS:1.7PB

        - Other:1.5PB

    - 5 PB tape storage:

        - Modified Castor 1.7.1.5

- Grid site

    - DPM:400TB

    - dCache:540TB
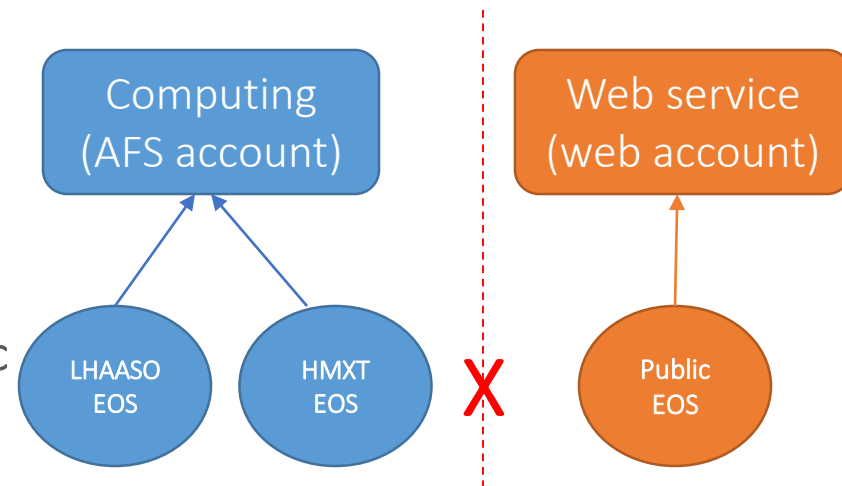
# EOS at IHEP

- 3 instances

  - 2 for experiment data storage providing computing service

    - ~1.7 PB capacity

    - ~60 million files

    - ~1.2 million directories

  - 1 for user data providing web and public services

    - Support for IHEPBox

    - ~160 TB capacity

    - ~7.4 million files

    - ~1 million directories



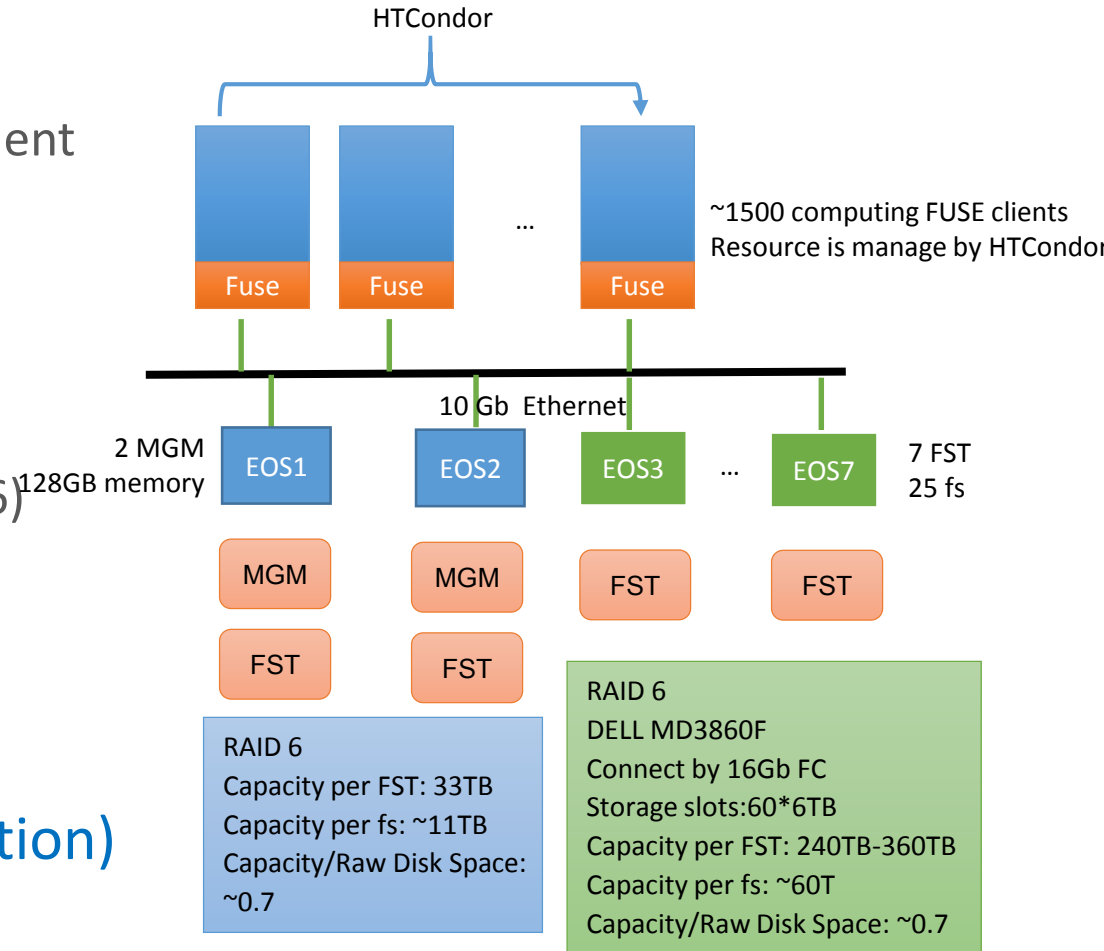Two separate cluster, based on different account system

