# Running Local Cluster Jobs at Remote Sites

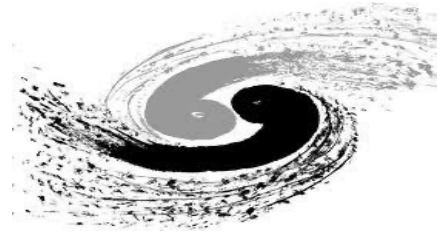Jingyan Shi

shijy@ihep.ac.cn

On Behalf of Computing Center, Institute of High Energy Physiscs

# Outline
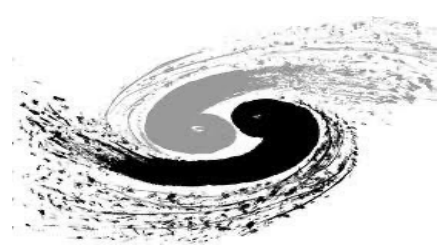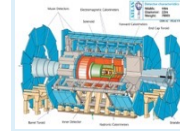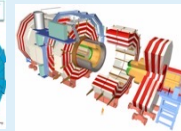
# Brief Introduction to IHEP

- The largest fundamental research center in China with research fields:
  - Experimental Particle Physics
  - Theoretical Particle Physics
  - Astrophysics and cosmic-rays
  - Accelerator Technology and applications
  - Synchrotron radiation and applications
  - Nuclear analysis technique
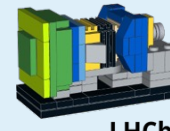  - Computing and Network application
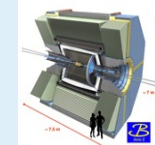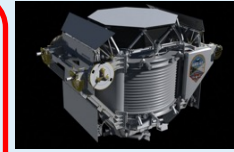
Grid

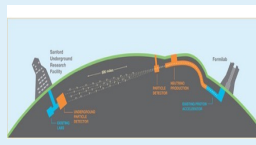**International collaboration**

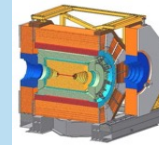ATLAS    CMS    **LHCb**    BELLE II    AMS02    DUNE
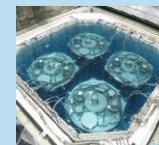
Particle Physics experiments
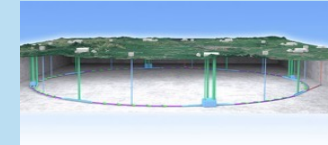
HTC cluster

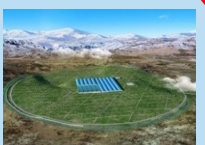**BESIII**    DYB    JUNO    CEPC    LHAASO

Particle Physics experiments

**IHEP Leading**

AliCPT    ASγ    HXMT    GECAM

Cosmic ray and astrophysics experiments

HPC cluster

CSNS    BSRF    **HEPS**

Neutron Source and Synchrotron Radiation Facilities

# Brief Introduction to IHEP-CC

- **Provide large-scale scientific computing environments for the HEP experiments**
  - facilities
  - computing
  - storage
  - network
- **Research on computing technologies to benefit high energy physics research**
- **2 Region centers**
  - North Region Center in Beijing (~45k cpu cores , ~80PB storage)
    - High Throughput Computing
    - High Performance Computing
    - Lustre file system
    - EOS file system
    - Tape Library
    - Tier 2  grid site of WLCG: WLCG grid middle ware deployed for the international collaboration
  - South Region Center in Dongguan
    - High Performance Computing
    - OceanStor9000 support by domestic vendor
    - Cloud Computing



