# Large-Scale Data Visualization Applications on Tianhe-1A



## National Supercomputer Center in Tianjin
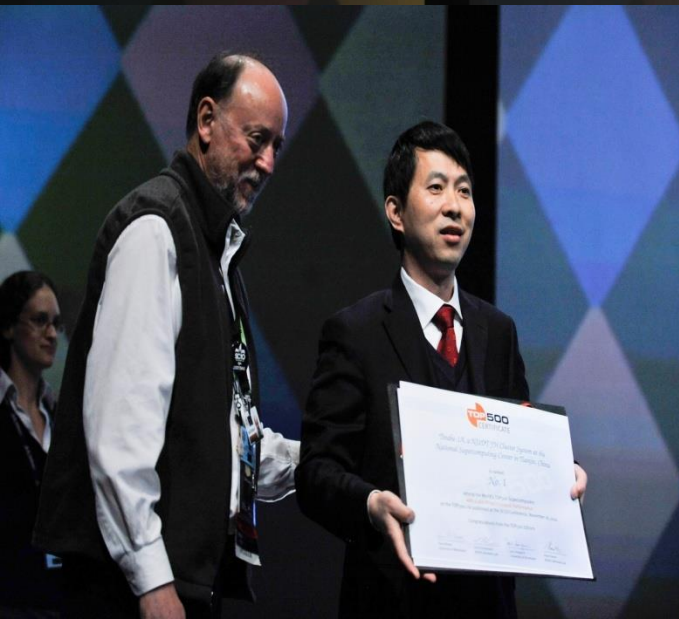
Liu Guangming

# Outline

- TH-1A system and its application

- Large-Scale data Visualization

  - ➢ Large-Scale Flow Visualization

  - ➢ Multi source geological data visualization graphics engine——OpenProbe

- Summary

**Tianhe -1A supercomputer got the world's top 500 ranked first at November 16, 2010 .**

## 36th List: The TOP10

| Rank | Site | Manufacturer | Computer | Country | Cores | Rmax [Tflops] | Power [MW] |
|------|------|--------------|----------|---------|-------|---------------|------------|
| 1 | National SuperComputer Center in Tianjin | NUDT | Tianhe-1A NUDT YH MPP, Xeon 6C, NVidia | China | 186,368 | 2,566 | 4.04 |
| 2 | Oak Ridge National Laboratory | Cray | Jaguar Cray XT5, HC 2.6 GHz | USA | 224,162 | 1,759 | 6.95 |
| 3 | National Supercomputing Centre in Shenzhen | Dawning | Nebulae TC3600 Blade, Intel X5650, NVidia Tesla C2050 GPU | China | 120,640 | 1,271 | 2.58 |
| 4 | GSIC, Tokyo Institute of Technology | NEC/HP | TSUBAME-2 HP ProLiant, Xeon 6C, NVidia, Linux/Windows | Japan | 73,278 | 1,192 | 1.40 |
| 5 | DOE/SC/ LBNL/NERSC | Cray | Hopper Cray XE6, 6C 2.1 GHz | USA | 153,408 | 1,054 | 2.91 |
| 6 | Commissariat a l'Energie Atomique (CEA) | Bull | Tera 100 Bull bullx super-node S6010/S6030 | France | 138.368 | 1,050 | 4.59 |
| 7 | DOE/NNSA/LANL | IBM | Roadrunner BladeCenter QS22/LS21 | USA | 122,400 | 1,042 | 2.34 |
| 8 | University of Tennessee | Cray | Kraken Cray XT5 HC 2.36GHz | USA | 98,928 | 831.7 | 3.09 |
| | Forschungszentrum | IBM | Jugene Gene/P Solution | Germany | 294,912 | 825.5 | 2.26 |
| | | | Cielo 2.4 GHz | USA | 107,152 | 816.6 | 2.95 |

3

# TH-1A  system and its application

- **TH-1A  main hardware parameters :**
  - Computation：4.7 petaflop
    - ✓ Number of computing nodes：7500
    - ✓ Computing node partition：4 computing  partitions according to 4 storage partitions
  - Storage：two parts: On line and Near Line
    - ✓ On Line：4 sets Lustre, each set corresponds to a storage partition, the capacity are 430TB, 420TB, 1.3PB, 700TB respectively.
    - ✓ Near  Line：dual copy storage，available capacity 4PB
    - ✓ There was one set Lustre before：430 TB

