



### HPC Infrastructure and Applications in Chinese Academy of Sciences

#### Xuebin CHI, Haili XIAO, Rongqiang CAO, Yining ZHAO

#### (chi@sccas.cn)

Computer Network Information Center (CNIC) Chinese Academy of Sciences (CAS) Jan. 28, 2016, Geneva, Swiss





### Outline

Supercomputing center in CNIC (SCCAS)
 HPC infrastructures in China
 Applications in CAS
 Collaborations





### Supercomputing Center in CNIC



# **Supercomputing Center (SCCAS)**

Subsidiary branch of CNIC, CAS in Beijing, 100+ staffs

Missions

- Operation and maintenance of the Supercomputing Environment of CAS (China ScGrid)
- Development of visualization, HPC application software
- HPC service provider

Our roles in the national HPC infrastructure of China

- Operation and Management Center of CNGrid (announced in 2005)
- The northern major node of CNGrid
- Management of Supercomputing Innovation Alliance



















## New Petascale Supercomputer - Era

• ERA - 元(Yuan)

**CAS HPC from T to P - new period** 

**D** Peak performance - 2.36 Petaflops

□ The 6<sup>th</sup> generation supercomputer in SCCAS

Installation

□ Site: Huairou Branch Center of CNIC

□ Two stages

■ Stage 1: announced on June 19, 2014

□Stage 2: will be announced on March, 2016











### **Huairou Branch Center of CNIC**





### Hardware @ Stage 1

- 303.4 Teraflops ( CPU + GPU/Intel Xeon Phi )
- Storage capacity 3.041 PB
- Integrated bandwidth 64.5GB/s
- 56Gbps FDR InfiniBand Interconnection
- New Gridview 3.0 Cluster Management System
- Highly efficient horizontal air-flow water cooling system



### Hardware @ Stage 2

![](_page_8_Figure_1.jpeg)

![](_page_9_Picture_0.jpeg)

![](_page_9_Picture_1.jpeg)

- Compiler , Math Libs, OpenMP, MPI
- HPC Software automatic installation tool Clussoft
- Matlab、MolCAS、Q-Chem、Amber、CHARMM、Gaussian

![](_page_9_Figure_5.jpeg)

![](_page_10_Picture_0.jpeg)

![](_page_10_Picture_1.jpeg)

![](_page_10_Figure_2.jpeg)

![](_page_10_Picture_3.jpeg)

![](_page_10_Picture_4.jpeg)

### Part II

### HPC Infrastructures in China

![](_page_11_Picture_2.jpeg)

![](_page_12_Picture_0.jpeg)

### Supercomputing Environment of CAS – China ScGrid

- Three-tier grid
  - 1 head center
  - 9 regional centers (NEW: USTC + Guangzhou)
  - **18** institution centers, **11** GPU centers
- Applications 120
  - Computational Chemistry, Physics, Material science, Life science, CFD, Industrial computing
- Status (by 2015)
  - □ User 500
  - □ #Job > 550 000
  - □ Walltime >100 000 000 CPU Hours

![](_page_12_Figure_12.jpeg)

![](_page_12_Picture_13.jpeg)

![](_page_12_Picture_14.jpeg)

![](_page_13_Picture_0.jpeg)

## **Overview of System Usage**

![](_page_13_Figure_2.jpeg)

![](_page_13_Picture_3.jpeg)

![](_page_13_Picture_4.jpeg)

息中心

![](_page_14_Picture_0.jpeg)

## **CNGrid environment**

- 14 sites
  - SCCAS (Beijing, major site)
  - SSC (Shanghai, major site)
  - NSCTJ (Tianjin)
  - NSCSZ (Shenzhen)
  - NSCJN (Jinan)
  - Tsinghua University (Beijing)
  - IAPCM (Beijing)
  - USTC (Hefei)
  - XJTU (Xi'an)
  - SIAT (Shenzhen)
  - HKU (Hong Kong)
  - SDU (Jinan)
  - HUST (Wuhan)
  - GSCC (Lanzhou)