Tianhe -2

Scientific Cloud

Commercial Cloud

Collaboration Resources

**Particle physics**

**Cosmic physics**

North(Beijng)Region Center

South (CSNS)Region Center

**Space Astronomy**

**Multidisciplinary intersection**

**Accelerator Design**

**Theoretical physics**

# Motivation

- IHEP-CC would like to be the central site for domestic HEP Exp. computing
  - HTCondor cluster of IHEP-CC is the main place of offline data process
  - A "HEP Job Tool" based on HTCondor API developed for user to simplify user job management
    - Examples: hep_sub job.sh # no submit file needed；
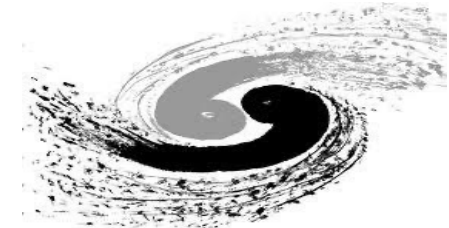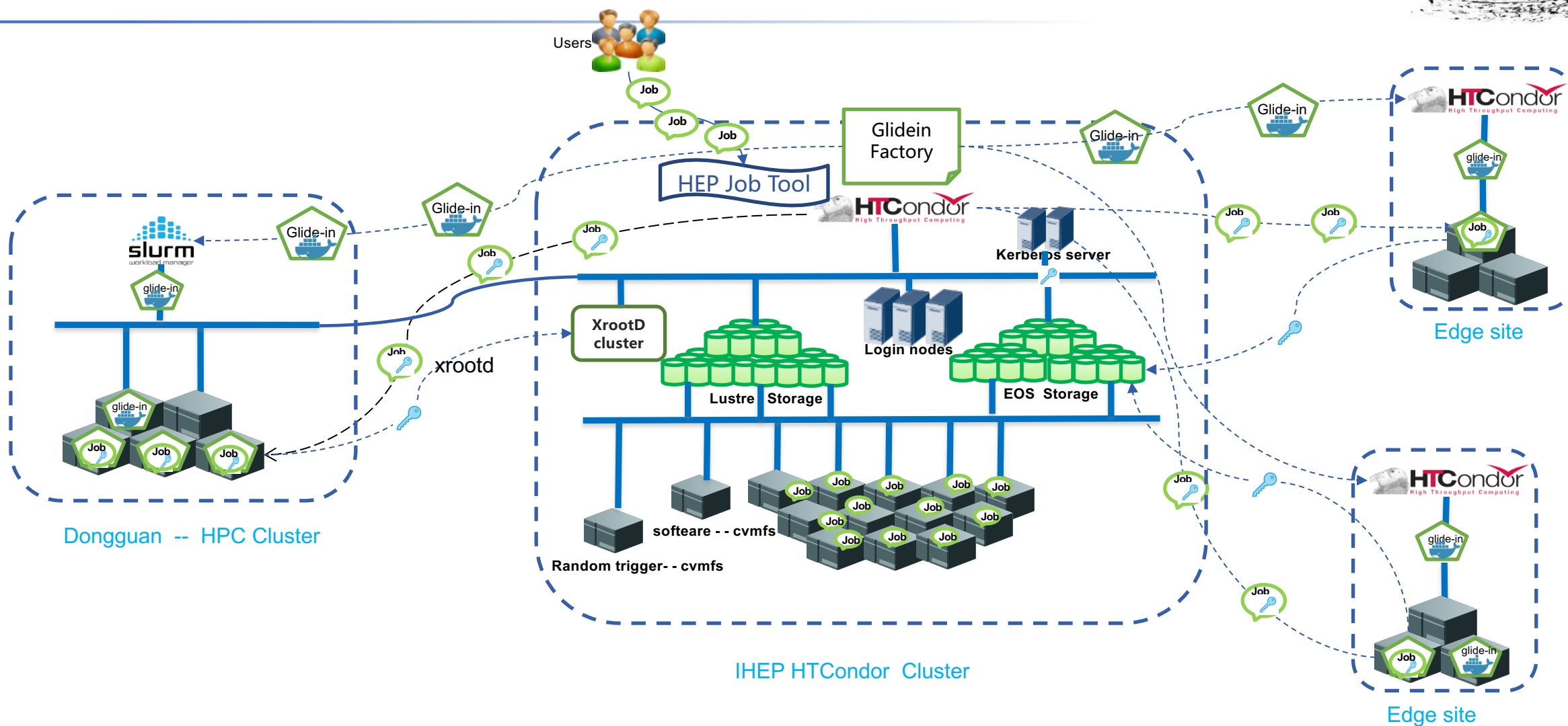      - hep_q jobid
    - Our user has become accustomed to the "cluster way" rather than the "grid way".
    - CPU resources remain highly utilized, resulting in significant job queuing time
  - Jobs running at local HTC directly access data files stored on the public file system
    - Lustre file system
    - EOS file system
- Remote sites
  - Dongguan offers 8k CPU cores and 10k arm CPU cores
    - 20Gbps dedicated link between Dongguan and IHEP → big data center
    - No storage space provided to IHEP
  - Edge site: collaboration member sites → small scale
    - Limited network connection without stable storage
  - Some super computing center is the potential resource provider
  - No extra manpower to maintain grid site at remote site
- Try to expand IHEP local HTCondor cluster to the remote site
  - Keep "the IHEP HTCondor cluster" way for the user