- **In addition to the TH-1A, we also constructed**
  - TH Cloud Computing Center
  - TH electronic government affairs center
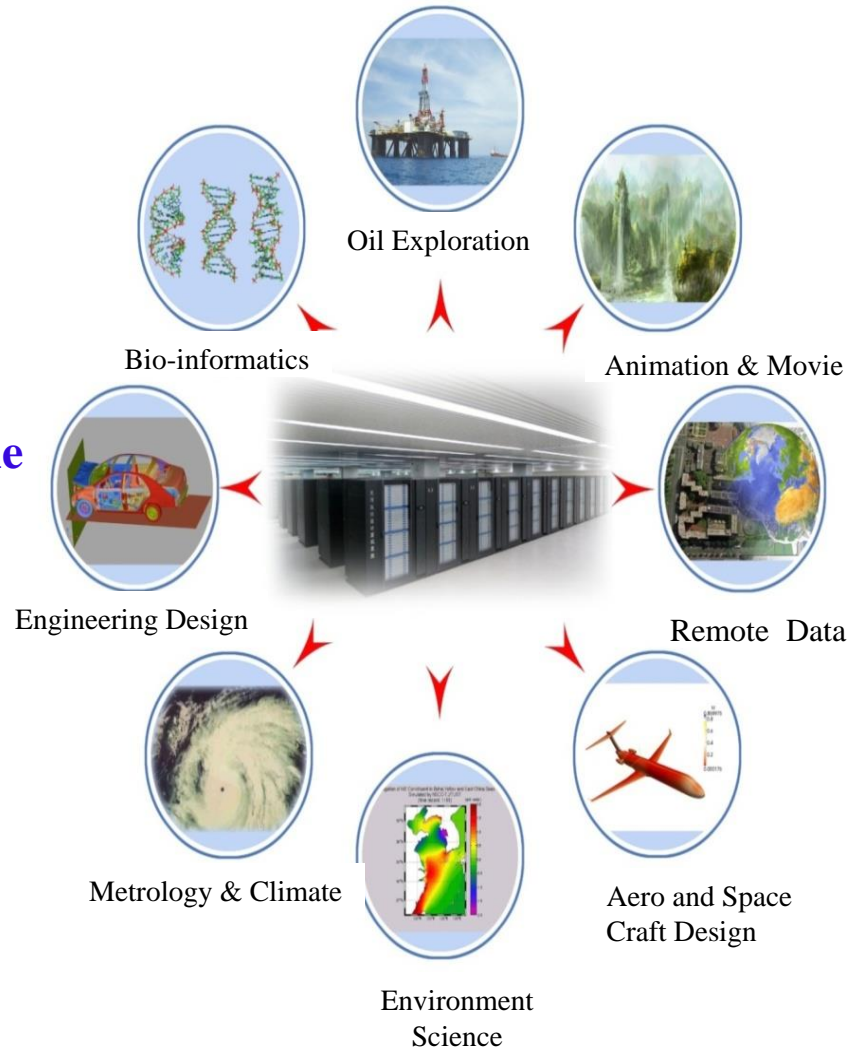  - TH big data processing environment

# HPC Application & Innovation Alliance

NSCC

Cooperation with Institutes

Cooperation with Universities

Collaboration with Companies

China First Heavy Industries

CNPC BGP

718th CSIC

Tsinghua University

Fudan University

Nan Kai University

Jilin University

SAMSUNG

Peking University

State Oceanic Administration

IoM, CAS

DUT

BGI

IoB, CAS

HKU

HKPU

SIMM, CAS

NUDT, CAS
NSCC-TJ

IoP, CAS

AoMMS

CUP

Tianjin University

CATRC

IMR, CAS

SINOPEC GRI

USTC

LASG, CAS

Great Wall Motors

Tianjin FAW Automobile

Zhejiang University

Nanjing Hydraulic Research Institute

Xi'an Jiaotong University

Shanghai Jiaotong University

Northeastern University

Tianjin Institute of Pharmaceutical Research

China Oilfield Services Co. Limited

National Animation Industry Park

中华人民共和国科学技术部
The Ministry of Science and Technology of the People's Republic of China

中华人民共和国国家发展和改革委员会
National Development and Reform Commission

国家自然科学基金委员会
National Natural Science Foundation of China

# Application Domains and some typical applications

■ Some typical application domains

➢ Geophysics：oil exploration

➢ Environmental Science ： Oceans, weather and climate

➢ Aerodynamics ： Aircraft

**These three types of application can use the whole system of computing resources.**

➢ Life Sciences：Gene, protein, brain science

➢ Engineering simulation

Oil Exploration

Bio-informatics

Animation & Movie

Engineering Design

Remote Data

Metrology & Climate

Aero and Space Craft Design

Environment Science

# Typical application： oil exploration seismic data processing



- GeoEast：Developed a data processing software for oil seismic exploration with independent intellectual property rights

- Completed a number of data processing tasks, the typical parameters are ：

  - Work area：2,600 Km$^2$

  - Origin data volume：2.2TB

  - Processing technique ：RTM

  - Imaging data volume ：27GB

  - One of Imaging interpretation technology ：Visualization

**Imaging of underground geological structure**

■ System usage ：＞85%

■ Data storaged ：＞2.5PB

■ Running years：More than 4 years: (2010.12—2015.4)



Job status on Tianhe-1A from 2015/3/20 to 2015/4/30

# Projects and users

- **TH-1A supports national projects > 800 items**

  - ➢ **NSFC projects > 600 items**
  - ➢ **863, 973 projects > 100 items**
  - ➢ **Other key project (MIIT, NDRC, CNPC, CNOOC, etc) > 40 items**
  - ➢ **International collaboration projects > 10 items**

- **Users all over China, and the number of user teams is more than 600 up to 2014**

- **NSCC-TJ is an open public technology service platform**



**User distribution graph**

# Outline

■ TH-1A  system and its application

■ Large-Scale data Visualization

➢ Large-Scale Flow Visualization

➢ Multi source geological data visualization graphics engine——OpenProbe

■ Summary

# Large-Scale Flow Visualization

## Limitation of Previous Flow Visualization Methods

Flow data usually includes multiple variables, but their corresponding analysis methods are lacking, especially the integrated analysis of vector fields with scalar fields.

Existing ensemble visualization methods focus on the scalar fields, while the comparison of vector fields are in need in scientific domain.

New flow field analysis methods equipped with scalable computation

1.  **Coupled Ensemble Flow Line Advection and Analysis (eFLAA)**

    vector field + ensemble analysis

2.  **Scalable Lagrangian-based Attribute Space Projection (LASP)**

    scalar/vector field + multivariate analysis

3.  **Latent Dirichlet Allocation Based Unsteady Flow Analysis (FLDA)**

    scalar/vector field + multivariate analysis

4.  **Advection-Based Sparse Data Management for Visualizing Unsteady Flow**

    a fundamental data management to support flow-related analysis

Hanqi Guo, Xiaoru Yuan, Jian Huang, and Xiaomin Zhu, "Coupled Ensemble Flow Line Advection and Analysis." IEEE Transactions on Visualization and Computer Graphics (Vis '13), 19(12):2733–2742, 2013.

