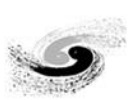




# Applications of Supercomputer Tianhe-II in BESIII

Jingkun Chen<sup>1</sup>, **Biying Hu**<sup>1</sup>, Xiaobin Ji<sup>2</sup>, Qiumei Ma<sup>2</sup>,  
Jian Tang<sup>1</sup>, Ye Yuan<sup>2</sup>, Xiaomei Zhang<sup>2</sup>, Yao Zhang<sup>2</sup>,  
Wenwen Zhao<sup>1</sup>, Wei Zheng<sup>2</sup>

1. Sun Yat-Sen University
2. Institute of High Energy Physics, China



# Outline

- 1 Introduction
- 2 BOSS deployment
- 3 Submission flow
- 4 Large-scale performance test
- 5 Summary

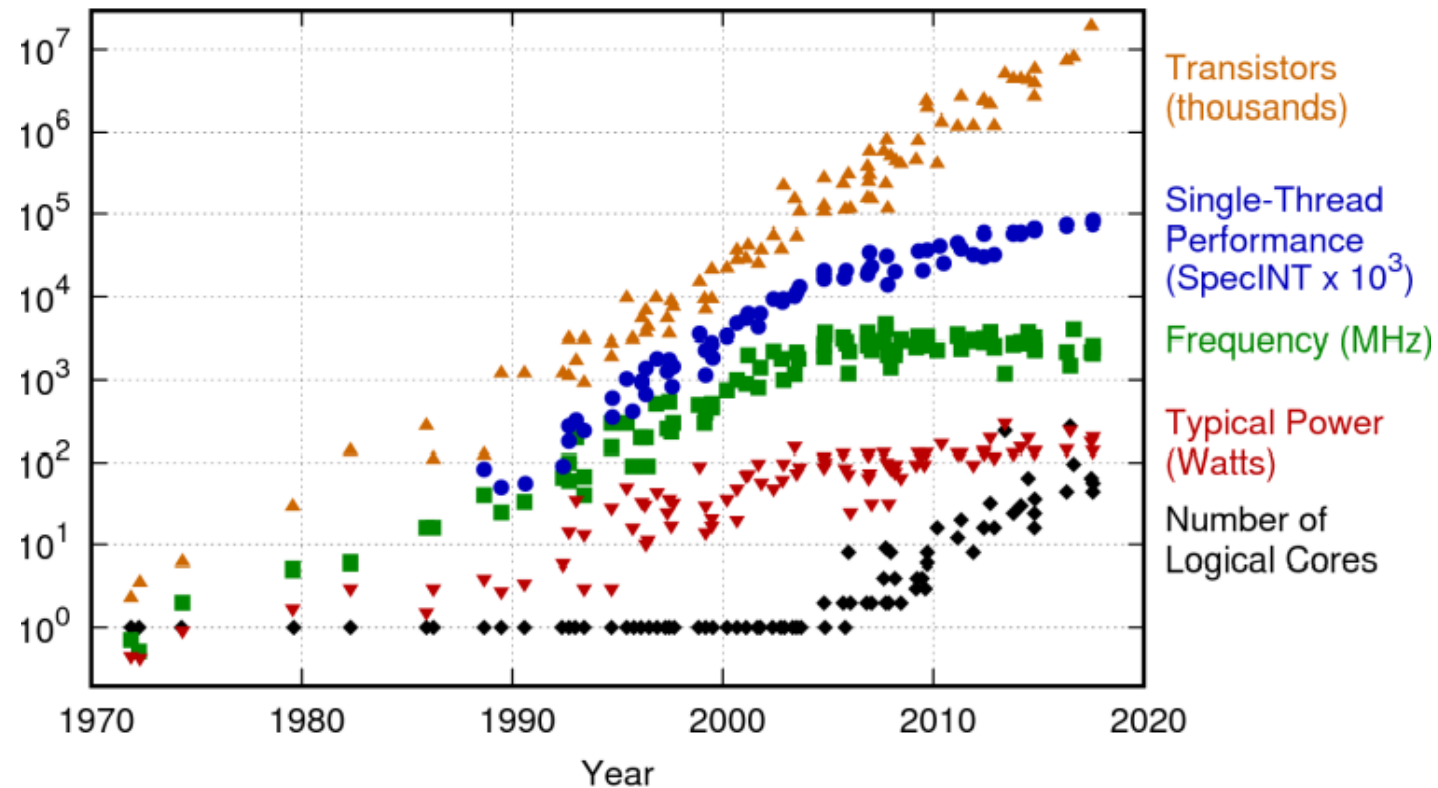


# 1

# Introduction

# Worldwide trend

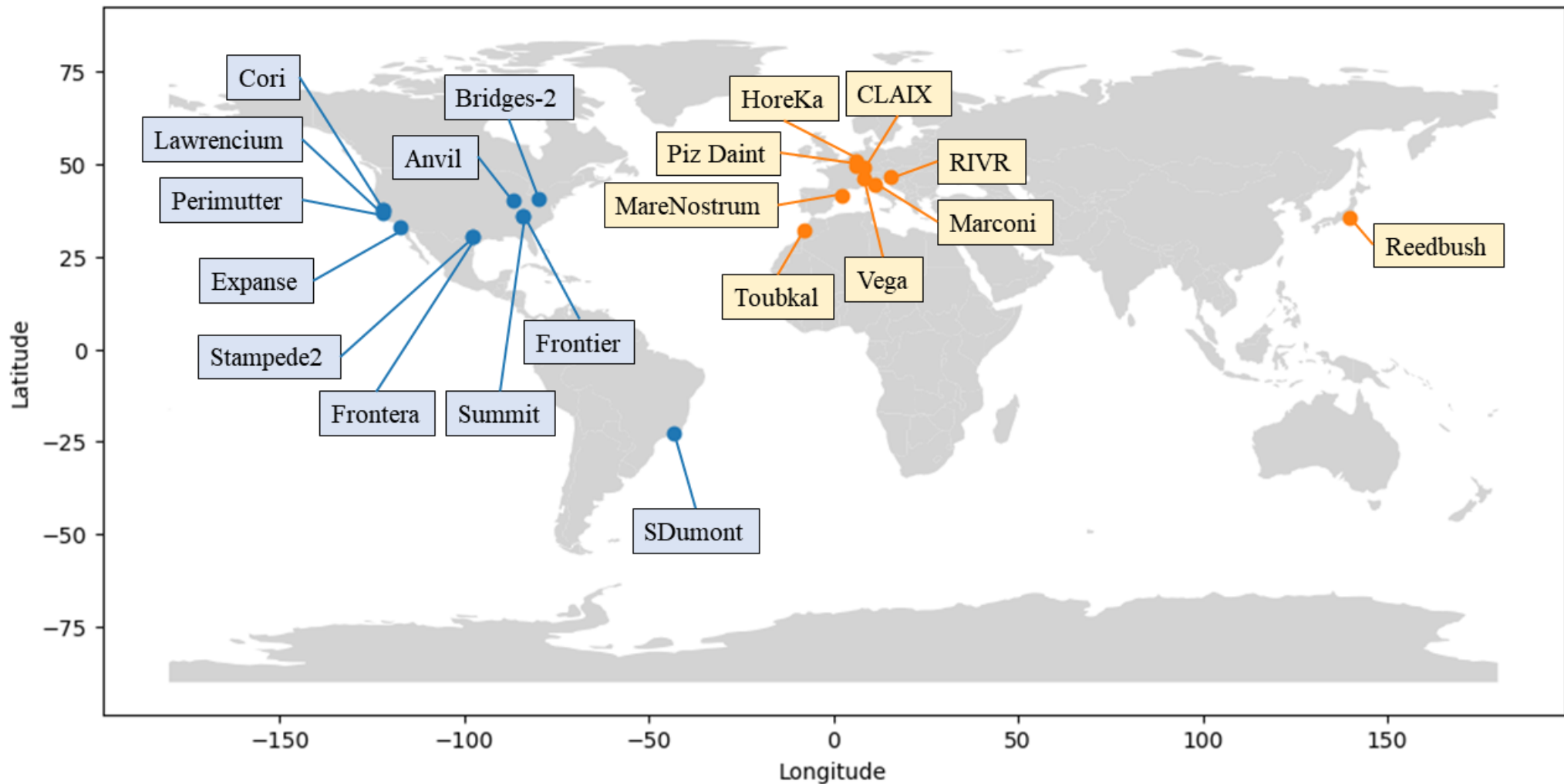
42 Years of Microprocessor Trend Data



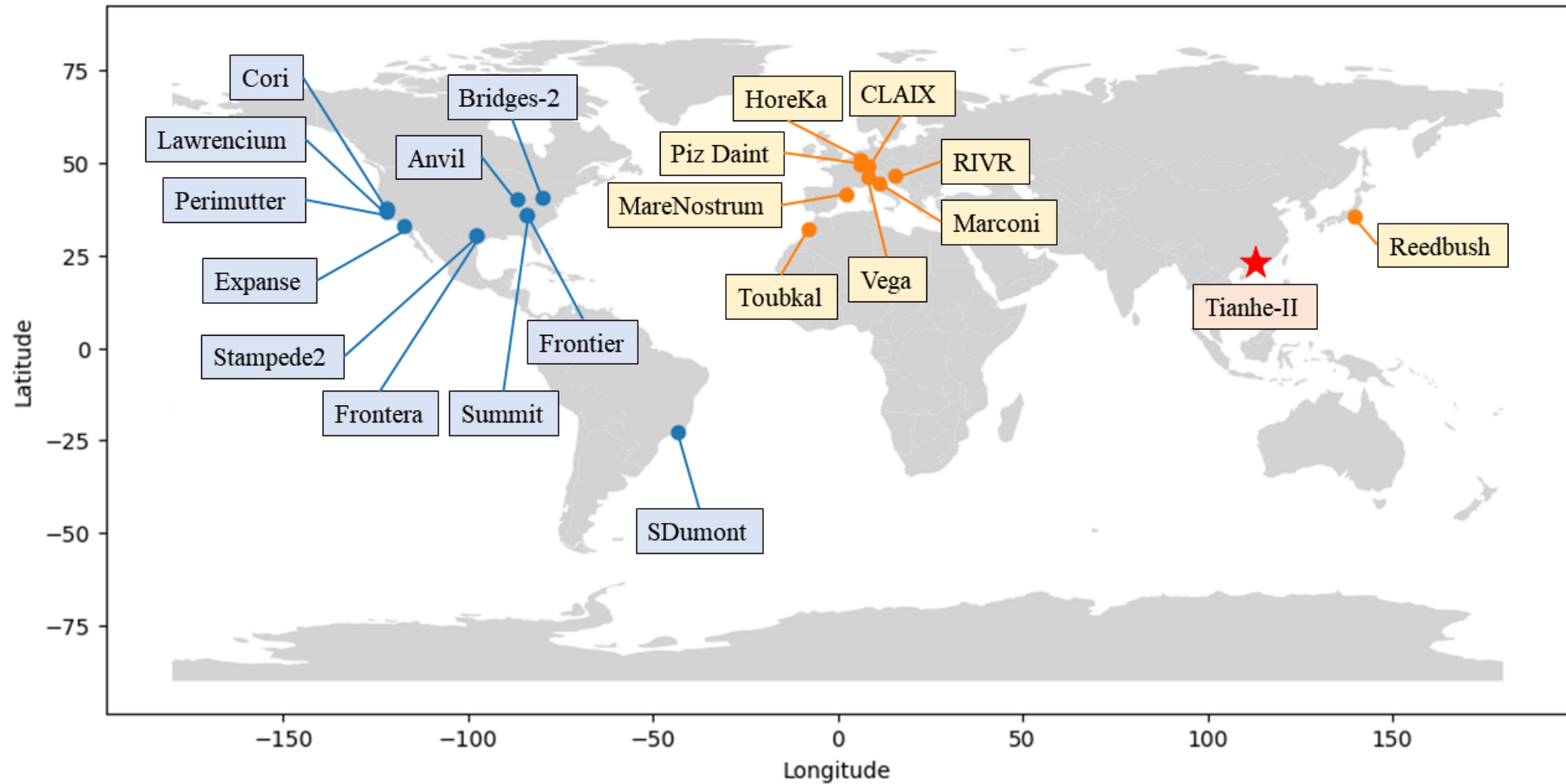
Original data up to the year 2010 collected and plotted by M. Horowitz, F. Labonte, O. Shacham, K. Olukotun, L. Hammond, and C. Batten  
New plot and data collected for 2010-2017 by K. Rupp

Strategic worldwide trend: massive investments into supercomputers.