# Issues to be Resolved

- Remote resources could be added to IHEP HTC cluster elastically and transparently
  - Drawing upon the "glidein" concept of grid, add remote resources to the IHEP LOCAL CLUSTER
- Scheduled the suitable jobs to the suitable remote worker node
  - Job and site classification
    - Find jobs with less IO and more cpu load
    - Tag site with size, stability and network link
  - Advanced scheduling algorithm
    - Send the necessary "glidein job" to the remote sitse from one "factory" based on the status of IHEP HTCondor cluster
- User authentication from remote ends
  - IHEP cluster authentication is based on kerberos
  - Kerberos token is the key to authenticate user from the remote site
- Access IHEP data file from remote ends
  - Transfer necessary files to/from the worker node's local disk via Xrootd protocol

# Design

# Site and Job Classification

- Site's characters
  - Scale: Big / small
  - Status: Stable / unstable
  - Network bandwidth： good / bad
- Cooperated with Exp. to classify and evaluate the job based on I/O and CPU load
- HEP Job Tool analysis the job and classify it as the attribute of the job
- Scheduler schedule the suitable jobs to run at the suitable remote worker node

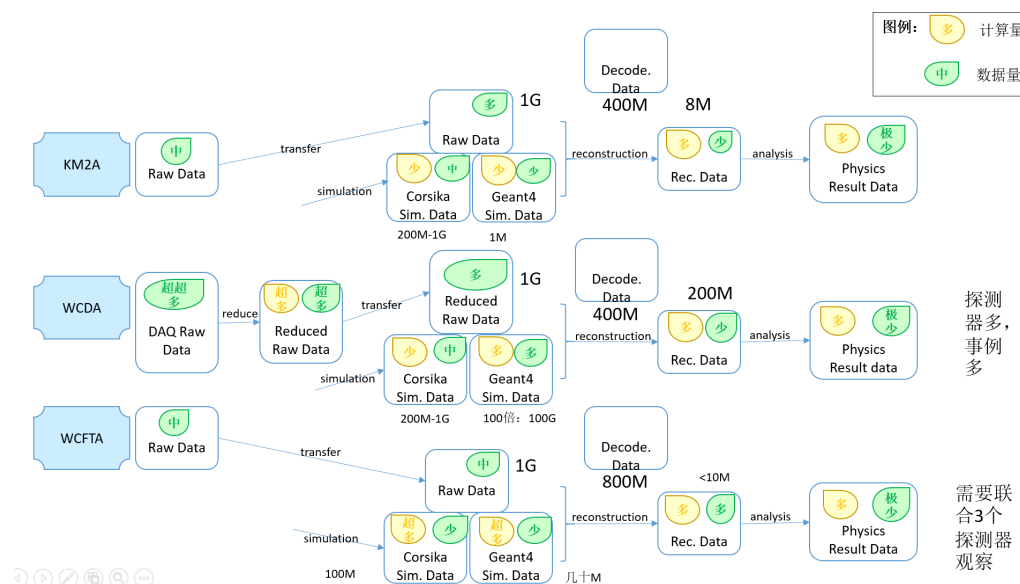- Site Classification
  - Region site – North/South region center
    - Long term and stable computing, storage resource provided
    - Good network connection
  - Edge site – collaboration member
    - No pledge storage
    - Limited network bandwidth
  - Temporary site – commercial resource
    - Short term resources for the peak usage