# Pipeline in Concept



- Ensemble data (large, 76 GB)
- Field line data (much larger than ensemble data, 5.8 TB)
- Variation field (small, less than 1 GB)
- Filtered lines (even smaller)

# Pipeline of the Parallel System

- Both data scale and problem size are often too large to handle in practice
- A streamed data management mechanism is used to make the system scalable, given the memory limits

# Application – GEOS-5 Simulation

- The metric: the differences of locations / $CO_2$ concentration along the pathline
- Findings
  - The variation of the wind field is high in the north hemisphere
  - However, The $CO_2$ difference is higher in south hemisphere and some places in the north
  - $CO_2$ concentration is not sensitive to wind in above regions

Hanqi Guo, Fan Hong, Qingya Shu, Jiang Zhang, Jian Huang, and Xiaoru Yuan, "Scalable Lagrangian-based Attribute Space Projection for Multivariate Unsteady Flow Data." In *Proceedings of IEEE Pacific Visualization Symposium (PacificVis 2014)*, pages 33-40, Yokohama, Japan, Mar. 4–7, 2014.

# Attribute Space Projection

**Eulerian-based Attribute Space Projection**

Data samples →
High-dimensional vector in attribute space →
Eulerian-based Attribute Space Projection →
Eulerian-based Attribute Space Projection (EASP)

**Lagrangian-based Attribute Space Projection**

Pathlines starting from data samples →
Pathlines in attribute space →
Lagrangian-based Attribute Space Projection →
Lagrangian-based Attribute Space Projection (LASP)

- It is unrealistic to implement LASP with a serial visualization pipeline

  ➢ The complexities of both particle tracing and projection are prohibitive

  ➢ The intermediate data is overwhelmingly large
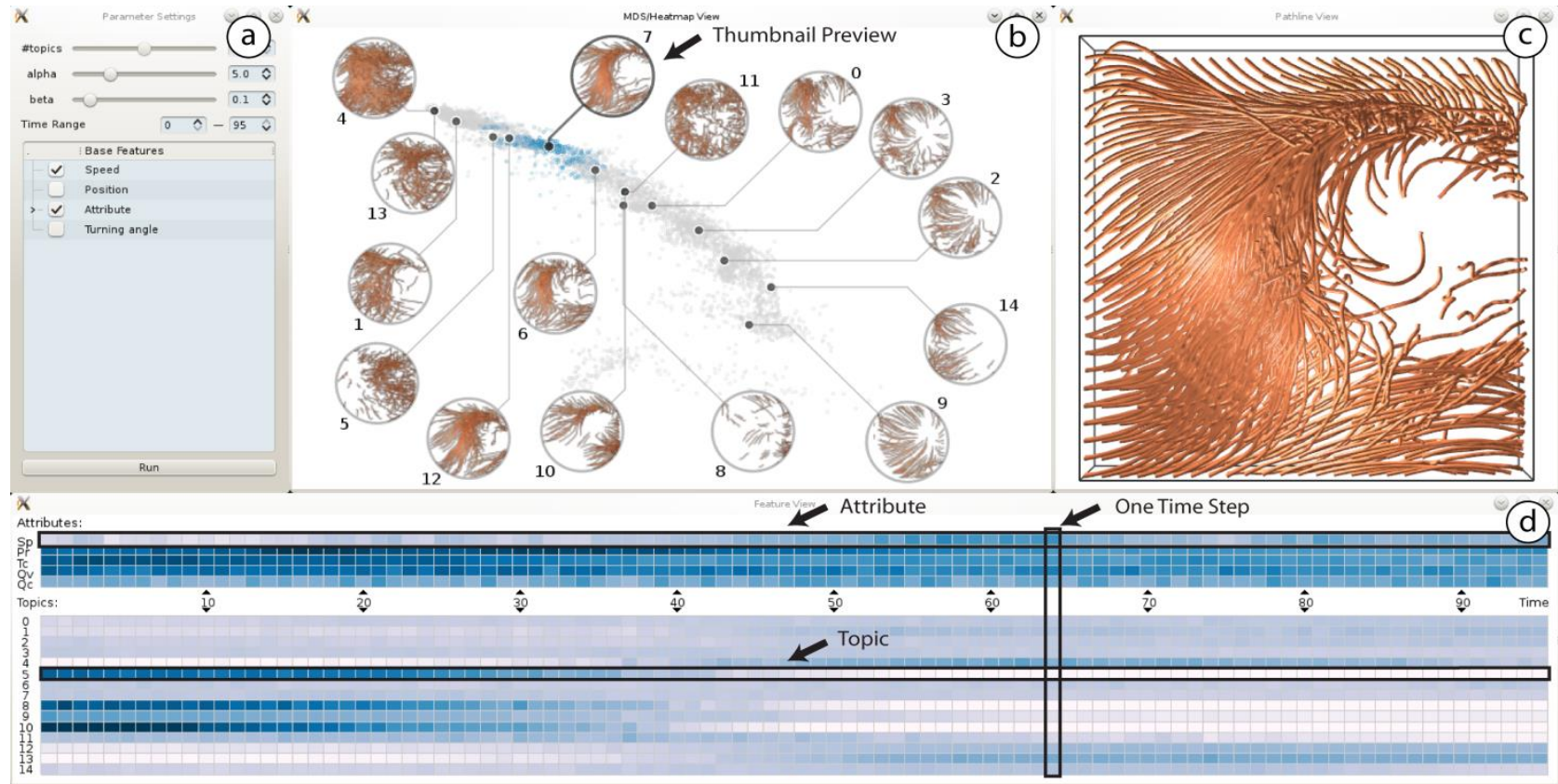
- The solution: integration of DStep and SPMDS

Attributes along pathlines are significantly different in the two group of selected features

# 3. FLDA: Latent Dirichlet Allocation Based Unsteady Flow Analysis



Fan Hong, Chufan Lai, Hanqi Guo, Enya Shen, Xiaoru Yuan, Sikun Li. "FLDA: Latent Dirichlet Allocation Based Unsteady Flow Analysis." In IEEE VIS 2014.