# Experiment data storage

- ## LHAASO EOS instance

  - Used for LHAASO experiment

  - Running for 2 years

  - 0.8PB -> 1.34PB

  - server  version: 0.3.195

  - 7 dell disk array box (raid6)

    - Dell MD3860F

    - 60*6TB

  - 10Gb network link

- ## Plain mode(single replication)

- ## Fuse access

HTCondor

... 

~1500 computing FUSE clients
Resource is manage by HTCondor

Fuse  Fuse  Fuse

10 Gb  Ethernet

2 MGM
128GB memory

EOS1  EOS2  EOS3  ...  EOS7

7 FST
25 fs

MGM  MGM  FST  FST

FST  FST

RAID 6
Capacity per FST: 33TB
Capacity per fs: ~11TB
Capacity/Raw Disk Space: ~0.7

RAID 6
DELL MD3860F
Connect by 16Gb FC
Storage slots:60*6TB
Capacity per FST: 240TB-360TB
Capacity per fs: ~60T
Capacity/Raw Disk Space: ~0.7

# Experiment data storage

- ## HXMT EOS instance

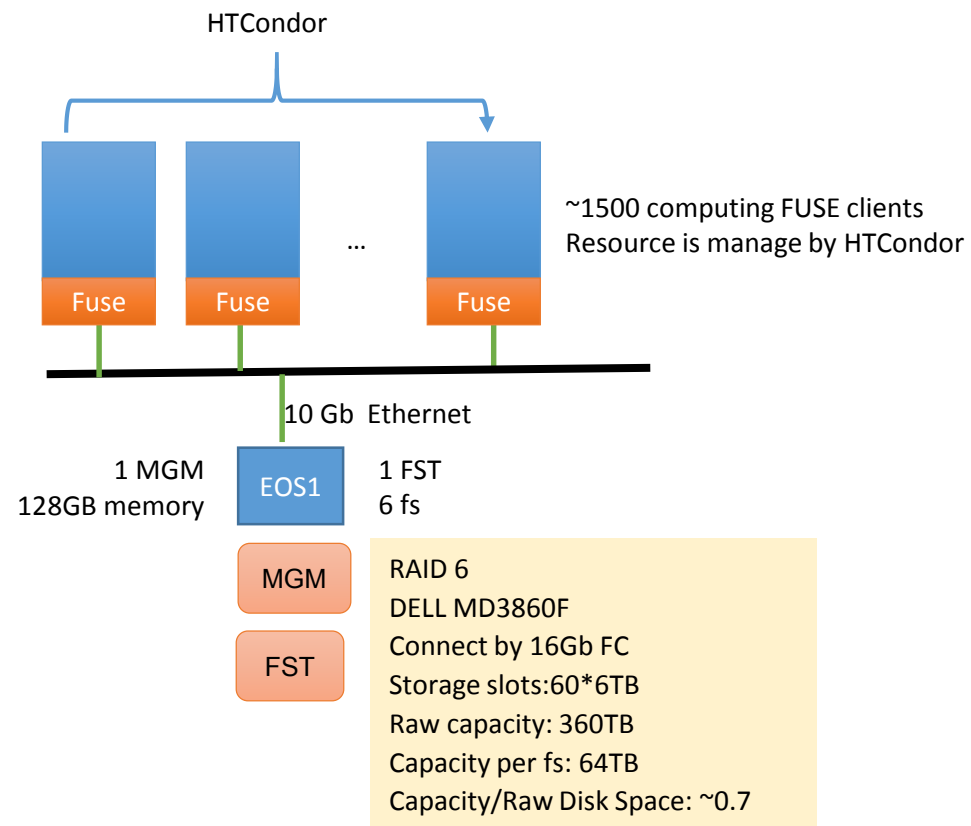  - New built for HXMT experiment in Dec 2017