Scheduler

Tag and Match
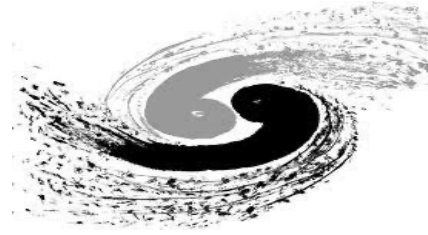
- Job Classification
  - Different Exp. Job has different characters
    - The percentages of simulation, reconstruction are different
    - IO requirement are different
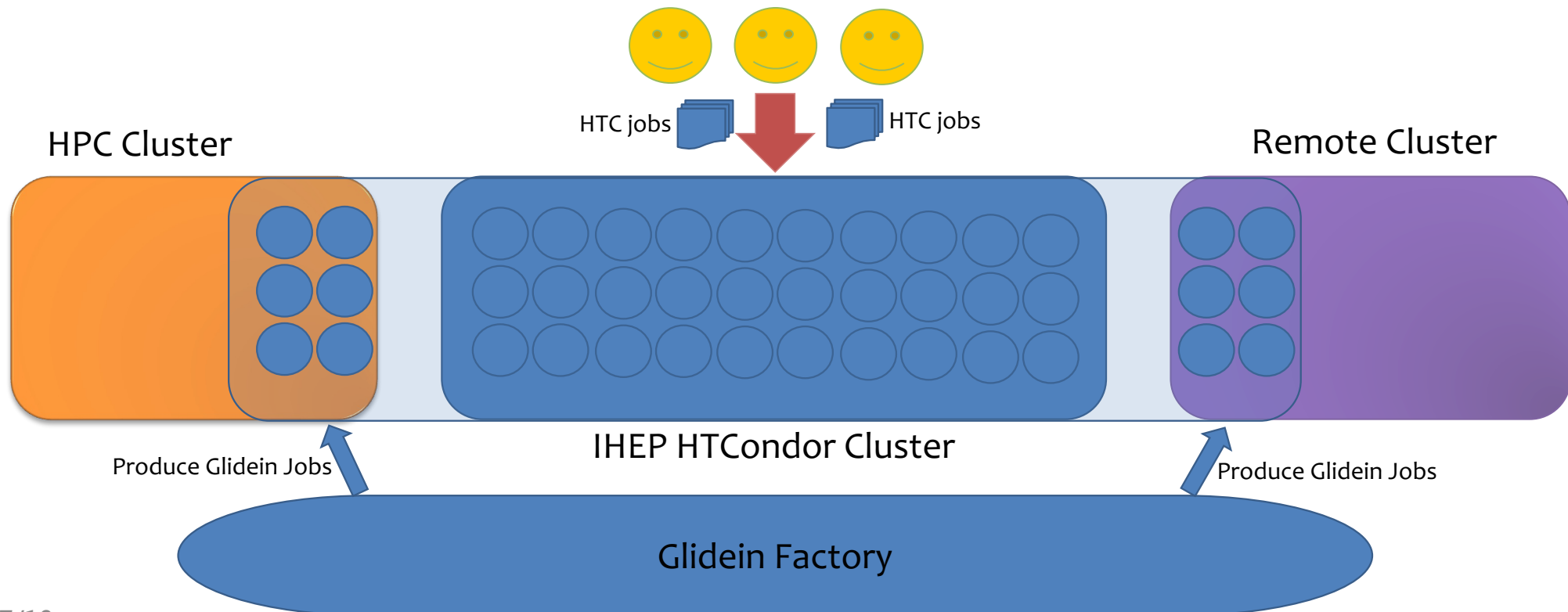    - Temporary storage requirement
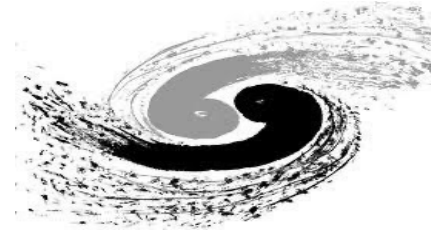