# LDA Topic Model vs Flow LDA Model

LDA (Latent Dirichlet Allocation) is a widely used topic model in text mining.
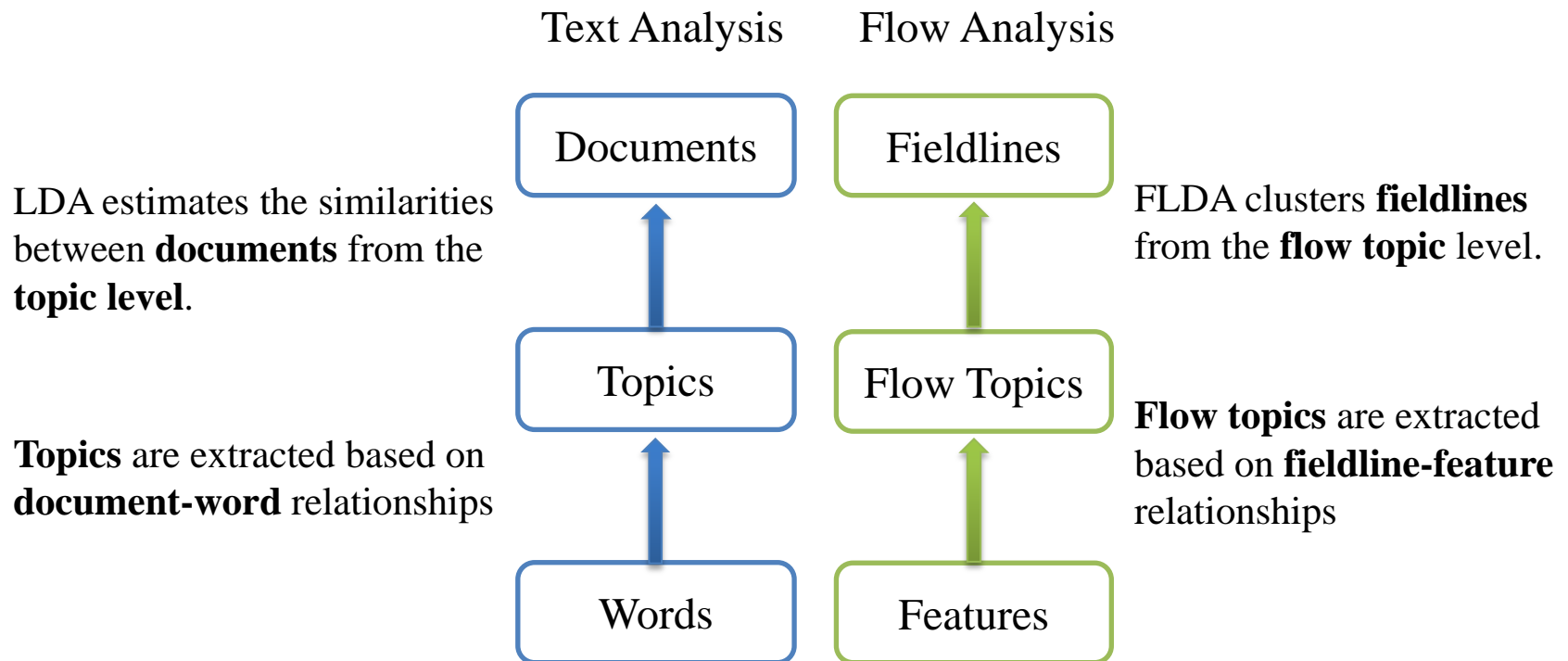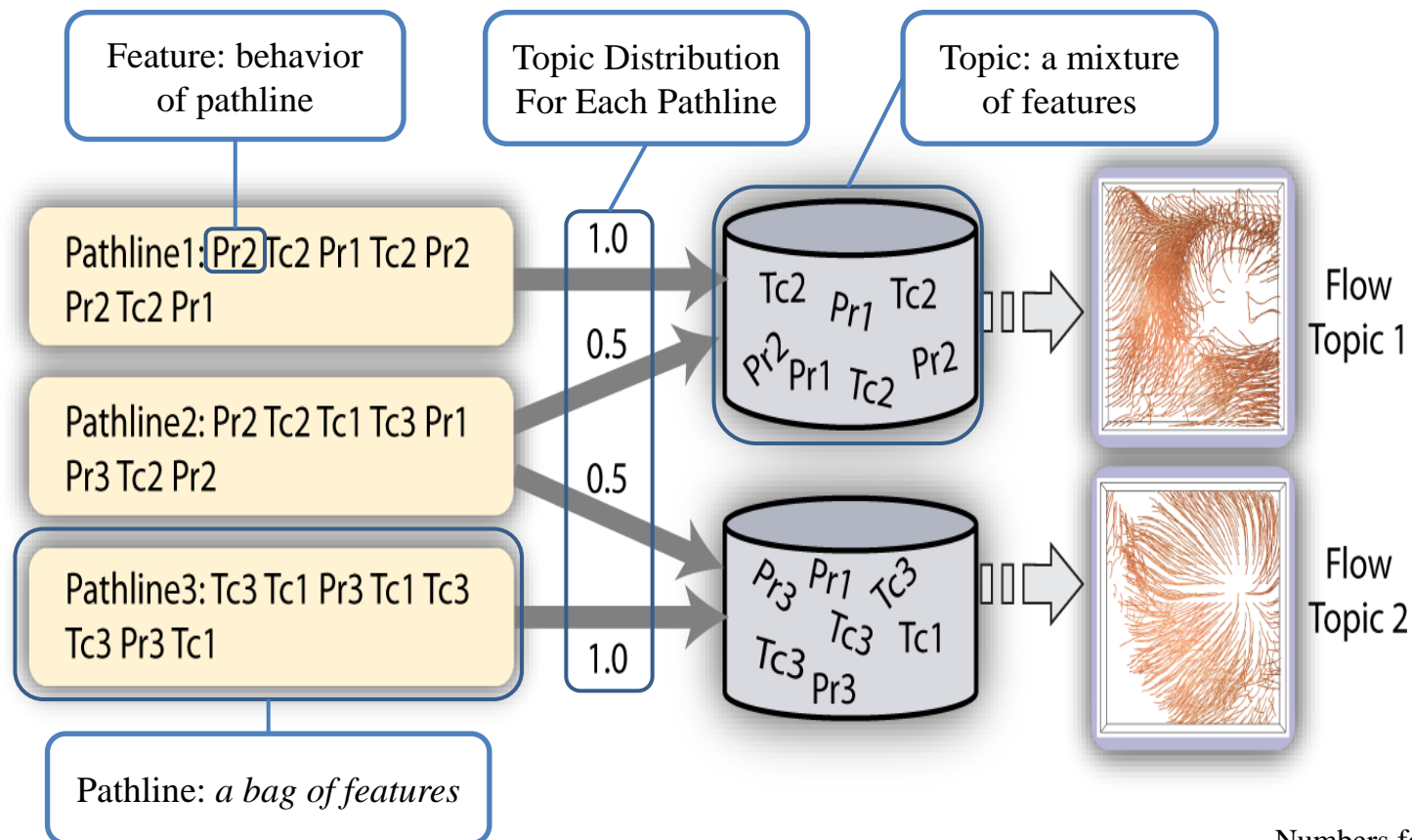
Text Analysis    Flow Analysis

