# Overview of Tianhe2 System and Application

Yutong Lu

Professor

School of computer science, NUDT &

State Key Laboratory of High Performance Computing, China

国防科学技术大学
*National University of Defense Technology*

HPCL

# Outline

天河

☐ **NUDT HPC Background**

☐ **Design and APP of Tianhe2 System**

☐ **Prospect of Tianhe-2A**

☐ **Summary**

国防科学技术大学
*National University of Defense Technology*

HPCL

## National university of defense technology



~2,000  Teachers
~15,000 Students
… Others

国防科学技术大学
*National University of Defense Technology*

# Overview

天河

## Supercomputers in NUDT, Changsha, China



100P System

2014 TH-2 33.86PF Top1

2010 TH-1A 2.67PF Top1

2009 TH-1 1st Chinese PFlops System (Heterogeneous)

2000 YH-3 1st Chinese TFlops System (MPP)

1992 YH-2 1st Chinese GFlops System (SMP)

1983 YH-1 1st Chinese Supercomputer (PVP)

国防科学技术大学
*National University of Defense Technology*

HPCL

# Overview

天河

## Supercomputing Centers in China



NSCC-Guangzhou,2013
Tianhe-2

NSCC-Wuxi
Shenwei-NG

NSCC-Changsha,2012
Tianhe-1A

NSCC-Jinan,2012
Shenwei-Bluelight

NSCC-Tianjin,2010
Tianhe-1A

NSCC-Shenzhen,2011
Dawning-6000

国防科学技术大学
National University of Defense Technology

HPCL

# Overview

Challenges(PSPRD)

- Performance
- Scalability
- Power consumption
- Reliability
- big Data

☐ Heterogeneous Trend

- Some of top-level supercomputers
- Tianhe-1/2, Titan …

☐ Compute Efficiency

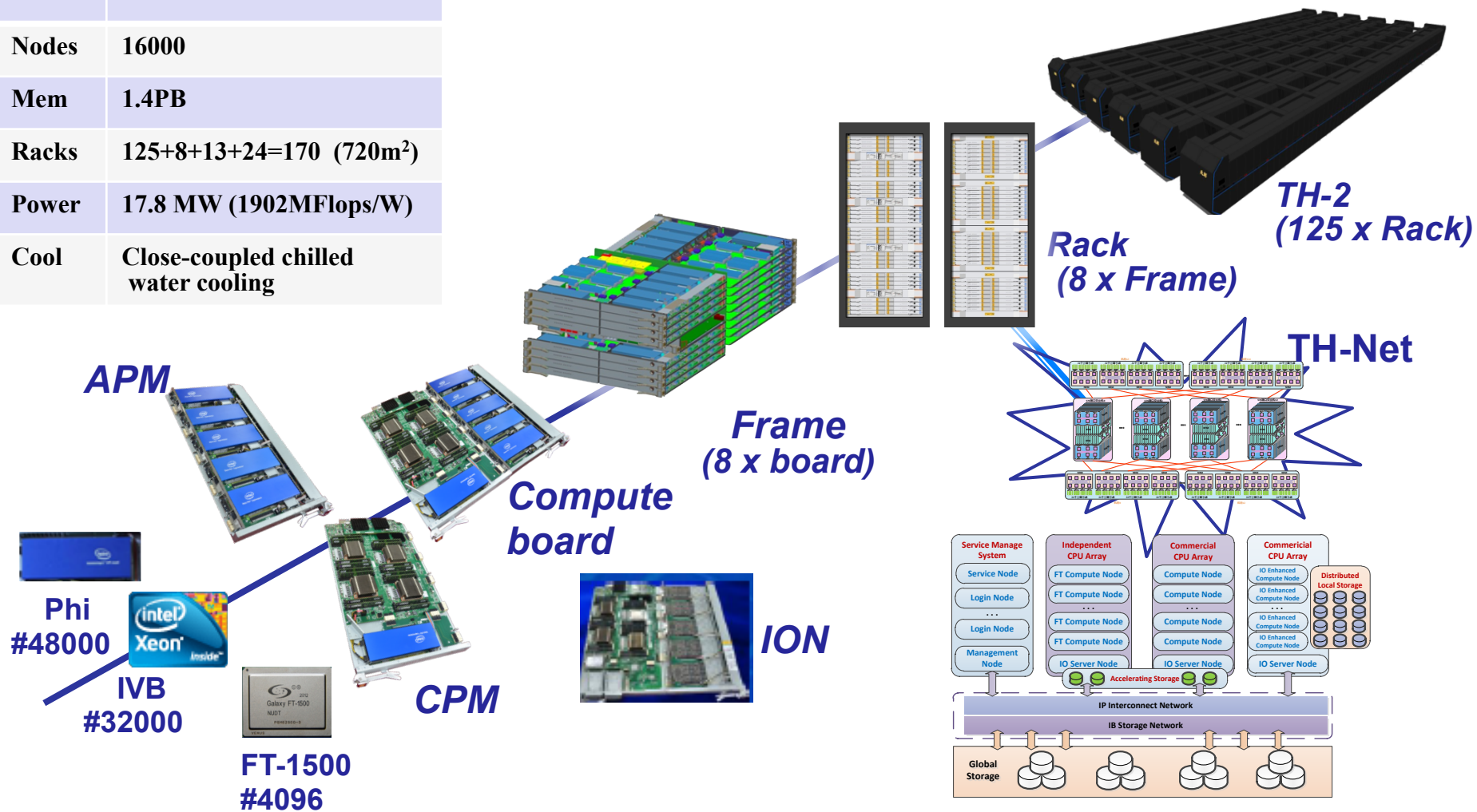- More computations per transistor
- More computations per joule