![](_page_14_Figure_17.jpeg)

![](_page_14_Picture_18.jpeg)

![](_page_14_Picture_19.jpeg)

![](_page_15_Picture_0.jpeg)

### **CNGrid new model**

![](_page_15_Figure_2.jpeg)

![](_page_16_Picture_0.jpeg)

## **SCE - Middleware for Science Cloud**

- Developed by SCCAS
- SCE
  - Scientific computing
  - Lightweight
  - Stable
- Diveristy
  - CLI
  - Portal
  - GUI
  - API

International Patent (PCT/CN2011/071640)

![](_page_16_Figure_13.jpeg)

#### All 学院超级计算中 SCEAPI - HPC Cloud API based on SCE

- RESTful API
  - Lightweight Web Service
  - OS independent
    - Windows, Linux
    - iOS, Android
  - Language independent
    - Java, C/C++
    - PHP, Python, Ruby
    - ...
  - Support App. Community
  - Support mobile APPs

![](_page_17_Figure_12.jpeg)

![](_page_17_Picture_13.jpeg)

![](_page_18_Picture_0.jpeg)

### **SCEAPI - HPC Cloud API**

![](_page_18_Figure_2.jpeg)

### **SCEAPI - HPC Cloud API**

![](_page_19_Picture_1.jpeg)

![](_page_19_Picture_2.jpeg)

Platform as a service

![](_page_19_Picture_4.jpeg)

Infrastructure as a service

![](_page_20_Picture_0.jpeg)

### **CNGrid & ATLAS**

- CNGrid support ATLAS experiment
  - SCEAPI works as a bridge between ARC-CE middleware and CNGrid resources
  - ATLAS simulation jobs run on Chinese HPCs including TianHe-1A and ERA

![](_page_20_Figure_5.jpeg)

![](_page_21_Figure_0.jpeg)

![](_page_21_Figure_1.jpeg)

900 800 700 600 500 400 300 200 200 200 2015-10-04 2015-10-18 2015-11-19 2015-11-19 2015-12-13 2015-12-27 2016-01-10 2016-01-24

### Part III

### HPC science and engineering applications in CAS

![](_page_22_Picture_2.jpeg)

![](_page_23_Picture_0.jpeg)

### **CCFD-** parallel CFD software

![](_page_23_Figure_2.jpeg)

Aerodynamic Computation

#### Multi-Body Separation

#### Aeroelastic Flutter

![](_page_23_Picture_6.jpeg)

![](_page_24_Picture_0.jpeg)

## Scale up to over 10,000 cores

#### **CCFD-MGMB**

- Multi-block structured grids
- Implicit time stepping(via pseudo-time iteration)
- Multi-grid acceleration

# **CCFD-MBS** Chimera grids Parallel grid assembling

#### **CCFD-AE**

- •Grid deformation
- Couple with structural analysis software

![](_page_24_Figure_10.jpeg)

TianHe II test **DLR-F6** model

![](_page_24_Figure_12.jpeg)

![](_page_24_Figure_13.jpeg)

![](_page_25_Picture_0.jpeg)

### **Phase Field Simulation**

![](_page_25_Figure_2.jpeg)

Chinese Academy of Science

![](_page_26_Picture_0.jpeg)

## **Phase Field Simulation with cETD**

#### compact Exponential Time Differencing(cETD)

- •Explicit large time step
- •Stable & accurate
- •High performance on CPU+MIC
- Supported by the Intel Parallel Programming Center(IPCC) program

![](_page_26_Figure_7.jpeg)

Over 1,300GFlops(DP) on 2CPU+2MIC, 52% peak.

![](_page_26_Picture_9.jpeg)

Phase separation on  $1024^3$  grids, 3hrs with 2CPU+2MIC.