# Worldwide HPC used for HEP experiment

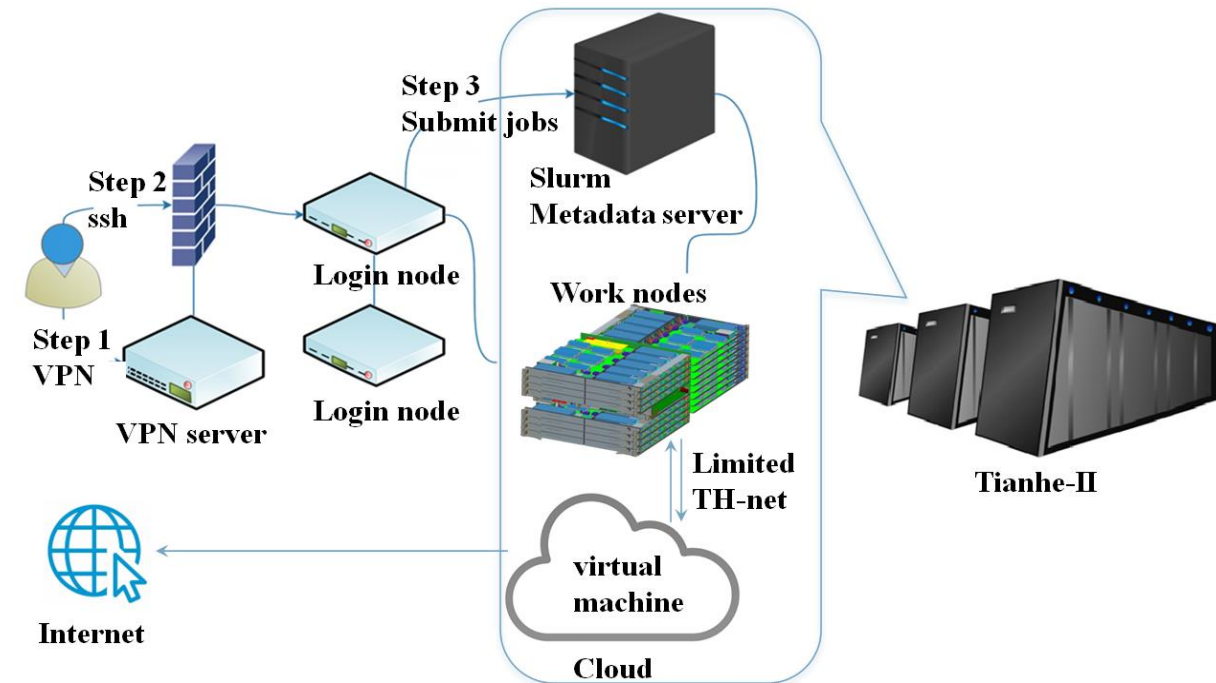


# Worldwide HPC used for HEP experiment



# Tianhe-II

- Located in Sun Yat-sen University, National Supercomputing center in Guangzhou, China
- Total **3,120,000 CPU cores**
- 30.65 Pflops (**world's fastest** at 2013 ~ 2015)
- Adapts serially SLURM system and metadata server
- 24 CPU cores per node and **15 PB** shared file system Lustre
- A new data center to support applications.



# HTC & HPC

## HTC (High-throughput Computing)

- Long timescale
- Data intensive
- Designed for HEP need

## HPC (High-performance Computing)

- Short timescale
- Computing intensive
- Not designed for HEP

### Challenges:

- Strict security policy.
- Limit of network.
- I/O crash caused by large input and output data.

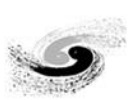
**All supercomputers will encounter these challenges,  
but no general solution for HPC to run HTC jobs.**





# Need to achieve

- BESIII Offline Software System (BOSS) deployment
- Remote submission
- Large-scale performance test



# 2 BOSS deployment



# Virtualization container

- **No root authority** on Tianhe-II
- Singularity is designed to run the container without a root authority.
- Two types of container:
  - Fat container with all the required software and relevant database.
  - Light container only includes the BOSS environment that applies to all BOSS versions.

# Deployment Solution 1 (done)

- Solution 1
  - Use container virtualization technology and a **fat container**
- Fat container disadvantages
  - Follow-up update of BOSS is very troublesome.
  - Requires **large storage space** at Tianhe-II
  - Requires **great human effort** to maintain many containers





# Deployment Solution 2

- Solution 2
  - Use Cern-VM File System (**CVMFS**) and a **light container**
- Two problems:
  - 1. Install the **cvmfs client** on Tianhe-2;
  - 2. Realize real-time access to the **database**



# Install CVMFS

Without root authority, we have failed in many approaches.

- General installation ❌
- Using a pre-installed executable CVMFS Cvmfsexec ❌
- Parrot-mount CVMFS ❌
- Our approach: **compile CVMFS from source code** ✓

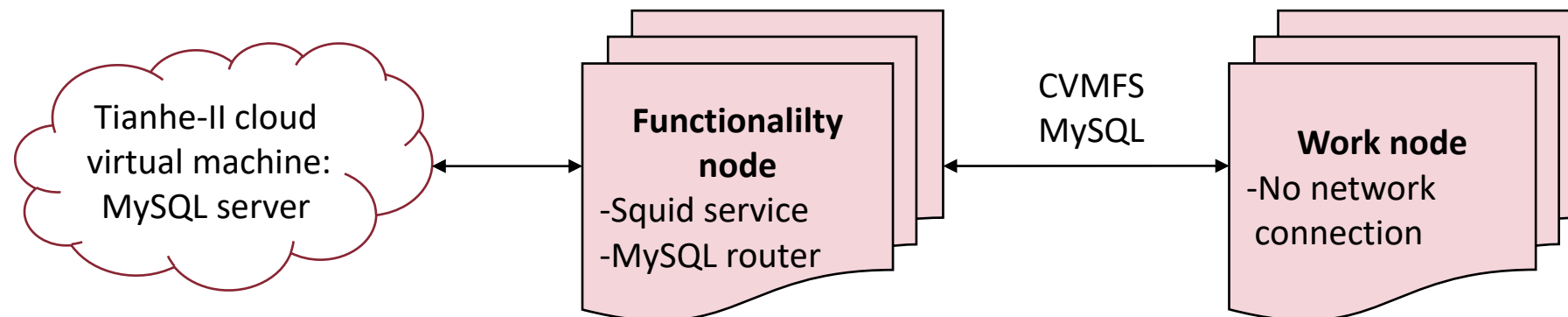
# Database connection

## Squid

- Squid is a caching and forwarding HTTP web proxy.
- Install Squid on functionality nodes.
- Allow CVMFS connect to the IHEP code database.