国防科学技术大学
*National University of Defense Technology*

HPCL

# Highlights of Tianhe-2

天河

| Perf | 54.9PFlops / 33.86PFlops |
|------|--------------------------|
| Nodes | 16000 |
| Mem | 1.4PB |
| Racks | 125+8+13+24=170 (720m$^2$) |
| Power | 17.8 MW (1902MFlops/W) |
| Cool | Close-coupled chilled water cooling |

**TH-2 (125 x Rack)**

**Rack (8 x Frame)**

**TH-Net**

**Frame (8 x board)**

**APM**

**Compute board**

**Phi #48000**

**IVB #32000**

**FT-1500 #4096**

**CPM**

**ION**

| Service Manage System | Independent CPU Array | Commercial CPU Array | Commericial CPU Array | |
|------|------|------|------|------|
| Service Node | FT Compute Node | Compute Node | IO Enhanced Compute Node | Distributed Local Storage |
| Login Node | FT Compute Node | Compute Node | IO Enhanced Compute Node | |
| ... | ... | ... | IO Enhanced Compute Node | |
| Login Node | FT Compute Node | Compute Node | IO Enhanced Compute Node | |
| Management Node | FT Compute Node | Compute Node | IO Enhanced Compute Node | |
| | IO Server Node | IO Server Node | IO Server Node | |

Accelerating Storage

IP Interconnect Network

IB Storage Network

Global Storage

**Hybrid Hierarchy shared storage System 12.4PB**

国防科学技术大学
*National University of Defense Technology*

HPCL

# Highlights of Tianhe-2

天河

□ Software Stack

| Intelligent Monitor & System Management | HPC Application Service Platform | Scientific Data Visualization System | Cloud Computing Platform | PAE |
|---|---|---|---|---|
| | Multi-Domain Framework | | | |
| | MPI / C/C++/Fortran | GA / OpenMP | OpenMC / Intel Offload | Tools & Library / PDE |
| | Hybrid Runtime System | | | |
| | Resource Management System | | | PSE |
| | H²FS Parallel File System | | | |
| | Kylin Operating System | | | |

国防科学技术大学
*National University of Defense Technology*

HPCL

# Adaptive Interconnect

天河

## TH-Express 2+

☐ Network Interface Chip (NIC)
- ➤ Bandwidth
  - ◆ 14Gbps x 8lane
  - ◆ 12GB/s
- ➤ Function extension
  - ◆ Low latency and high rate operations
  - ◆ Collective comm offload
  - ◆ Parallel processing using Multiple DMA engines
  - ◆ Enhanced reliability

Extending to >100PFlops system

☐ Network Router Chip (NRC)
- ➤ 24 ports
  - ◆ Up to 5.376 Tb/s of switching capacity
- ➤ Table-based routing
  - ◆ Multi-path adaptive routing
  - ◆ Oblivious routing

# Scalable MPI

天河

## Message passing system based on TH-Express

☐ Multiple Communication Protocols

➢ Performance oriented

➢ Scalability oriented

➢ Combine application model

➢ Zero-copy data transfer

☐ Collective Operation Offload

➢ Construct topology-aware algorithm tree dynamically

➢ Message pass automatically based on the trigger of NIC

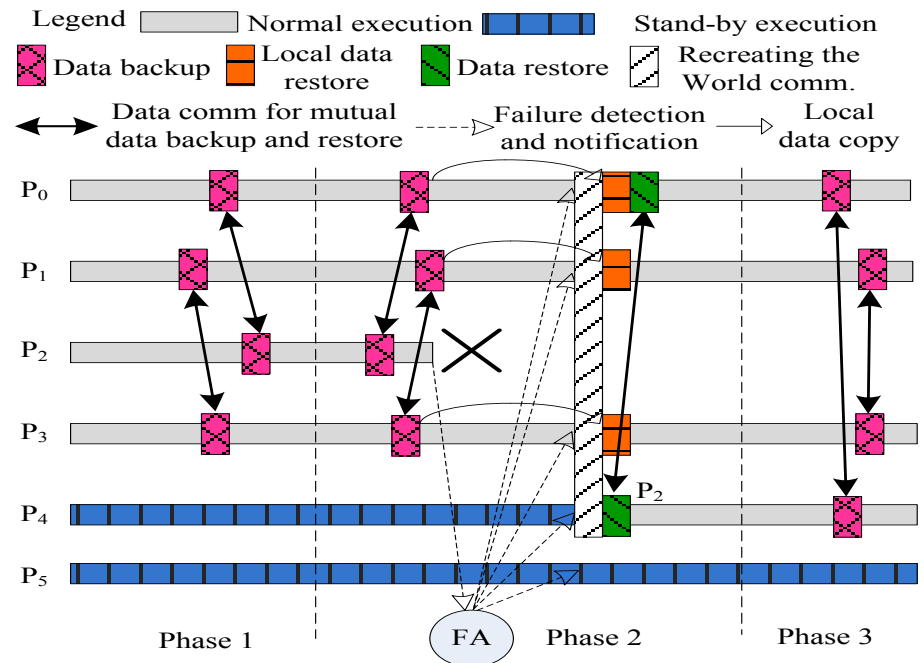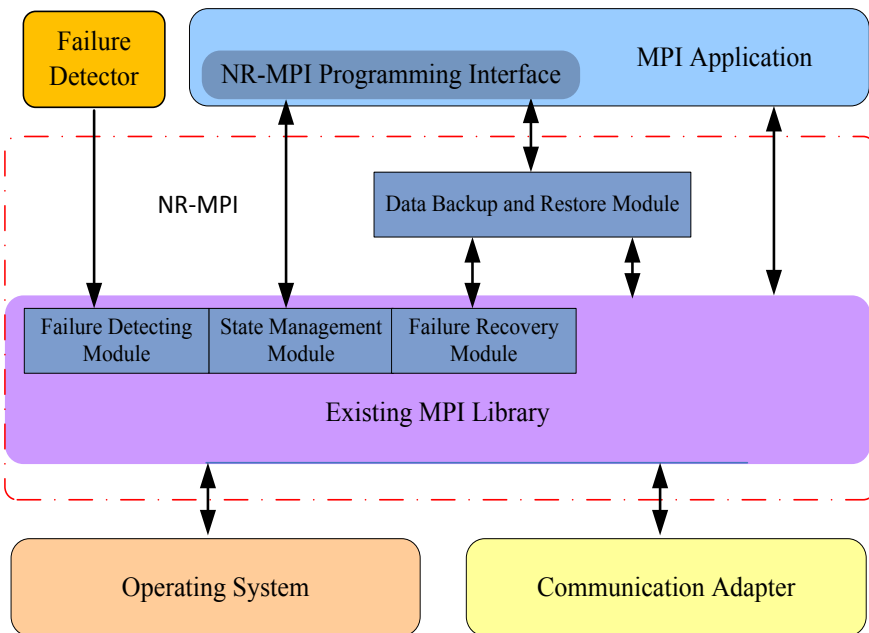➢ Bypass effect of OS noise, Reduce Latency of large scale