LHAASO job classification i
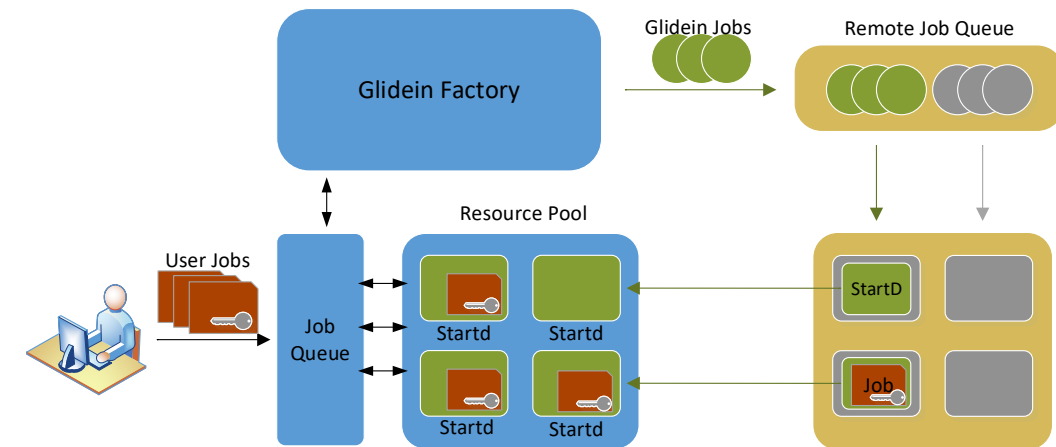
# Cluster expansion -- glidein

- One Glidein factory produces glidein job and submit it to the the remote cluster
  - Running "startd" of IHEP HTCondor cluster
  - Make the user keep the original way of IHEP HTCondor Cluster

HTC jobs                                    HTC jobs

HPC Cluster                                                      Remote Cluster

IHEP HTCondor Cluster

Produce Glidein Jobs                                            Produce Glidein Jobs
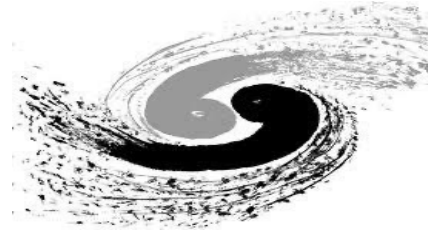
Glidein Factory
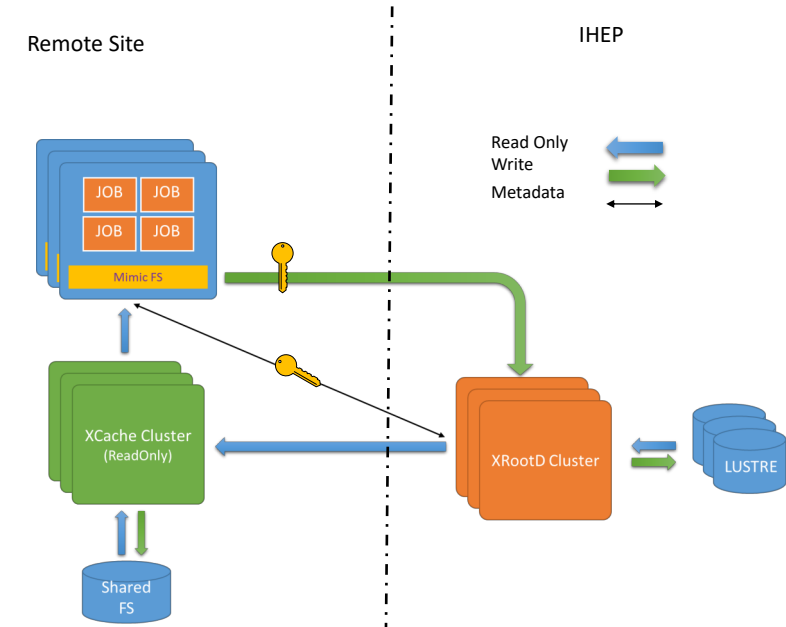
# Token based Authethication

- **HTCondor Service auth.**
  - CLAIMTOBE → IDTokens
- **User auth. from remote end**
  - User authentication at IHEP is performed via Kerberos.
  - Tried the way of Kerberos token auth. using HTCondor
    - Unsupported: user namespace inconsistency between the submission side and the execution side
  - Developed an automatic process of kerberos token
    - Token transmission, token lifetime prolong, token destroy
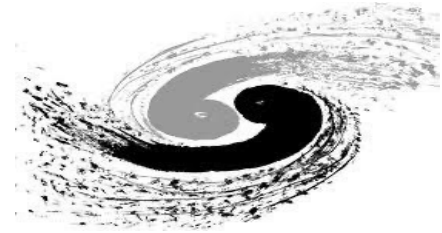    - Token ticket is initialized and valid inside job wrapper