LDA estimates the similarities between **documents** from the **topic level**.

**Topics** are extracted based on **document-word** relationships

| Documents | Fieldlines |
|-----------|-----------|
| Topics | Flow Topics |
| Words | Features |

FLDA clusters **fieldlines** from the **flow topic** level.

**Flow topics** are extracted based on **fieldline-feature** relationships

# Illustration of FLDA Model

**Input:**
- Definition of features
- Bags of features for each pathline

**Output:**
- Topic distribution for each pathline
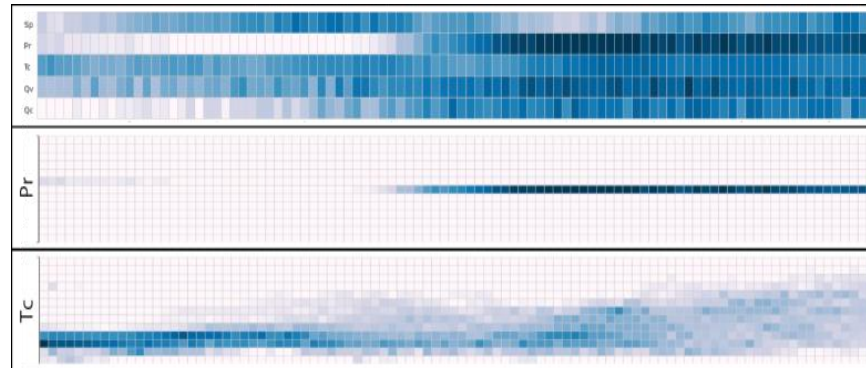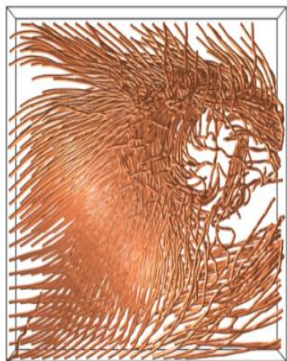- Feature distribution for each topic



Feature: behavior of pathline

Topic Distribution For Each Pathline

Topic: a mixture of features

Pathline1: Pr2 Tc2 Pr1 Tc2 Pr2 Pr2 Tc2 Pr1

Pathline2: Pr2 Tc2 Tc1 Tc3 Pr1 Pr3 Tc2 Pr2

Pathline3: Tc3 Tc1 Pr3 Tc1 Tc3 Tc3 Pr3 Tc1

1.0
0.5
0.5
1.0

Tc2 Pr1 Tc2 Pr2 Pr1 Tc2 Pr2

Pr3 Pr1 Tc3 Tc3 Tc1 Tc3 Pr3

Flow Topic 1

Flow Topic 2

Pathline: *a bag of features*

Numbers for illustration only.

Steps:

1. Define features to represent behaviors of interest.
2. Generate pathlines, and construct bags of features.
3. Feed feature bags into classical LDA model.
4. Visualize and analyze output topics and distributions.