  - 330TB capacity, 64TB used

  - Server version 4.2.7

  - 1 dell disk array box(raid6)

  - 10Gb network link
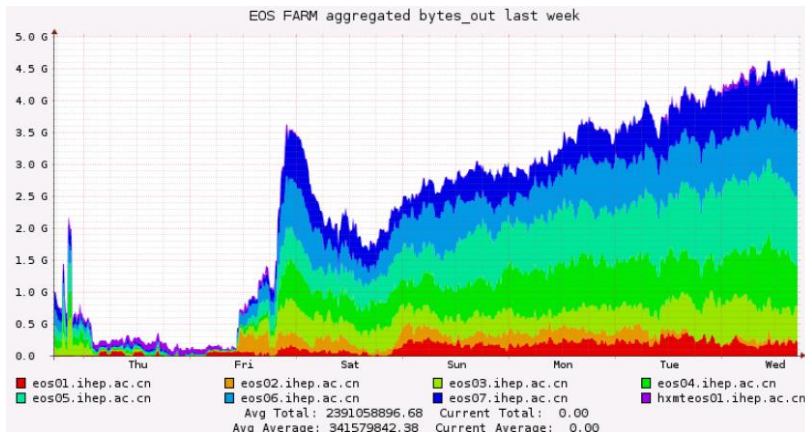
- ## Plain mode(single replication)

- ## FUSE access

- ## Many small files

  - ~4 KB per file

HTCondor

~1500 computing FUSE clients
Resource is manage by HTCondor

Fuse    Fuse    ...    Fuse

10 Gb  Ethernet

1 MGM          EOS1          1 FST
128GB memory                 6 fs

MGM

FST

RAID 6
DELL MD3860F
Connect by 16Gb FC
Storage slots:60*6TB
Raw capacity: 360TB
Capacity per fs: 64TB
Capacity/Raw Disk Space: ~0.7

# EOS running status

| Total Space | Free Space | Number of Files | Number of Directories |
|---|---|---|---|
| 1.500 PB | 870.21 TB | 48 M | 1.8 M |



EOS FARM aggregated bytes_out last week

eos01.ihep.ac.cn  eos02.ihep.ac.cn  eos03.ihep.ac.cn  eos04.ihep.ac.cn
eos05.ihep.ac.cn  eos06.ihep.ac.cn  eos07.ihep.ac.cn  hxmteos01.ihep.ac.cn

Avg Total: 2391058896.68  Current Total:  0.00
Avg Average: 341579842.38  Current Average:  0.00

- Read peak: 4.5GB/s

- Peak values allowed by the environment (mainly 5 FST each has 10Gb Ethernet)