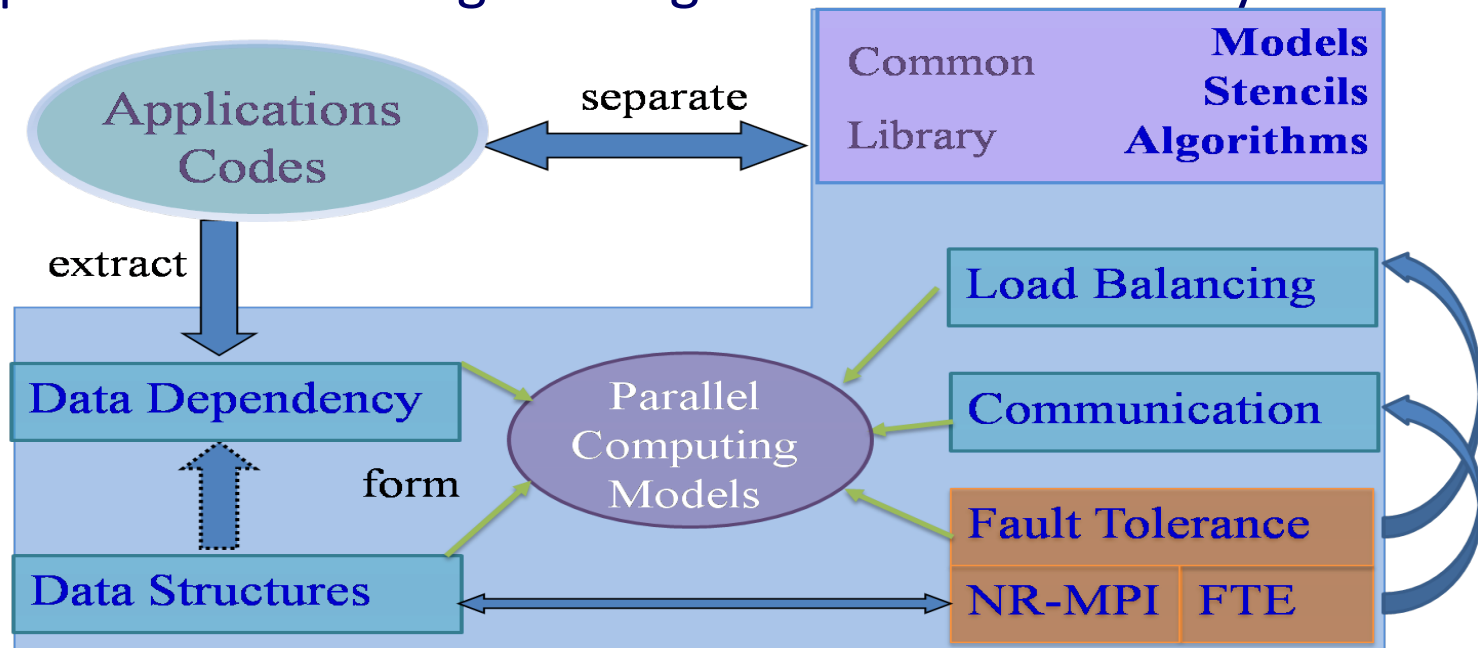

国防科学技术大学
National University of Defense Technology

# Scalable MPI

☐ Non-stop and fault Resilient MPI (NR-MPI)

  ➤ Application continue execution without being relaunched
  ➤ Failure detection and MPI state recovery done by runtime
  ➤ Data-backup by application-level diskless C/R
  ➤ Reconstruct of MPI communicator and channel

# Domain Specific Framework

天河

- ☐ Hide parallel programming complexity using millions of cores and the hierarchy of parallel computers
- ☐ Integrates the efficient implementations of parallel fast algorithms
- ☐ Provides efficient data structures and solver libraries
- ☐ Supports software engineering for code extensibility



国防科学技术大学
*National University of Defense Technology*

HPCL

# Scalable IO Structure

## □ IO Architecture on Tianhe-2

- ➤ Multiple Layers & Hybrid Storages
  - ◆ Local Disk
  - ◆ PCI-E SSD
  - ◆ Disk Array
- ➤ 6400 local Disks
  - ◆ Bus attached
- ➤ 256 IO nodes
  - ◆ Burst: above 1TB/s
  - ◆ TH-Express and IB QDR port
- ➤ 64 Storage Servers
  - ◆ Sustained：about 100GB/s