- Pathlines advect from hurricane eye to the periphery.
- Pathlines are similar in attributes pressure and temperature, in the first half of advection.
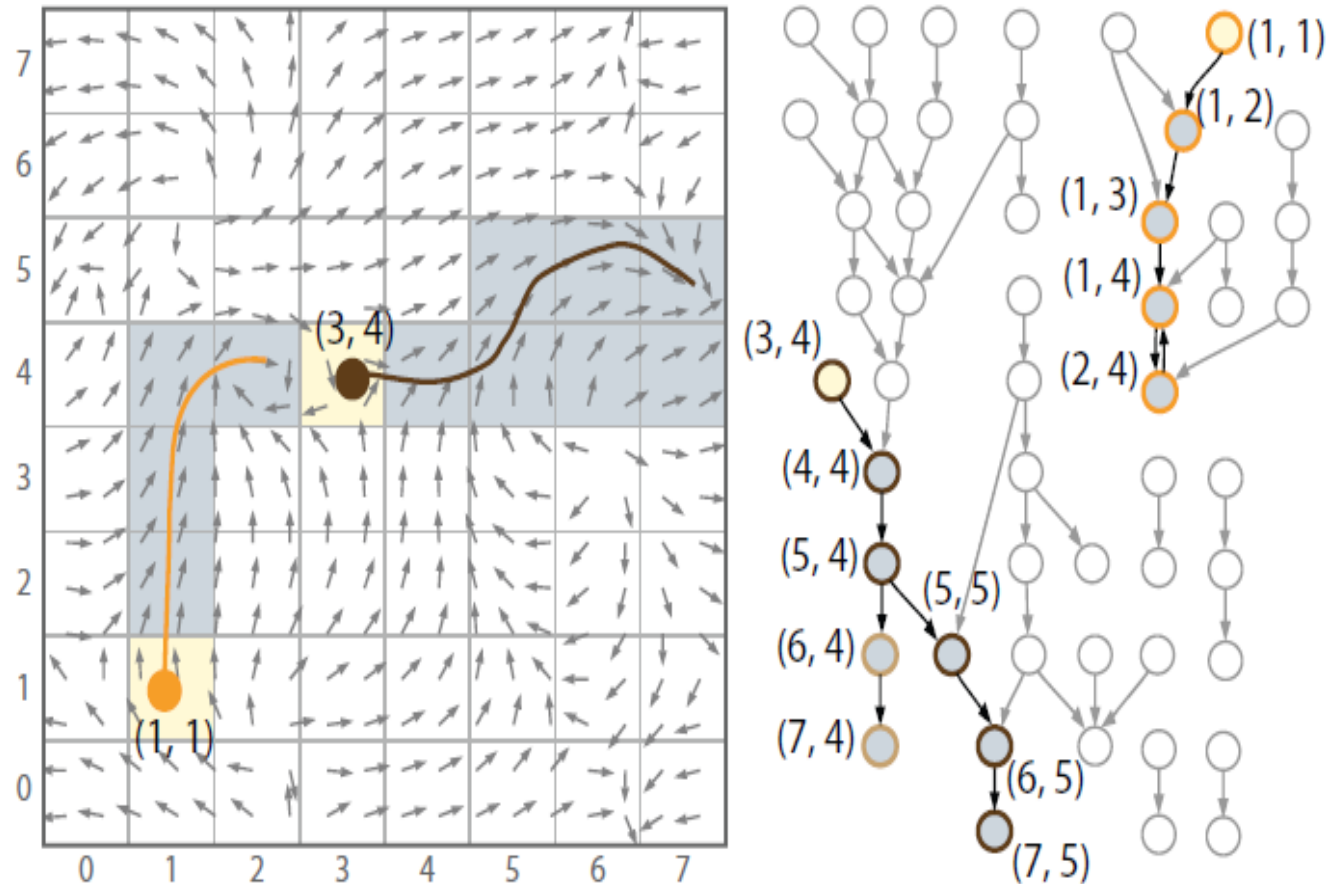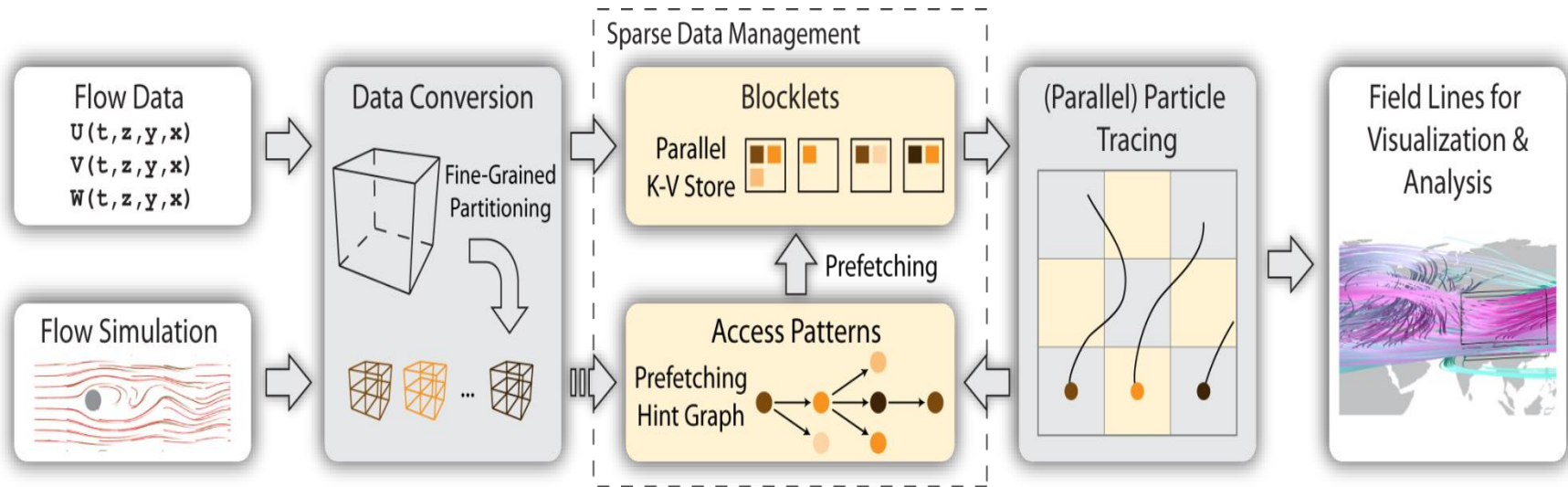- Pathlines have an increasing pressure, and stable temperature.



- Pathlines form a clockwise circulation around hurricane eye.
- Pathlines have focused pressure values in the last half of time.
- The pressure of pathlines are changing for focuses values to dispersed one.