## MySQL

- Deploy a MySQL server in cloud Virtual machine.
- Deploy a MySQL router as a bridge in functionality node

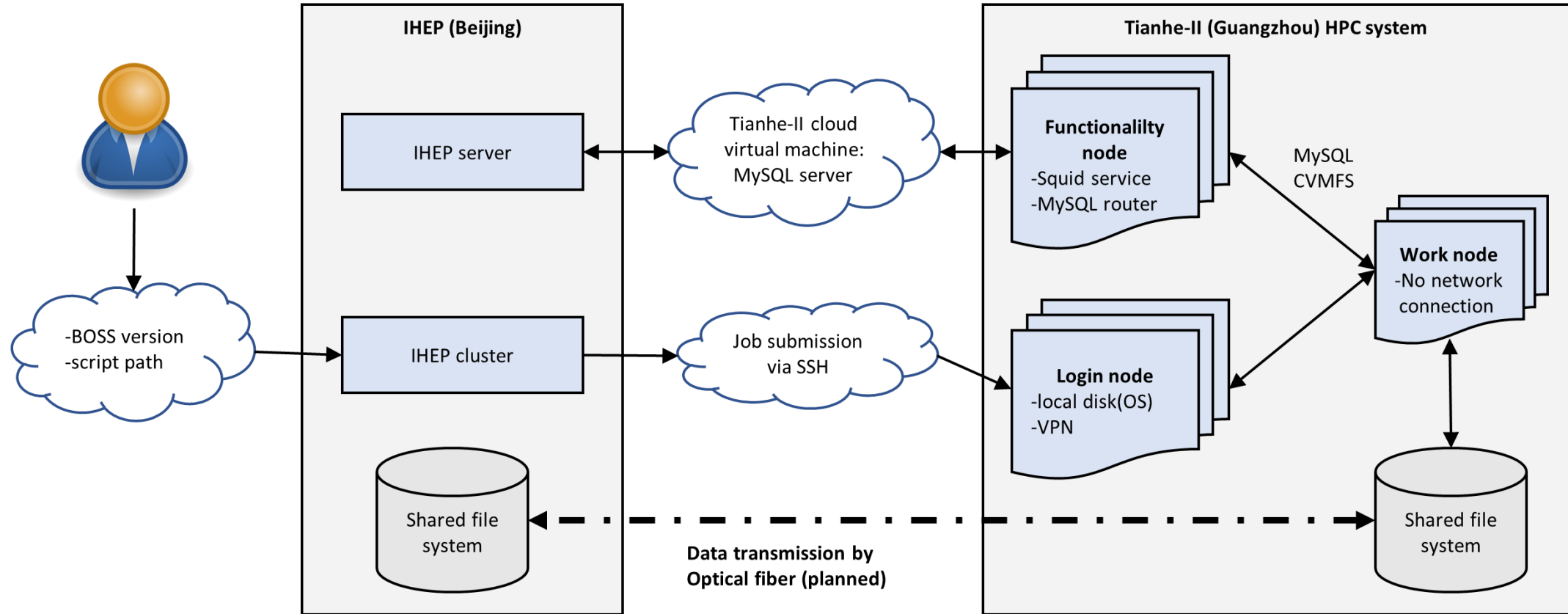




# 3 Submission flow

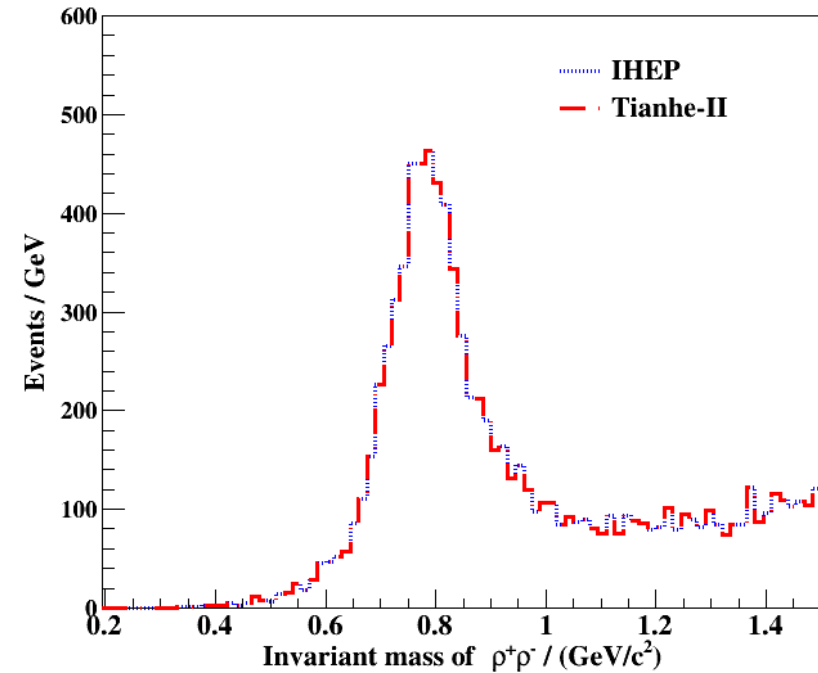
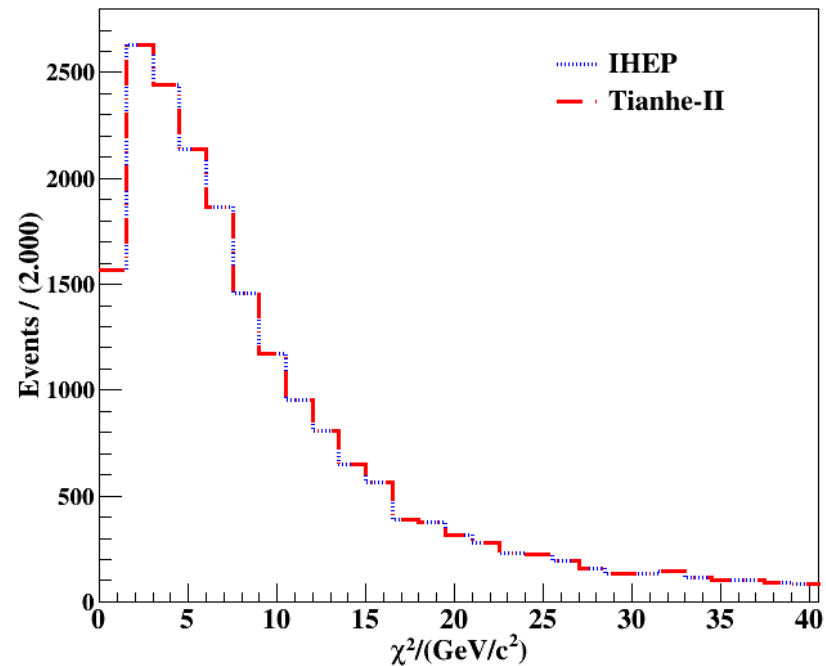