国防科学技术大学
National University of Defense Technology
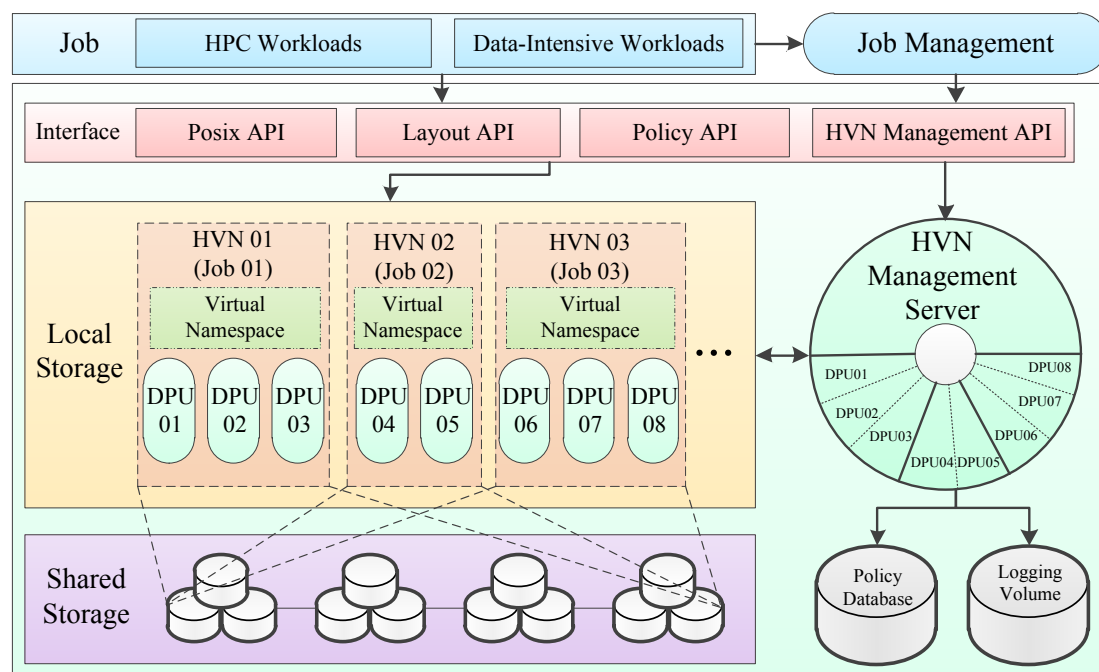
# Scalable IO Structure

天河

- ❑ H$^2$FS: Hybrid Hierarchy File System
  - ➢ DPU, A fundamental unit for data processing, tightly couples a compute node with its local storage
  - ➢ HVN, Hybrid, Unified and Isolated dynamic namespace maintained by centralized servers
  - ➢ Layered and enriched metadata, I/O hints as high level metadata
- ❑ I/O API
  - ➢ POSIX
  - ➢ MPI-IO
  - ➢ Extended API, layout and policy guide
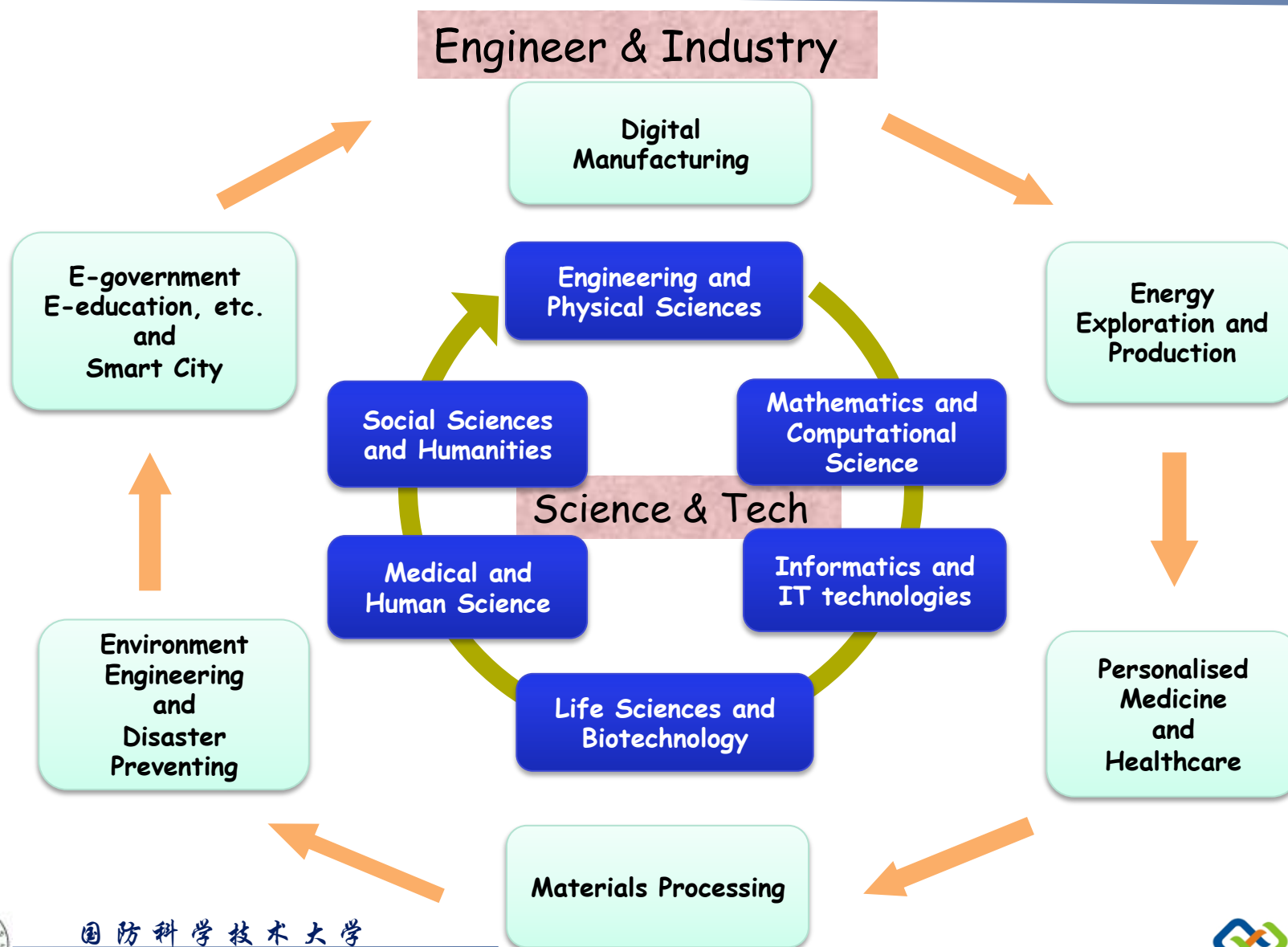  - ➢ HDF5 over POSIX and extended API