Hanqi Guo, Jiang Zhang, Richen Liu, Lu Liu, Xiaoru Yuan, Jian Huang, Xiangfei Meng, and Jingshan Pan. "Advection-based Sparse Data Management for Visualizing Unsteady Flow." In IEEE VIS 2014.
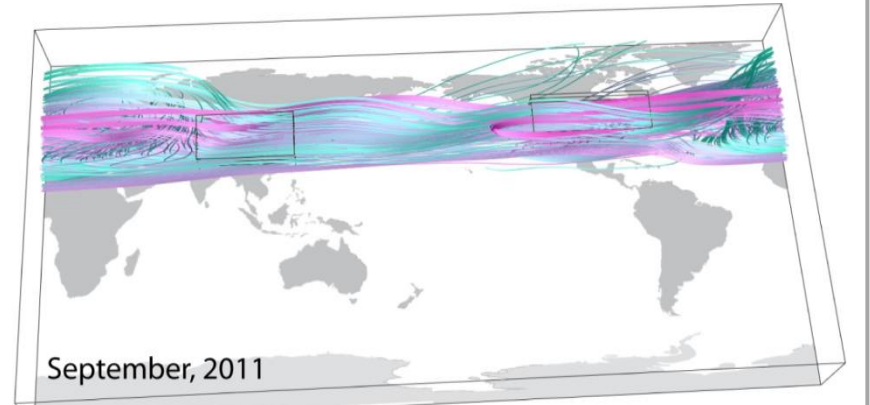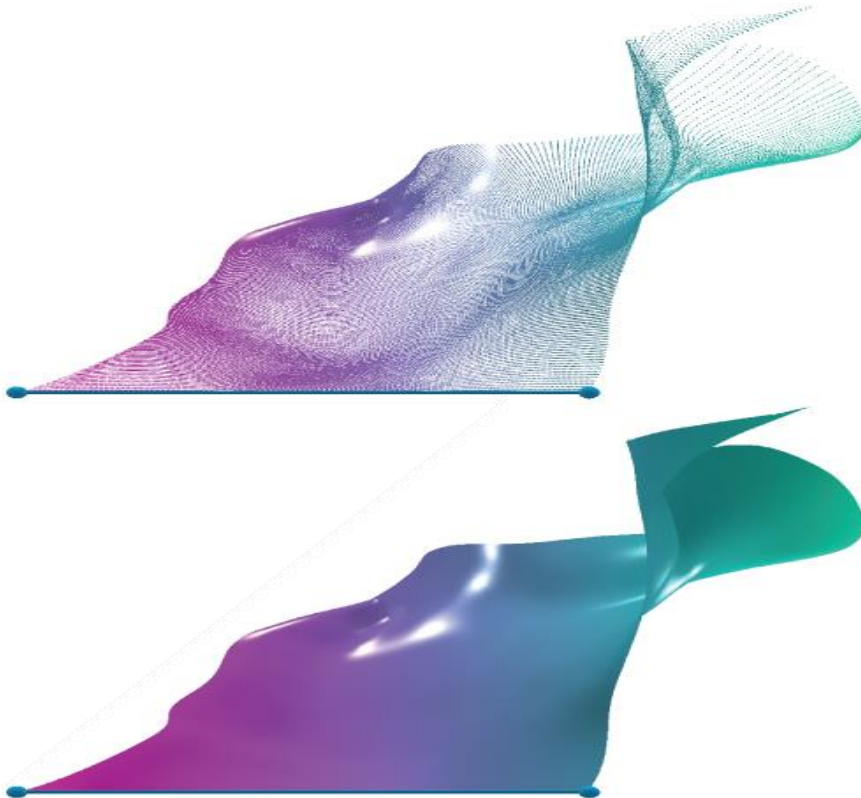
# Sparse Data Management



## Solution

- Data partition on granularity of blocklets
- Parallel key-value store based data management
- High-efficient data prefetching

## Benefits

- Enable large-scale unsteady flow analysis while requiring a very limited amount of hardware resources.
- Improve both performance and scalability of the naive task-parallel particle tracing

# Applications

## Streak Surface Computation

- Streak surfaces dipicts the flow field over the entire lifetime by continuously releasing particles from given seed curves.
- Data: TB-scale turbulence simulation data
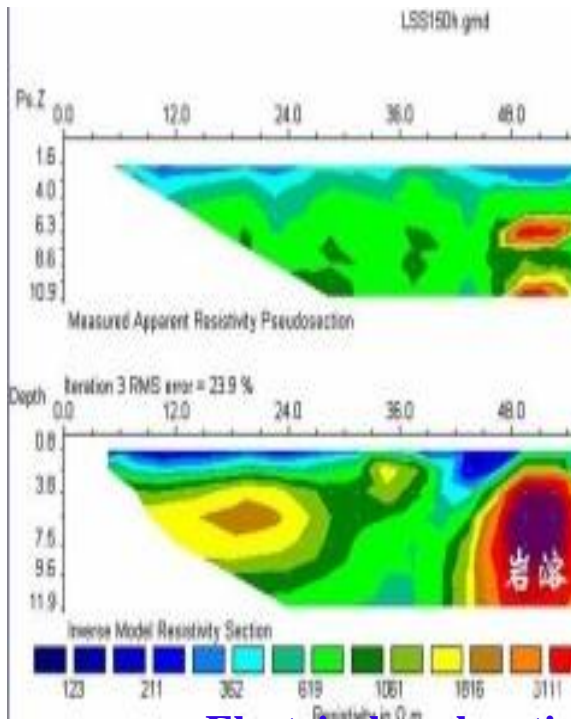  - curvilinear grid, with spatial resolution 1024x1024x720
  - 100 time steps





September, 2011



September, 2012

## Origin-Destination Query

- Study advection of massless particles (such as pullutants, etc.) in flow fields.
- Data: GEOS-5, global climate simulation data

# Outline

■   TH-1A  system and its application

■   Large-Scale data Visualization

    ➢  Large-Scale Flow Visualization

    ➢  Multi source geological data visualization graphics engine——OpenProbe
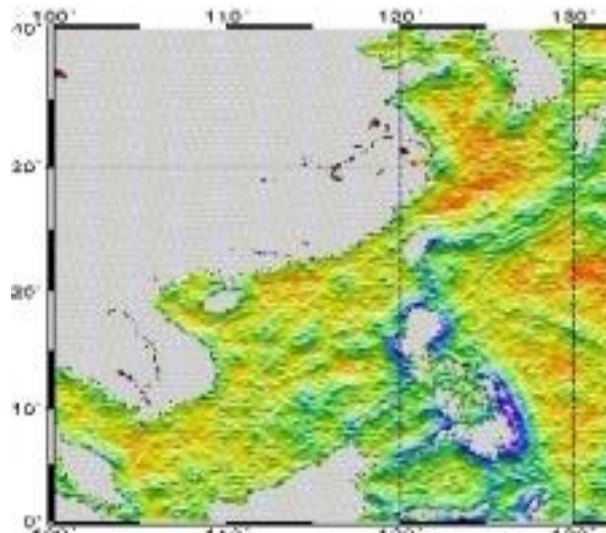
■   Summary