# Workflow



# Validation

$$e^+ e^- \rightarrow J/\psi \rightarrow \rho\pi$$



Invariant mass spectrums of  $\rho$  meson and chi-square values of the four momentum from IHEP and Tianhe-II simulation results are completely the same.



# What can we do on Tianhe-II?

## Network status:

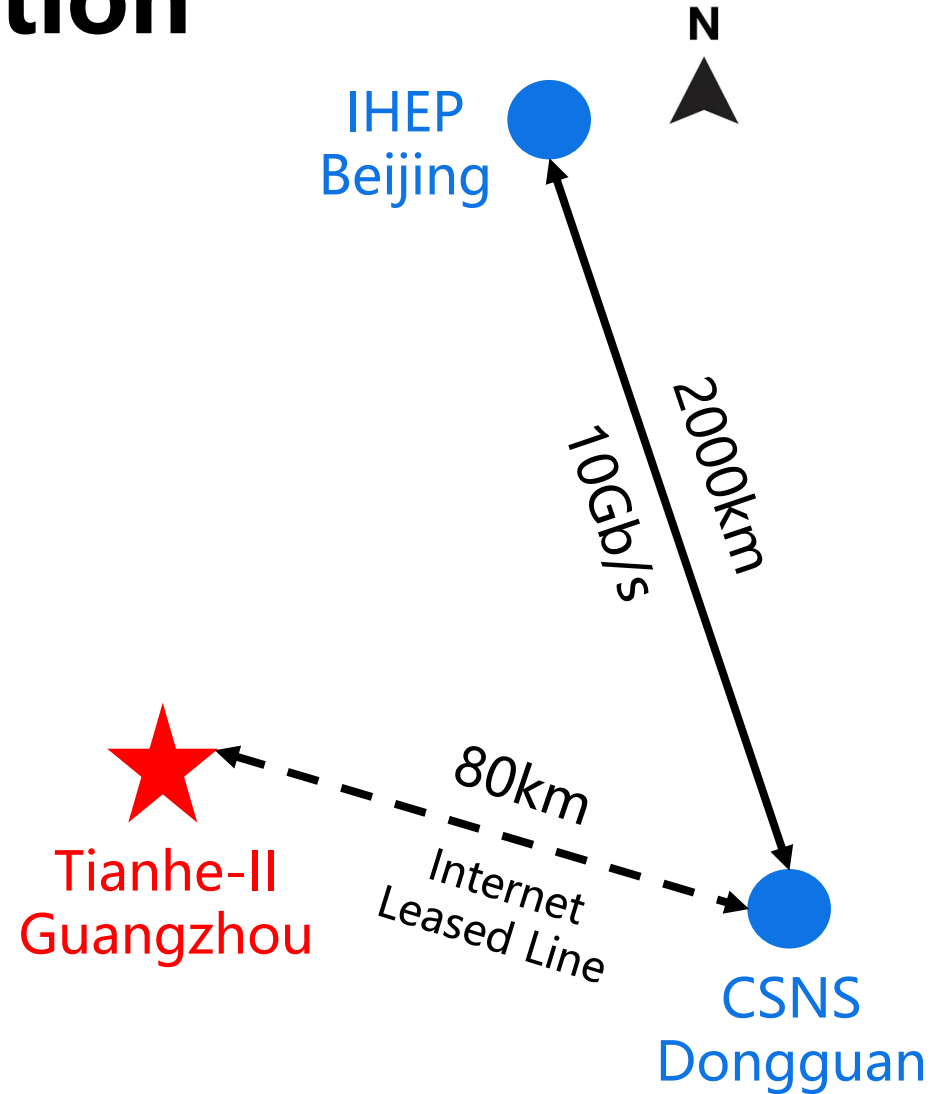
- The speed test for IPV4 only reach to **20MB/s**.
- Adequate for software deployment and database update, but **insufficient** for mass data transmission.
- The **operators** of IHEP and Tianhe-II are **different**, and IPV6 has only 50Mb bandwidth.

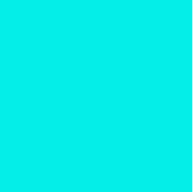



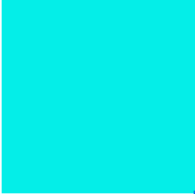
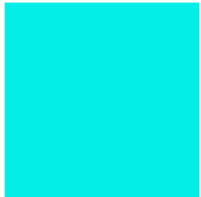
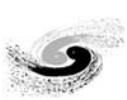
## How to use:

- Base on **simulation**, which requires very few input files.
- **Reconstruct** and **analyze** the mass simulation data on Tianhe-II instead of sending back IHEP.