| Job | HPC Workloads | Data-Intensive Workloads | Job Management |

| Interface | Posix API | Layout API | Policy API | HVN Management API |

Local Storage

HVN 01 (Job 01) — Virtual Namespace — DPU 01, DPU 02, DPU 03

HVN 02 (Job 02) — Virtual Namespace — DPU 04, DPU 05

HVN 03 (Job 03) — Virtual Namespace — DPU 06, DPU 07, DPU 08

...

HVN Management Server — DPU01, DPU02, DPU03, DPU04, DPU05, DPU06, DPU07, DPU08

Shared Storage

Policy Database

Logging Volume

国防科学技术大学
National University of Defense Technology

# Scalable Application

天河

**Engineer & Industry**

Digital Manufacturing

E-government E-education, etc. and Smart City

Engineering and Physical Sciences

Energy Exploration and Production

Social Sciences and Humanities

Mathematics and Computational Science

**Science & Tech**

Medical and Human Science

Informatics and IT technologies

Environment Engineering and Disaster Preventing

Life Sciences and Biotechnology

Personalised Medicine and Healthcare

Materials Processing

国防科学技术大学
*National University of Defense Technology*

HPCL

# Application Case -- CFD

☐ Hybrid RANS/LES simulation of scramjet combustion

- ➤ CPU+MIC version is 8.63X to 13.86X faster than the original CPU version

- ➤ The largest scramjet combustion simulation up to now: totally 998,400 cores on 26,880 million cells, parallel efficiency 79%.
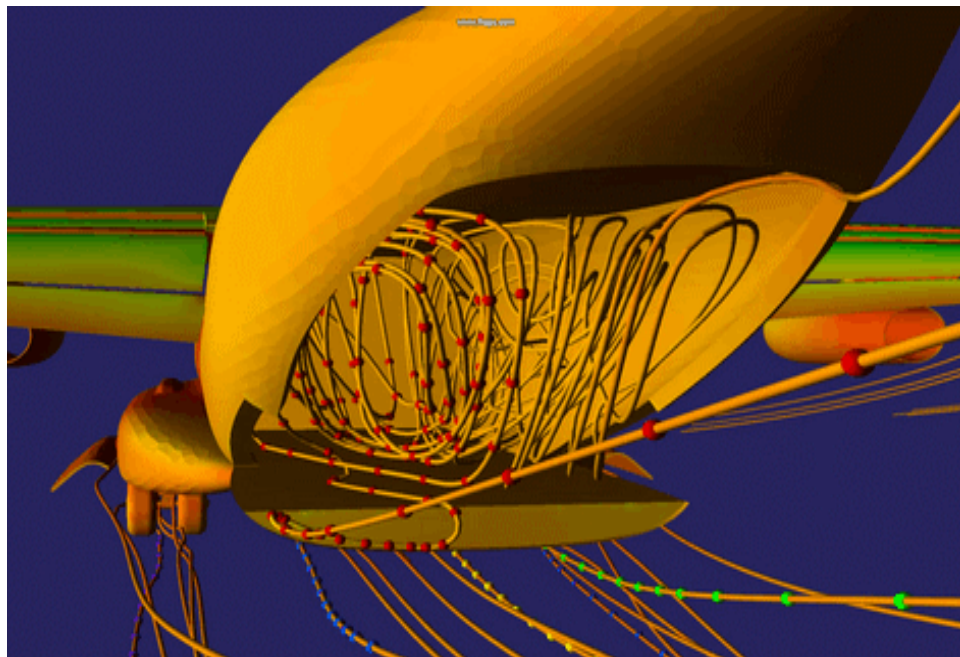




$t_1 = 1\Delta t$

packet 2    packet 1

# Application Case -- CFD

天河

□ Large transport aircraft
□ Large passenger aircraft

# Application Case – Bio-Medicine

- High Throughput Virtual Screening
- Applications: Computer Aided Drug Design, Molecular Docking, and Virtual screening
- mD$^3$DOCKxb,  Lamarckian Genetic Algorithm
- Data Scale, 40 millions molecules, 800TB
- Bottleneck, IO BW, Comm BW

Find out the best 100 molecules;

Do experiment;

Clinical

Aim: Finish docking of all the drug molecules on earth within one day.