![](_page_27_Picture_0.jpeg)

### **MIC-accelerated DPD**

Flowchart of DPD simulation

Progress in IPCC

![](_page_27_Figure_4.jpeg)

![](_page_27_Figure_5.jpeg)

• Implement our DPD code on a single MIC, and achieve more than 5 times speedup than a single CPU core.

![](_page_28_Picture_0.jpeg)

# **New Energy Power Generation**

- New energy simulation system
  - new energy time series modeling
  - time series power generation simulation
  - stochastic power generation simulation
- save at least ¥10 billion every year
- Increase new energy at least 1 billion kwh
  - = saving the coal nearly 400,000 tons

reducing carbon dioxide emissions by 800,000 tons

![](_page_28_Picture_10.jpeg)

![](_page_28_Picture_11.jpeg)

Computer Network Information Center Chinese Academy of Sciences

![](_page_28_Picture_13.jpeg)

Academy of Mathematics and Systems Science Chinese Academy of Sciences

![](_page_28_Picture_15.jpeg)

![](_page_28_Picture_16.jpeg)

![](_page_29_Picture_0.jpeg)

# **New Energy Power Generation**

- CMIP ( Chinese Mixed Integer Programming Solver )
- Mid long term wind power forecast
- Power network topology visualization

![](_page_29_Figure_5.jpeg)

![](_page_29_Picture_6.jpeg)

# **EXAMPLE C\_Tangram: a Charm++-Based Parallel Framework for Cosmological Simulations**

- Hide complex parallel technologies.
- Provide a platform for composing components together into a complex application.

![](_page_30_Figure_3.jpeg)

1. Modularity

componentization collaboration

- 2. Runtime Adaptivity Fault Tolerance Load Balance
- 3. Domain Specification

Cosmological hydrodynamics N-body

![](_page_30_Picture_9.jpeg)

![](_page_30_Picture_10.jpeg)

FIGURE 1. MULTI-LAYERS DESIGN OF SC\_TANGRAM

# 「国科学院超级计算中・SC\_Tangram: a Charm++-Based Parallel Framework for Cosmological Simulations

Domain Specific Data Types

![](_page_31_Figure_2.jpeg)

![](_page_31_Figure_3.jpeg)

![](_page_31_Figure_4.jpeg)

DATA TYPES "GF" FOR UNIFORMED MESH AND "GP" FOR PARTICLES IN MESH ON CPU-CLUSTERS

Applications

![](_page_31_Figure_7.jpeg)

![](_page_31_Figure_8.jpeg)

![](_page_31_Figure_9.jpeg)

EXECUTION TIME OF ONE STEP ON THE SCALE OF  $134217728 \ \mbox{particles}$  on the mesh

![](_page_32_Picture_0.jpeg)

# **CAS** earth system model

- We participate in the development of CAS earth system model
- Run CMIP6 experiments on "era"

![](_page_32_Figure_4.jpeg)

![](_page_32_Picture_5.jpeg)

![](_page_32_Picture_6.jpeg)

![](_page_33_Picture_0.jpeg)

# High performance integrated computing platform

- Assemble, compile and run Scripts
- Cas-coupler
  - 2-D coupler(based on MCT)
  - 3-D coupler
  - Coupler creator
- Standard and unified component models
  - Parallel and optimization algorithms for IAP AGCM,CoLM,LICOM,RIEMS,...
  - Standard and unified rules and interface
- Tools and Library
  - High performance communication library
  - Parallel I/O library
  - Performance debugging tools

#### Prototype system of ESM simulation facilities

![](_page_33_Figure_15.jpeg)

![](_page_33_Picture_16.jpeg)

![](_page_33_Picture_17.jpeg)

![](_page_34_Picture_0.jpeg)

# Application of HPSEPS in first principle calculation software MESIA

(b)

HPSEPS, a parallel eigenproblem solver developed by SC,CAS is adopted in a multiscale first principle calculation software MESIA, developed by the Key Lab. Of Quantum Information, CAS.