# Data File Access from Remote End

- Necessary files are found and transferred
  - Doe at remote worker node
  - File access need to be authenticated by token
  - File transferred to the local worker node disk via Xrootd protocol
    - Add XRootd cluster in front of lustre file sytem
    - EOS support XRootd naturally
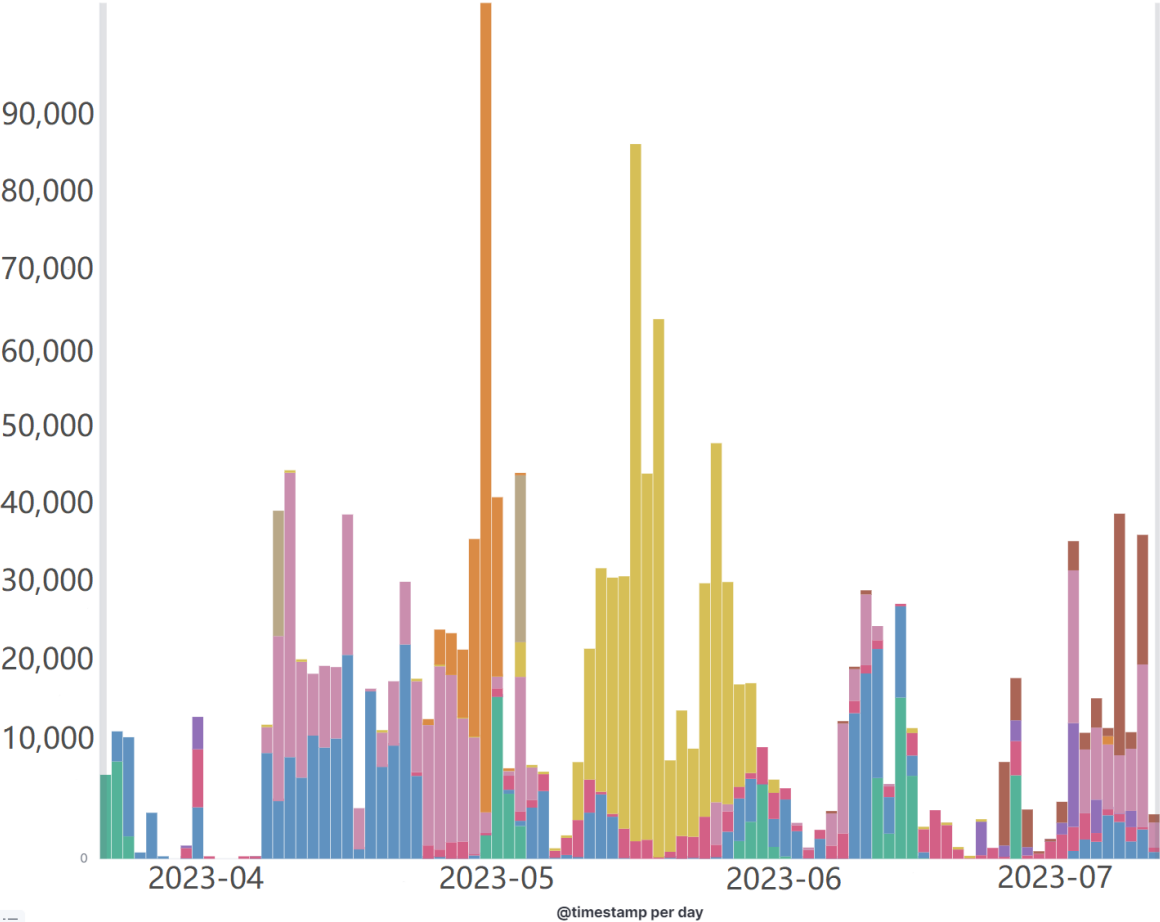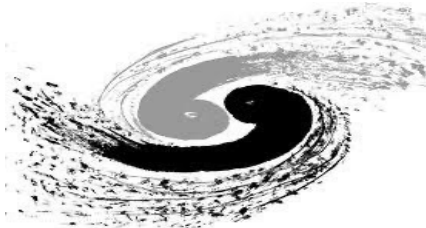- No public storage needed at remote site
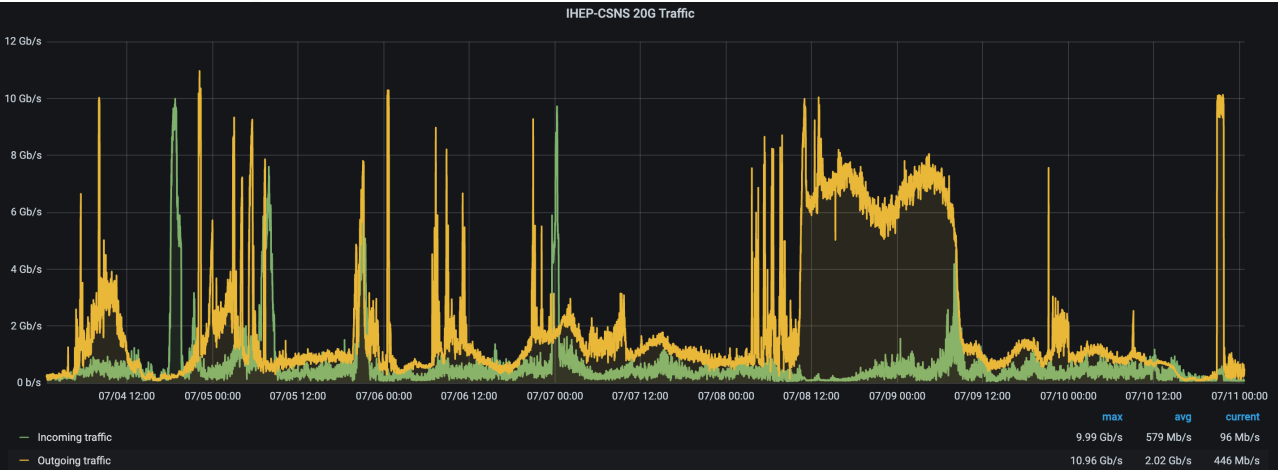
# Current Progress