HPCL
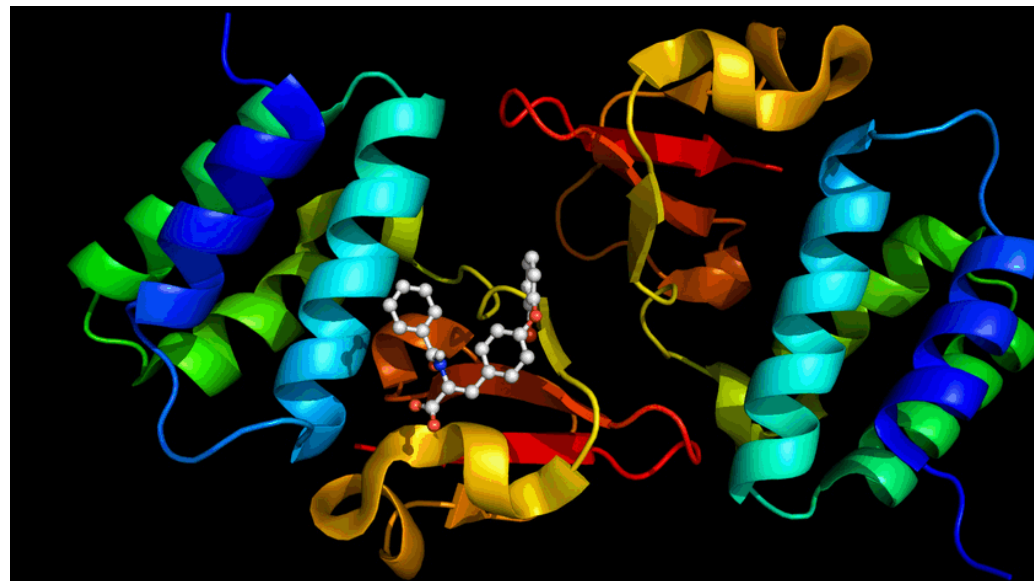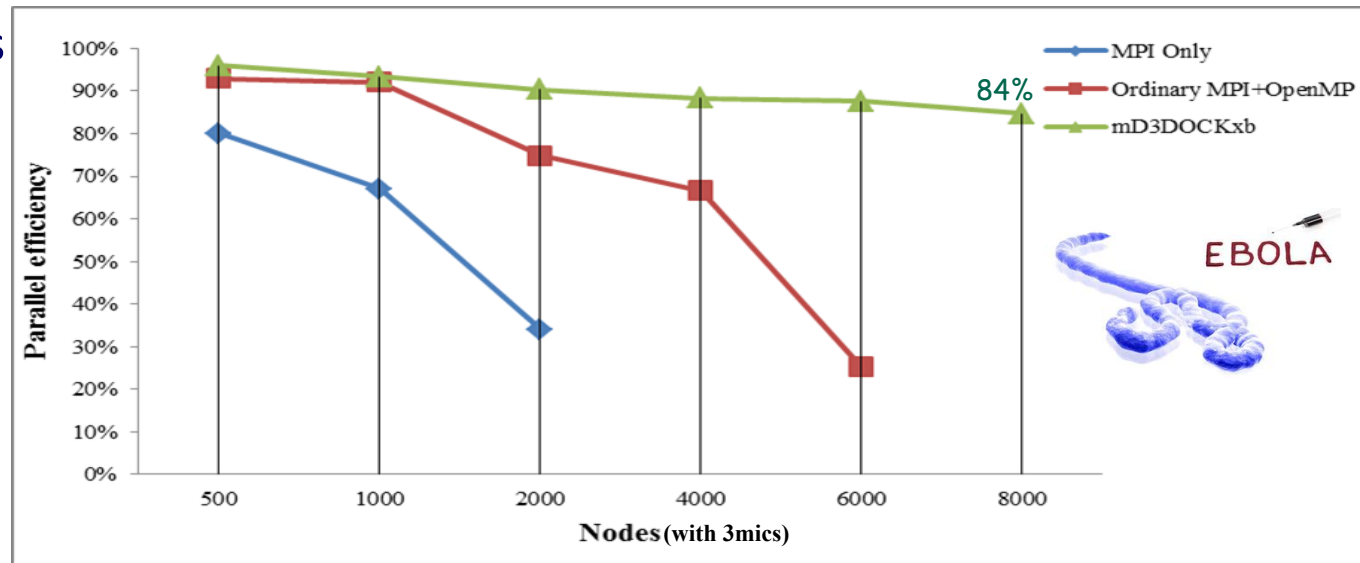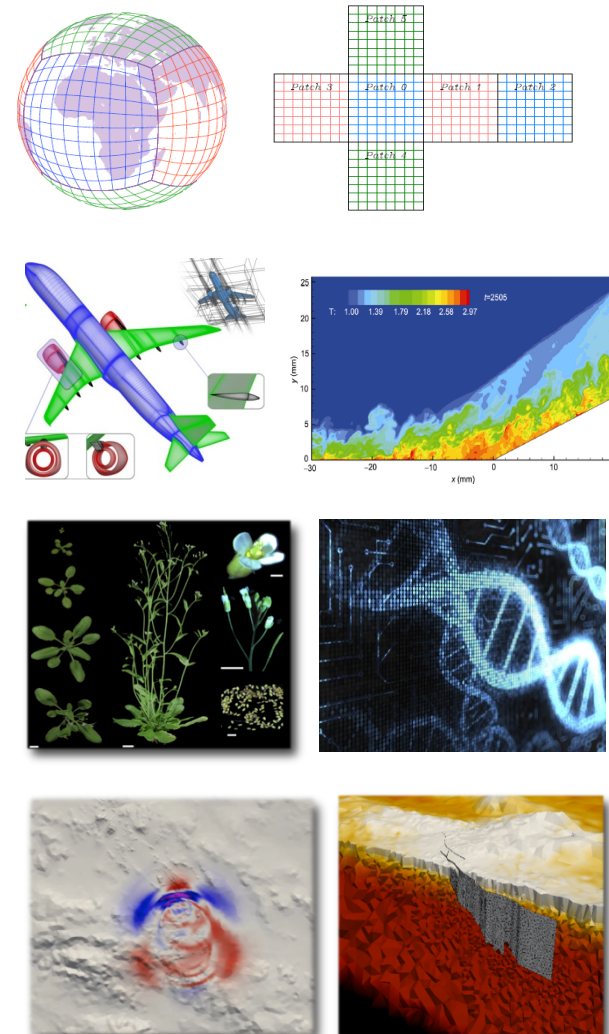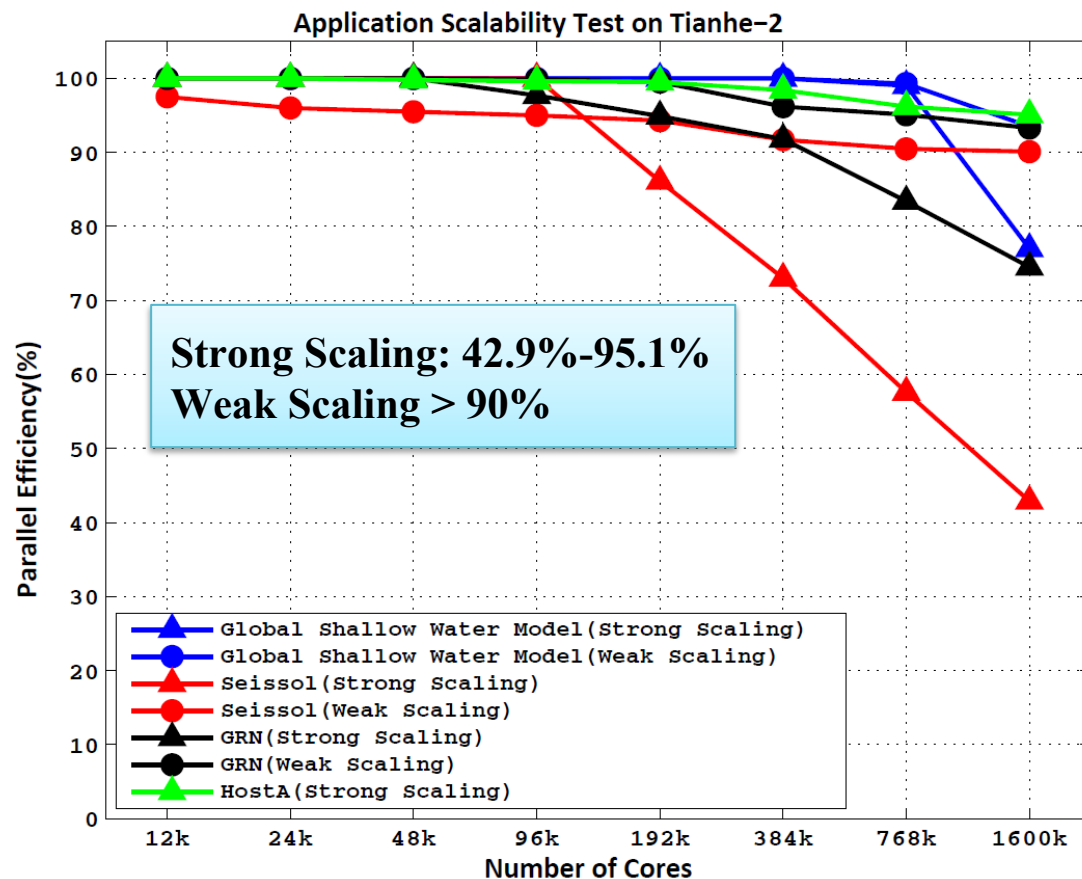
# Application Case – Bio-Medicine