MESIA produced the correct energy sequences for B3(Zinc Blend) and B4 (Wurtzite)

(a)

• B20 cluster

![](_page_34_Figure_5.jpeg)

	Plane wave	DZP	SIESTA
double ring (a)	0.00 eV	0.00 eV	0.00 eV
candidate1 (b)	2.74 eV	2.52 eV	-0.13 eV

The energy sequence calculated using SIESTA is incorrect, MESIA gave the correct one.

We developed MPI-GPU eigensolver for dense eigenproblem Computational Throughput: 16GPUs= 512CPUs

![](_page_34_Picture_9.jpeg)

![](_page_34_Picture_10.jpeg)

![](_page_35_Picture_0.jpeg)

# Large scale three dimensional fragment method on GPU

**Chemical accuracy** 

**RPA, DMFT** 

hybrid

Meta-GGA

GGA

LDA

1. Accuracy

(climb Jacob's ladder)

- 2. Temporal scale (from fs to seconds) (new algorithms, like accelerated MD)
- 3. Size scale(mesoscale problems)

(divide & Conquer methods)

![](_page_35_Figure_7.jpeg)

Boundary effects are (nearly) cancelled out between the fragments

$$System = \sum_{i,j,k} \left\{ F_{222} + F_{211} + F_{121} + F_{112} - F_{221} - F_{212} - F_{122} - F_{111} \right\}$$

![](_page_35_Picture_10.jpeg)

Titan GPU: 88% of total computing powerBut NO plane wav e code on GPU.

### LS3DF: Linear Scaling Three Dimensional Frag ment Method

Collaborate with Lin-Wang Wang, LBNL

This project is supported by INCITE program and CSC

![](_page_35_Picture_15.jpeg)

![](_page_35_Picture_16.jpeg)

![](_page_36_Picture_0.jpeg)

### Large scale three dimensional fragment method on GPU

![](_page_36_Figure_2.jpeg)

![](_page_36_Figure_3.jpeg)

LS3DF data distribution

![](_page_36_Figure_5.jpeg)

LS3DF algorithm compared with LDA algorithm

#### On Titan Supercomputer:

3877 atom Si system, 1500 computing nodes(total 24000 CPU cores) compared with 1500 GPU cards, LS3DF\_GPU h as a speedup of 10.5x.

![](_page_36_Picture_9.jpeg)

![](_page_36_Picture_10.jpeg)

LS3DF-GPU speedup compared with CPU code

![](_page_37_Picture_0.jpeg)

### Fast Parallel Direct Solver for Large linear system

#### HSS algorithm

✓ hierarchically semiseparable matrix, Chandrasekaran, Gu, Xia, et al

✓ three steps: HSS compression, ULV factorization & ULV solver

✓ Complexity: O(kN<sup>2</sup>) for step 1, O(kN) for step 2&3; storage: O(kN), k: block rank

✓ Recursive Low-rank compression by tree (c1&c2: children of node j)

$$D_{j} = \begin{pmatrix} D_{c_{1}} & U_{c_{1}}B_{c_{1}}V_{c_{2}}^{T} \\ U_{c_{2}}B_{c_{2}}V_{c_{2}}^{T} & D_{c_{2}} \end{pmatrix}, \quad U_{j} = \begin{pmatrix} U_{c_{1}}R_{c_{1}} \\ U_{c_{2}}R_{c_{2}} \end{pmatrix}, \quad V_{j} = \begin{pmatrix} V_{c_{1}}W_{c_{1}} \\ V_{c_{2}}W_{c_{2}} \end{pmatrix}$$

$$\boxed{D_{1} \quad U_{1}B_{1}V_{2}^{T}} \quad B_{3} \quad V_{6}^{T} \\ U_{3} \quad U_{3} \quad U_{3} \quad U_{3} \quad B_{6} \quad G_{6} \\ R_{1},W_{1} \quad B_{1} \quad R_{2},W_{2} \quad R_{4},W_{4} \quad B_{4} \quad S_{5} \\ 1 \quad B_{2} \quad 2 \quad 2 \quad 4 \quad B_{5} \quad 5 \quad 5 \\ U_{1},V_{1} \quad U_{2},V_{2} \quad U_{4},V_{4} \quad U_{5},V_{5} \\ D_{1} \quad D_{2} \quad D_{4} \quad D_{5} \\ \hline \Psi \equiv H \neq \mathbb{R} \\ 1 \text{ or mode the two sets of t$$