# Network Solution

The **China Spallation Neutron Source (CSNS)** is the pulsed neutron source facility located at **Dongguan City**, which is next to **Guangzhou** City. To meet the need for data services, the National HEP Science Data Center (NHEPSDC) in **Beijing** set up its Branch Center at the CSNS. The Branch Center has more than **10 PB** of storage space, an international network link of **10,000 megabits per second**, as well as a complete information support system.

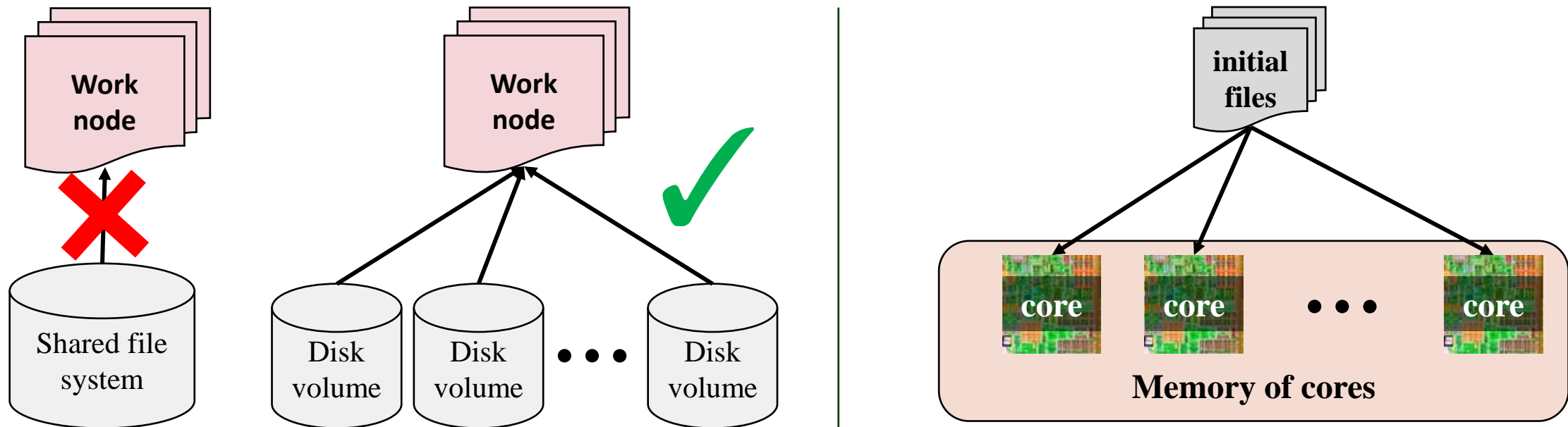




# 4 Large-scale performance

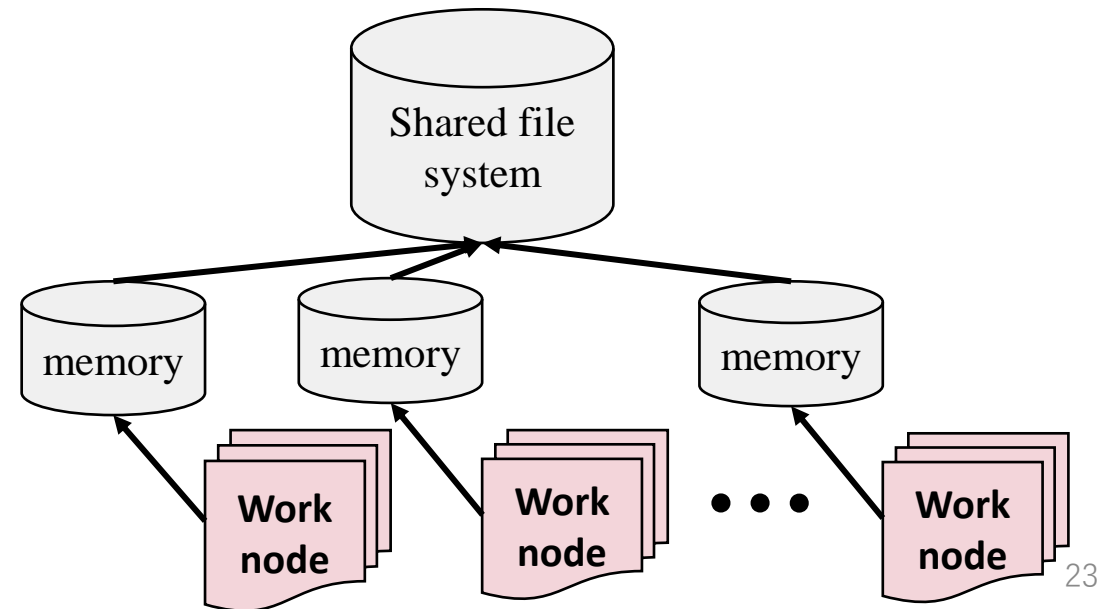
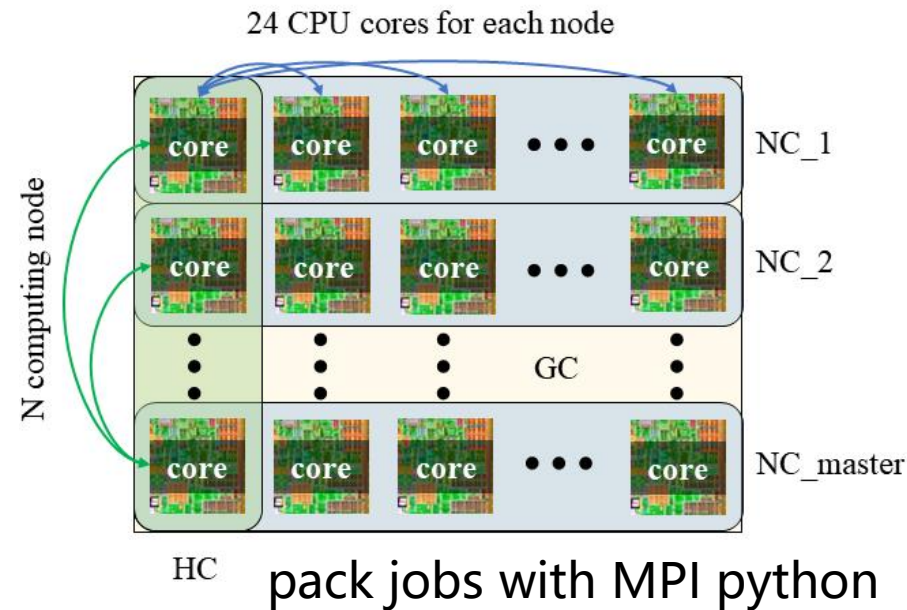
# Input issue

- Input too many files from the share file system may cause IO congestion.
- Storage mirror files in **each disk volume**
- Create **initial files in memory** instead of reading in file system.