- Focus on the expansion to 8k intel CPU cores and 10k arm cpu cores at Dongguan
  - HEP exp. job classification has been done
  - Automatic kerberos token process has been developed and deployed
  - File transfer optimization via Xrootd has been done
  - No glidein factory till now but static glidein job running at Slurm cluster at Dongguan
- Some of the two HEP exp. Jobs are scheduled to Dongguan from IHEP HTCondor cluster
  - BES: Simulation and reconstruction jobs are **transparently** scheduled to run at Dongguan
    - Replace the pat at fixed job option template with the workernode local disk file path
    - Completely transparent to users → user does not care where his job runs
  - LHAASO: Corsika and simulation jobs are scheduled to run at Dongguan
    - User submit the job with –rmt
    - Could be run on both Intel and Arm machine at Dongguan
- Small sites test has been done

# Running Status since 2023-04
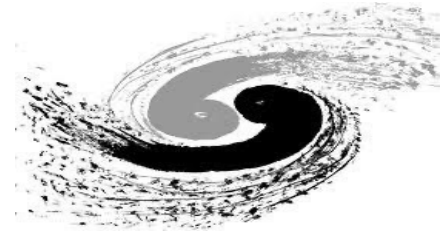


Traffic of IHEP-Dongguan network link

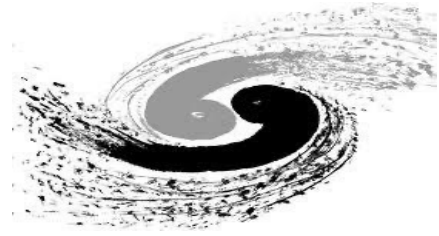****BES and LHAASO jobs has been run at Dongguan last 3 month

# Next Plan

- Glidein factory with scheduling algorithm

- More exp. Jobs and more sites would be added

- Performance optimization
  - Compress the size of the files to be transmitted
    - User's code and lib would be transferred once in one Cluster job
  - Xrootd performance evaluation and optimization
  - Big input file issues
    - Cvmfs or XCache ?