# User data storage

- IHEPBox use-case

  - IHEPBox is based on owncloud integrated with EOS
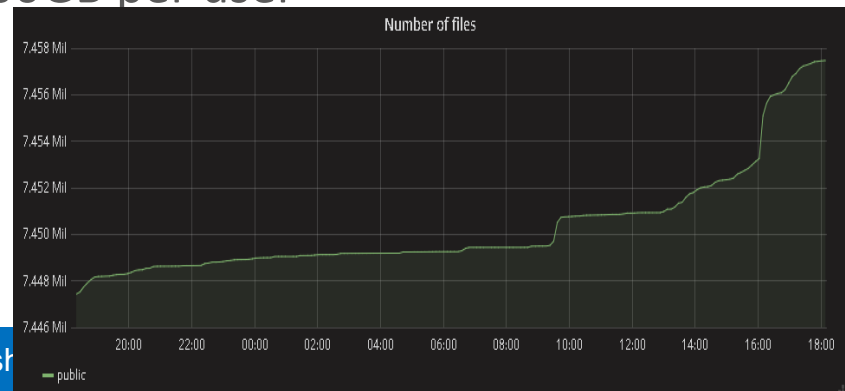
- Fast growing service

  - Integrate some plug-in to interact with IHEPBox , such as documents\galleries\calendar\activity…