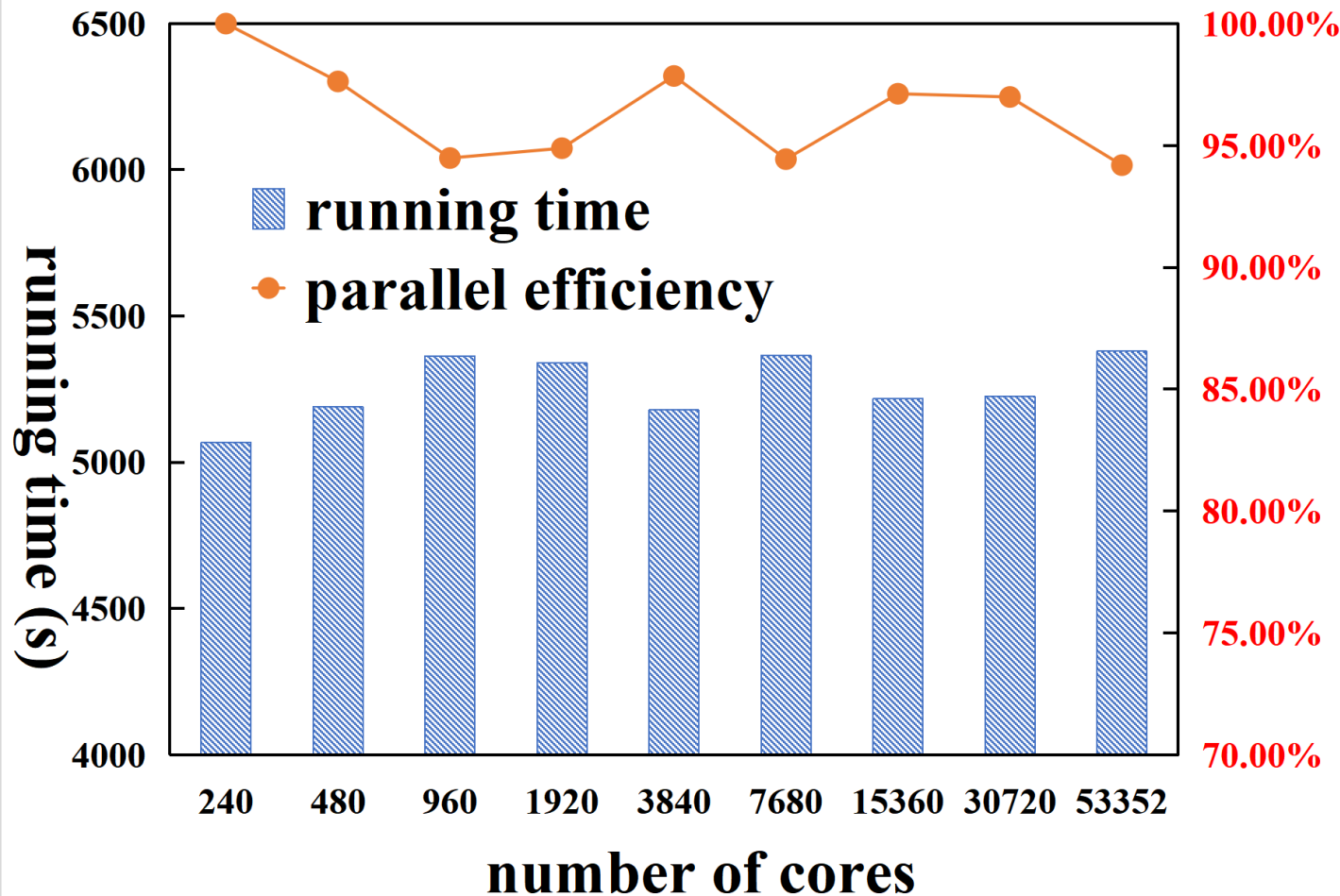
# Output issue

- The Tianhe-II SLURM system could become **crash** if numerous HTC jobs are submitted in a short time.
- Develop a **MPI python submitting script** to pack the HTC jobs and control the I/O.
- Save the **output in memory** first and move to file system **in order**.



# Large-scale performance result

## Tianhe-II test results



Parallel efficiency is used to measure the **performance of parallel systems**.

Parallel efficiency define as:

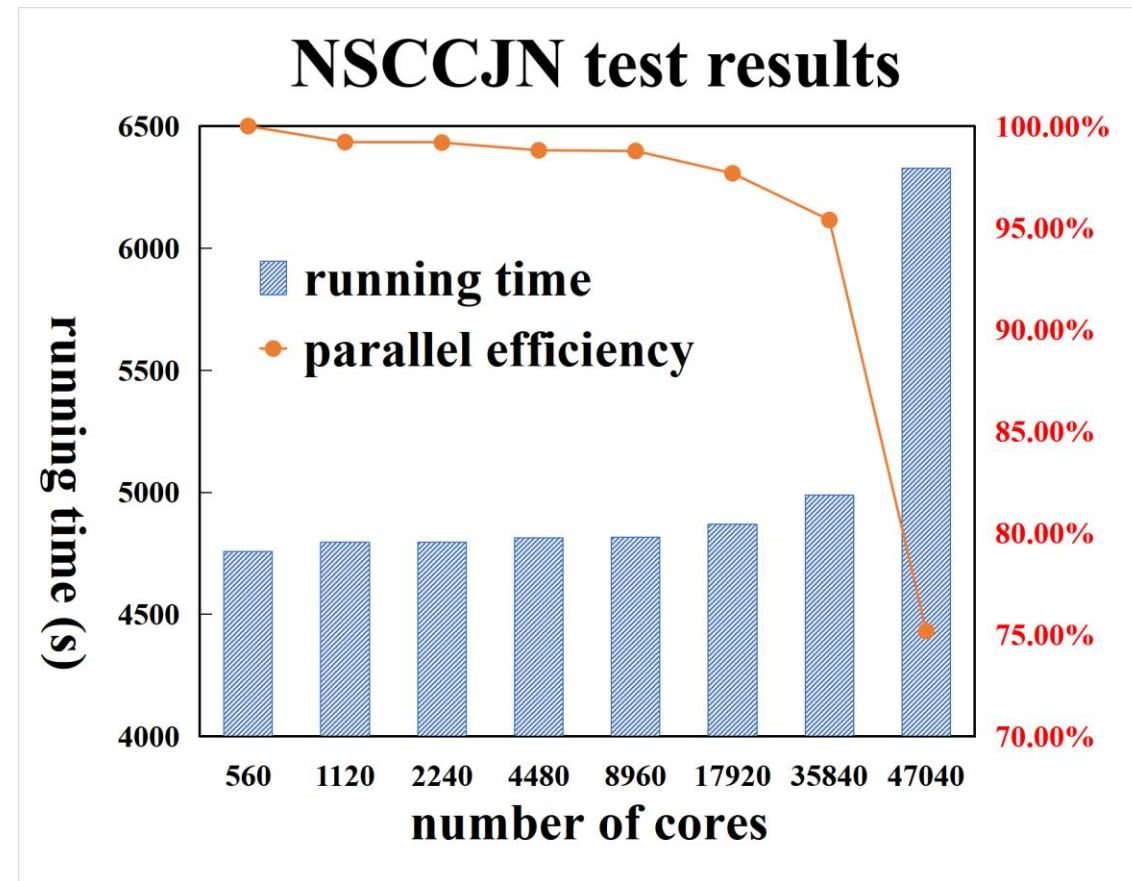
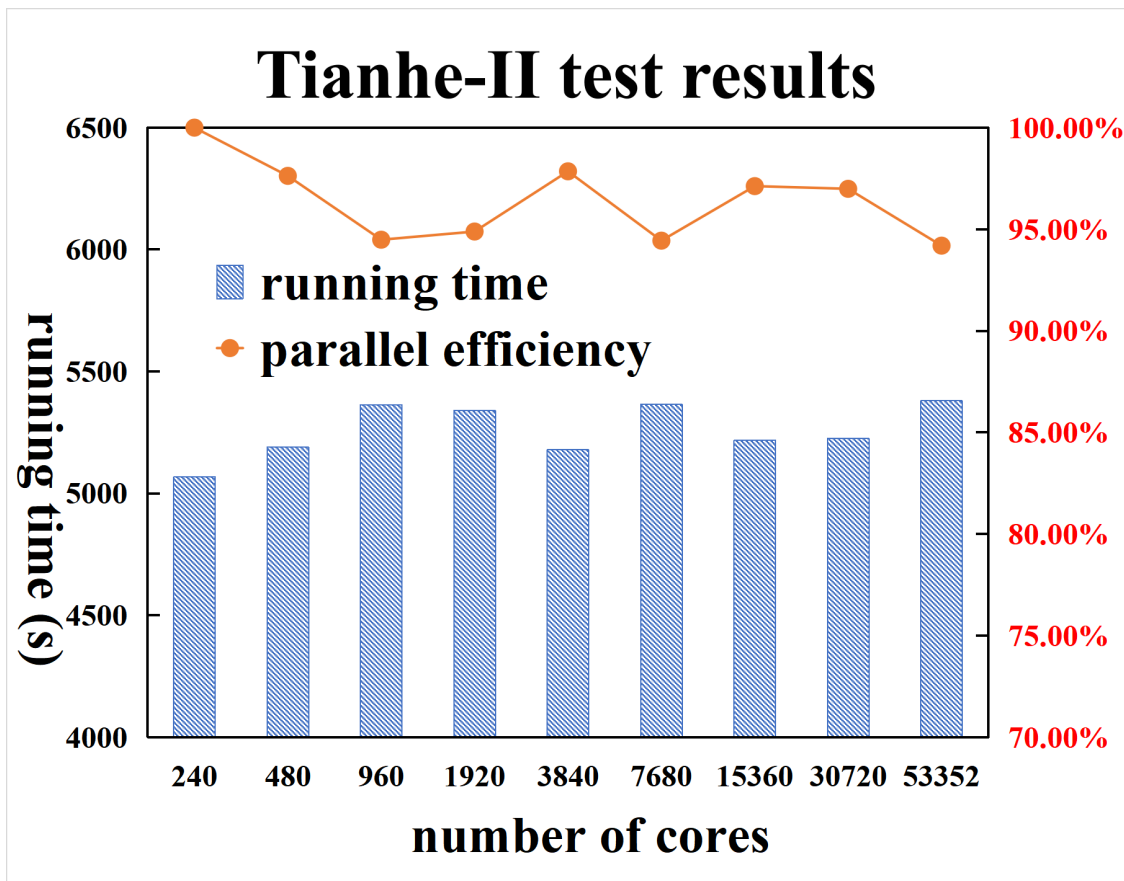
$$E(P, N) = \frac{t_1}{t_p}$$

$t_1$  is the running time on single nodes,  
 $t_p$  is the running time on P nodes.

The ideal parallel efficiency is 100%, and **80%** is considered as good performance.



# Large-scale performance result



We ran BOSS with a light container on Tianhe-II and a fat container on NSCCJN.

The parallel efficiency of NSCCJN drops to 75% due to the **IO crash**.



# Summary & Prospects

- Summary

- Real-time updates of BOSS on Tianhe-II is realized.
- Remote job submission to Tianhe-II from IHEP farm is scheduled.
- Large-scale performance running is done.

- Prospects

- Look forward to cooperating with the physics groups.
- Improve the data transfer speed, for example, set up a bare optical fiber between Tianhe-II and CSNS.
- More sites and resources to join in the future.



中山大學  
SUN YAT-SEN UNIVERSITY



中国科学院高能物理研究所  
Institute of High Energy Physics Chinese Academy of Sciences



国家超级计算广州中心  
NATIONAL SUPERCOMPUTER CENTER IN GUANGZHOU

# Thank you!