# Summary

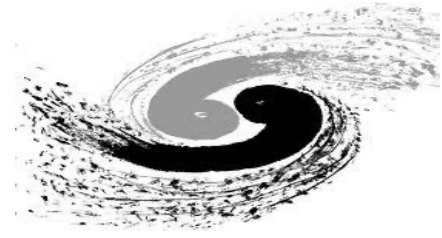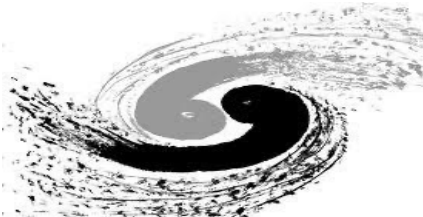- Expand IHEP HTCondor cluster rather than grid to the remote site
  - Run IHEP HTCondor cluster as the central site and expand it
  - Two HEP exp. jobs have scheduled to run at Dongguan site from IHEP HTCondor cluster
- Still a lot of works need to be done
- Thank the HTCondor team for their assistance, and a special thanks to Greg for providing us with a wealth of technical support and guidance

# Questions and Comments?
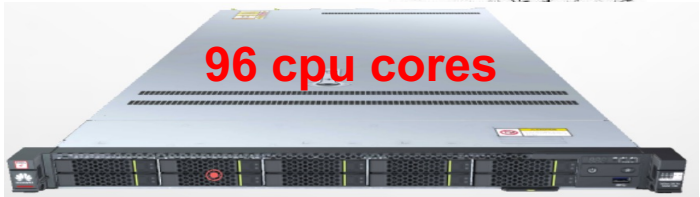
# Backup – Arm Evaluation



Kunpeng 920



Intel 5218

- **Evaluation**
  - LHAASO Corsika job
  - ARM : Kunpeng 920@2.6 GHz, 48*2 cpu cores
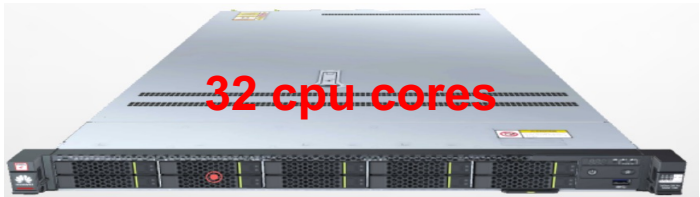  - X86 : Intel 5218@2.3 GHz, 16*2 cpu cores

- **Single core**

| cpu | job | walltime(m) | idle power consumption(w) | Operating power consumption（w) |
|-----|-----|-------------|---------------------------|----------------------------------|
| ARM | 1 | 109 | 300 | 306 |
| X86 | 1 | 90 | 180 | 240 |

- **Full cores**

| cpu | job | walltime | Average walltime of a single job. | Power consumption(w.h) |
|-----|-----|----------|-----------------------------------|-------------------------|
| ARM | 96 | 4h6m | 103.75m | 1355.51 |
| X86 | 32 | 3h20m | 83.09m | 967.36 |

Power consumption variations during the execution of Corsika jobs

# Backup – Arm Evaluation

- HS06 (HEP-SPEC06) test
  - Benchmark of HEP computing
  - ARM : Kunpeng 920@2.6G Hz, 48 cores*2
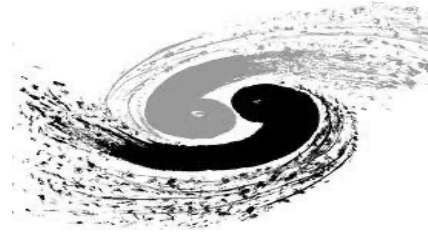  - X86 : AMD EPYC 7773x@2.2G Hz, 64 cores*2
  - Intel Xeon 8352Y@2.20GHz, 32 cores*2
  - Intel Xeon 6258R@2.70GHz, 28 cores*2



**HS06 of Single core**

- Kunpeng 920-6246: 13.4
- AMD 7773x: 24.14
- Intel 8352Y: 20.92
- Intel 6258R: 19.43



**Overall performance and power consumption**

- Kunpeng 920-6246: HS06 1715, consumption 350
- AMD 7773x: HS06 3090, consumption 800
- Intel 8352Y: HS06 1339, consumption 650
- Intel 6258R: HS06 1088, consumption 650

Legend: HS06 ◼ , comsumption ●—