  - Large and growing user base

    - 4K+ users, 1K+ active users,100GB per user

    - 160 TB ,used 21TB

      7M+ files
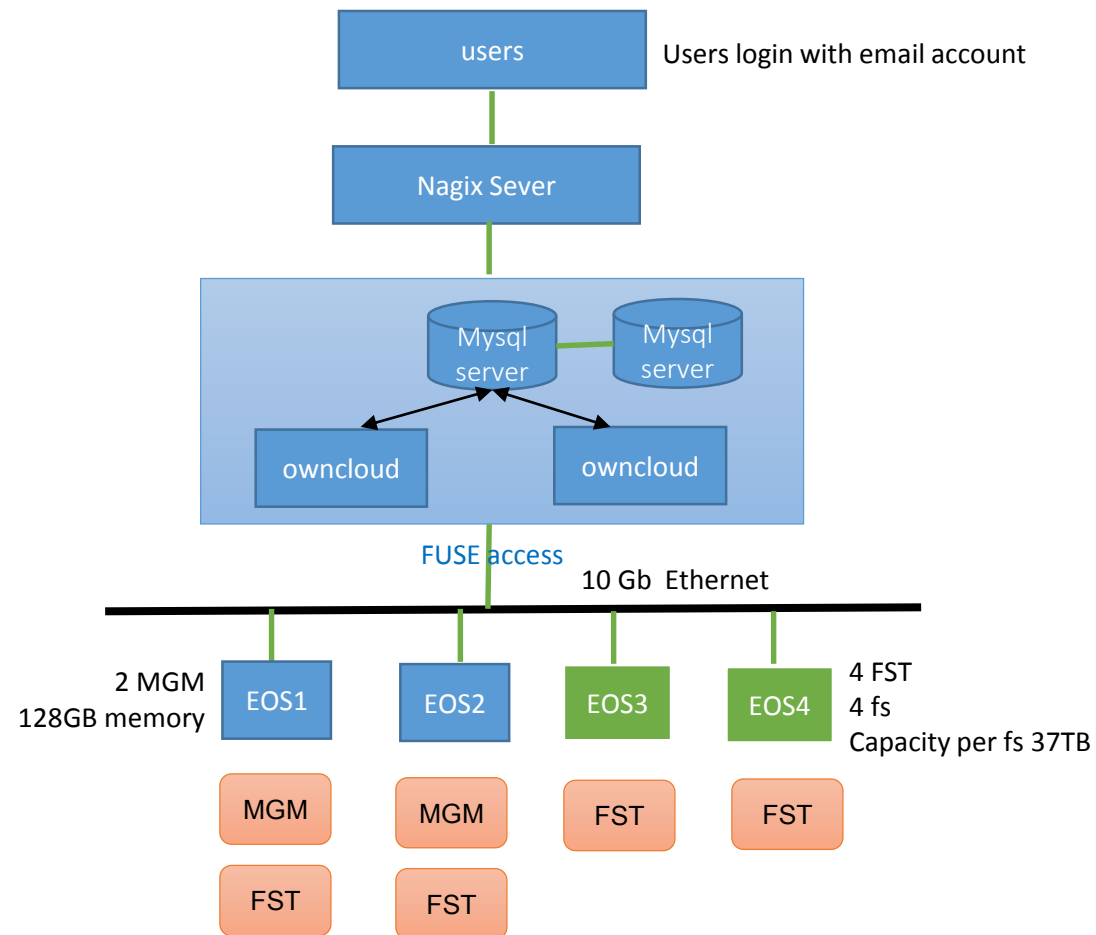
# IHEPBox

- Users
  - IHEP LDAP users
  - Local users

- Application Service
  - 1 * Nginx
  - 2 * IHEPBOX Servers

- Data Storage
  - 4 * fst
  - Two replicates
  - 2 * Mysql servers



users — Users login with email account

Nagix Sever

Mysql server — Mysql server

owncloud — owncloud

FUSE access

10 Gb Ethernet

2 MGM
128GB memory — EOS1 — EOS2 — EOS3 — EOS4 — 4 FST 4 fs Capacity per fs 37TB

MGM — MGM — FST — FST

FST — FST

# IHEPBox demands from users

- Demand for supporting more office documents online editing, such as .docx \ .ppt \.excel

- Demand for supporting one-way synchronization rules, which means files don't delete on client while the delete operation is initially by server.

# Contents

- IHEP Introduction

- EOS at IHEP

- Issues encountered

- Future Plan

- Summary

# Issues encountered

- Most issues are related to fuse.

- Some of them has been solved with the help of EOS team.

# Fuse crash

- "transport endpoint is not connected" error

  happens frequently on client.

  - Eosd Client :4.1.27

  - Server:0.3.195

- Solution: the eosd client plans to update to 4.2.12

# Nodes Held

- When a large number of jobs (using fuse to access data in eos) are submitted, the nodes managed by HTCondor sometimes become Held state.

- The possible reason is HTCondor detected the /eos is unavailable in that situation.



- Next: Do more stress tests.

# Concurrent reading error in Fuse client

- When the number of jobs exceed 11 in one node (16 cores), the job will fail with "bus error" or "segment fault". While it runs well in lustre.

- Job type: software repository, read data

- Server version: 4.2.7

- Client version: 4.1.27

- The reason is unknown now. Next we will do more tests.

# Output file synchronization problem

- When a file is being written on node1, the other node cannot read the data until the file is closed on node1.

- As a result
  - The parallel computing jobs will not run correctly.
  - Users can not read the output files until the jobs complete, which will mistaken users that the job did not run or failed

- The problem exists in all versions of EOS eosd.

- FUSEX will fix this bug, expecting.

# Atime lost

- Access time (atime) is not recorded in namespace

- The latest access time cannot be found.

- It is needed in some cases, such as statistic the access frequency of the files.

# RAIN issue

- In the rain mode, if a data block is damaged, the block can not be repaired automatically.

- The needs for using it as RAID6.

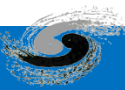- The recommended configuration now uses the replica pattern with JBOD.

# Large-small sites in storage federation

- If the master site has a large capacity, while the slave site has a small capacity, how to use EOS build storage federation?

- The replication mode may not suitable, because what data the remote site need may be uncertain.

- Is there similar needs on other sites?

# Contents

- IHEP Introduction

- EOS at IHEP

- Issues encountered

- Future Plan

- Summary

# Future plan

- Add 1 PB capacity this year

- The OS of computing node  will migrate to Centos 7

- Provide SWAN for IHEP users based on EOS in 2018

- Use EOS to build Storage federation between Beijing and Chengdu, the distance is ~2000 km and network latency is about 60ms.

# Contents

- IHEP Introduction

- EOS at IHEP

- Issues encountered

- Future Plan

- Summary

# Summary

- EOS has been deployed for 2 years in IHEP.

- LHAASO mainly use EOS as storage systems, 1PB capacity will be increased in 2018.

- We will support more experiments in IHEP, but have to do more performance stress tests.

- The functionality and stability of FUSE are critical.

- More manpower on the monitoring and operation.

- More attempts will be conducted based on EOS, such as storage federation, supporting more applications.

# Thanks for your attentions!

# 谢谢！