天河

- Against Ebora virus 40 millions real drug molecules docking in 20hours ON TH-2

# Other Applications

☐ Other Applications

**Application Scalability Test on Tianhe-2**



Strong Scaling: 42.9%-95.1%
Weak Scaling > 90%

Legend:
- Global Shallow Water Model(Strong Scaling)
- Global Shallow Water Model(Weak Scaling)
- Seissol(Strong Scaling)
- Seissol(Weak Scaling)
- GRN(Strong Scaling)
- GRN(Weak Scaling)
- HostA(Strong Scaling)

Y-axis: Parallel Efficiency(%)
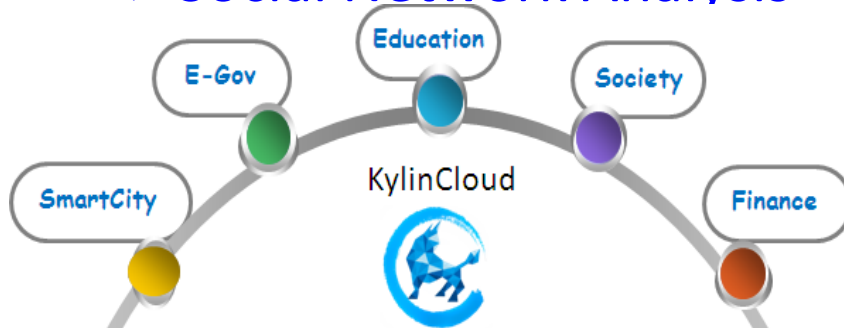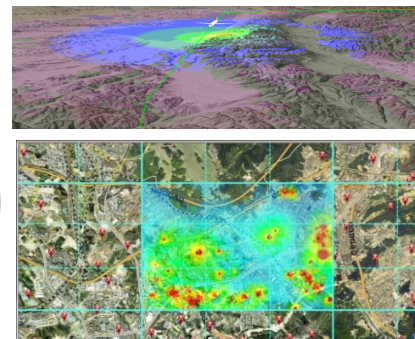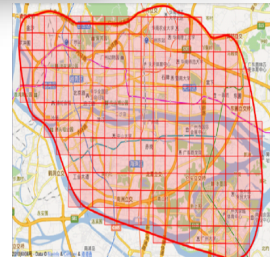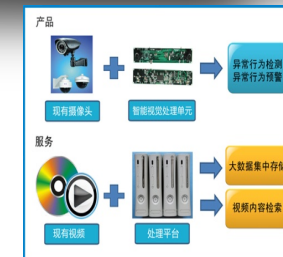X-axis: Number of Cores (12k, 24k, 48k, 96k, 192k, 384k, 768k, 1600k)