HSS matrix structure & HSS tree

![](_page_38_Picture_0.jpeg)

### **Result: HSS solver vs ScaLapack**

Dense linear system:

$$\left[E_{inc}(u_m)\right] = \left[Z_{mn}\right]\left[j_n\right], Z_{mn} = \frac{k\eta}{4} d_n H_0^{(2)}(kR_{mn})$$

✓ in which H is cylinder Hankel function, inc means incident field
 ✓ N: 32768, HSS tree level: 8, block rank: 32

✓ Total runtime, compared with pzgesv of ScaLapack (using MKL)

![](_page_38_Figure_6.jpeg)

![](_page_39_Picture_0.jpeg)

![](_page_39_Picture_1.jpeg)

### Collaboration inside and outside China

### CAS Computational Science Application Research Center

Approved by CAS and launched in 2014

Computational Science		Algorithm Optimization &Implementation		Heterogeneous Platforms		Software Production Line		
Academy of Mathematics & Systems Science	Dali Institu Chem Phys	an Ite of Iical Iics	Institute of Atmospheric Physics	National Astronomical Observatories	Bei Instit Gen	jing ute of ome	Institute of Modem Physics	
Computer Network Information Center								

![](_page_41_Picture_0.jpeg)

### **Supercomputing Innovation Alliance**

 The Alliance, approved by the Ministry of Science and Technology, is the Industry-University-Research-Application cooperation organization. The Alliance was established in September 25, 2013, initiated by the national or local Supercomputing Centers, high-performance computing application research institutes, and related enterprises of total 55 units.

![](_page_41_Picture_3.jpeg)

![](_page_41_Picture_4.jpeg)

![](_page_41_Picture_5.jpeg)

![](_page_41_Picture_6.jpeg)

![](_page_41_Picture_7.jpeg)

![](_page_42_Picture_0.jpeg)

### Supercomputing Innovation Alliance - Organization Structure

![](_page_42_Picture_2.jpeg)

Honorary chairman: Academician Jin Yilian

![](_page_42_Picture_4.jpeg)

Secretary-General: Research Professor Chi Xuebin

![](_page_42_Picture_6.jpeg)

#### **Chairman: Professor Qian Depei**

![](_page_42_Picture_8.jpeg)

Vice Secretary-General: Research Professor Xie Xianghui

### **Relying institutions: CNIC, CAS**

![](_page_42_Picture_11.jpeg)

![](_page_42_Picture_12.jpeg)

![](_page_42_Picture_13.jpeg)

### Supercomputing Innovation Alliance - Membership (50+)

![](_page_43_Picture_1.jpeg)

# The second of th

- The only Intel Parallel Computing Center (IPCC) in Mainland China
- Intel & CNIC, started in Apr. 2015
- Focusing on applications using MIC

![](_page_44_Picture_4.jpeg)

![](_page_45_Picture_0.jpeg)

# **CHANGES 2016**

- CHANGES (CHinese-AmericaN-German E-Science and cyberinfrastructure Workshop)
  - □ JSC-NCSA-CNIC collaborations
  - □ Since 2012
  - □ HPC/BigData/Vis.
  - **CHANGES16** 
    - Fall of 2016, Juelich, Germany
    - Preparation meeting at ISC15

![](_page_45_Picture_9.jpeg)

![](_page_45_Picture_10.jpeg)

![](_page_45_Picture_11.jpeg)

![](_page_45_Picture_12.jpeg)

# Thank you!

chi@sccas.cn