HPCL

# Cloud Computing & Bigdata

天河

□ Additional Applications

- ➢ E-Gov
- ➢ micMR
- ➢ RenderCloud
- ➢ Health care
- ➢ Smart city
- ➢ Video Processing
- ➢ Education
- ➢ Social Network Analysis

E-Gov  Education  Society  SmartCity  Finance

KylinCloud

国防科学技术大学
*National University of Defense Technology*

HPCL

# Application scale in next 5 years

天河

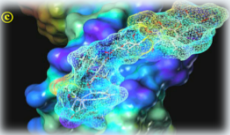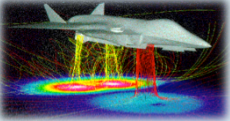| Applications | Current Scale in China | Scale in next 5 years |
|---|---|---|
| **Seismic Exploration** | 2600km$^2$ , 5km depth 217900 shots 2.2TB data | Millions of shots |
| **Genomics Research** | 2PB bioinformatics data | 100PB bio data |
| **New Energy** (Magnetic Confinement Fusion) | 2 billion ions 0.83 billion electrons | 100 billion electrons |
| **Drug Design** | 200-300ns Molecular Dynamics simulations | 10 Million molecular 1000ns/day |
| **CFD** (Aircraft Design) | 3.5 billion mesh points | 100 billion mesh points |
| **Universal Evolution** (neutrinos) | 110 billion particles | Trillion particles |
| **Smart City** (Urban Electromagnetic Spectrum Monitoring System) | Area (Guangzhou city): 200km$^2$ Grid size:1.0km*1.0km | Grid Size: 100m*100m |

# Status of Tianhe System

天河

| System | Tianhe-1A | Tianhe-2 | Tianhe-2A |
|---|---|---|---|
| System Peak(PF) | 4.7 | 54.9 | ~100 |
| Peak Power(MW) | 4.04 | 17.8 | ~18 |
| Total System Memory | 262 TB | 1.4 PB | ~3PB |
| Node Performance(TF) | 0.655 | 3.431 | ~6 |
| Node processors | Xeon X5670 Nvidia M2050 | Xeon E5 2692 Xeon Phi | Xeon or China CPU China Accelerator |
| System size(nodes) | 7,168 nodes | 16,000 nodes | ~18,000 |
| System Interconnect | TH Express-1 | TH Express-2 | TH Express-2+ |
| File System | 2 PB Lustre | 12.4PB $H^2$FS +Lustre | ~30PB $H^2$FS+TDM |

# China Accelerator

天河

## Matrix2000 GPDSP

☐ High Performance

- ➢ 64bit Supported
- ➢ ~2.4/4.8TFlops(DP/SP)
- ➢ 1GHz, ~200W

☐ High Throughput

- ➢ High-bandwidth Memory
- ➢ 32~64GB
- ➢ PCIE 3.0, 16x

# China Accelerator

□ Software stack

- ➤ OpenMP4.0
- ➤ OS, Compiler and Math Library on GPDSP
- ➤ GPDSP Driver, Communication Library

| OpenMP4.0 Application | |
|---|---|
| **Software Stack** | OpenMP4.0 Compiler system |
| | GPDSP Communication Library |
| | GPDSP Driver / GPDSP Compiler / GPDSP Math Library / GPDSP OS |
| **Node** | Intel Xeon CPU / GPDSP |

国防科学技术大学
*National University of Defense Technology*

HPCL

# HPC & Bigdata convergence Stack

天河

## Starlight Software Stack

**Application Workflow Management**

| Starlight-DSL | EHVN Ctrl | Workload monitor | Visualization system |

**HPC Software Stack**

Domain Framework — ADE

Tools | Numerical Lib

MPI | Global Array — PDE

C/C++/Fortran | OpenMP

**Bigdata Cloud Stack**

| Machine Learning | Statistical Analysis | Business Intelligence |

Hadoop | Spark

Big Data Analizing Framework

**Data Management**

Hybrid Hierarchy File system

| HDFS | SQL Database | NoSQL Database |

**Resource Management**

| System Monitor | Fault tolerant Engine | Deployment | Orchestration | Virtualization | Authentication | Scheduling |

Operating System        Kylin Cloud Platform

HPCL

# Co-design Eco-system

天河

**Application**



| Domain Models | Algorithms | Proxy Apps | Benchmarks | Data Analysis |

Solutions

Requirements

**Software**

**Bridge**

| Domain Framework | Data Management | Tools |

| MPI | OpenMP | GA | CUDA /OpenAcc | Spark | New Emerged Programing Interface |

Hybrid Runtime

Constraints

| OS | Compiler | Library | File System |

Tradeoff

**Hardware**

System Architecture

| CPU/Accelerator Hybrid Node | Memory | Interconnection | Storage Device |

# Summary

天河

☐ Heterogeneous Architecture

☐ Adaptive system and software design

☐ HPC & Bigdata Convergence

☐ Supercomputing Eco-system



国防科学技术大学
*National University of Defense Technology*

HPCL

# Thanks