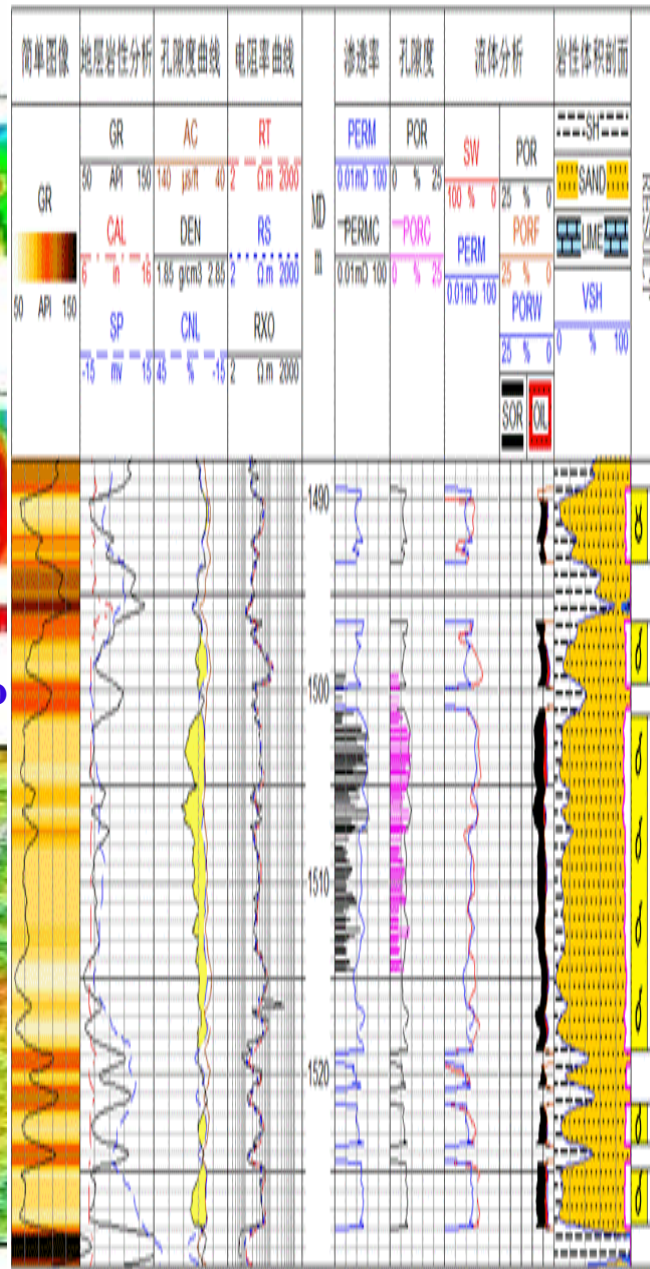# Multi-source geologic data visualization graphics engine—OpenProbe

- Supports  seismic, logging, gravity , magnetic, electric and other   data

- Supports Windows, Linux and other operating systems

- Supports multi-touch interactive operation mode

- Supports plug-in type secondary development

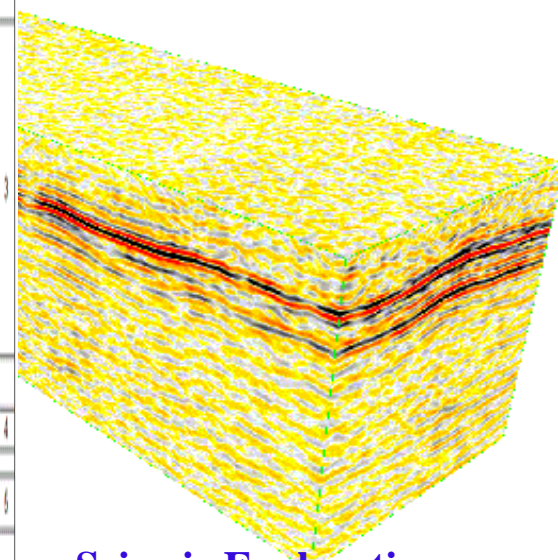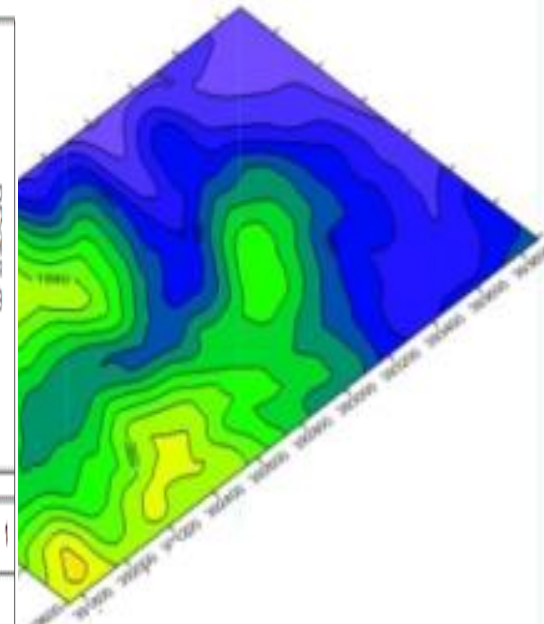- Supports large-scale data 3D stereoscopic visualization
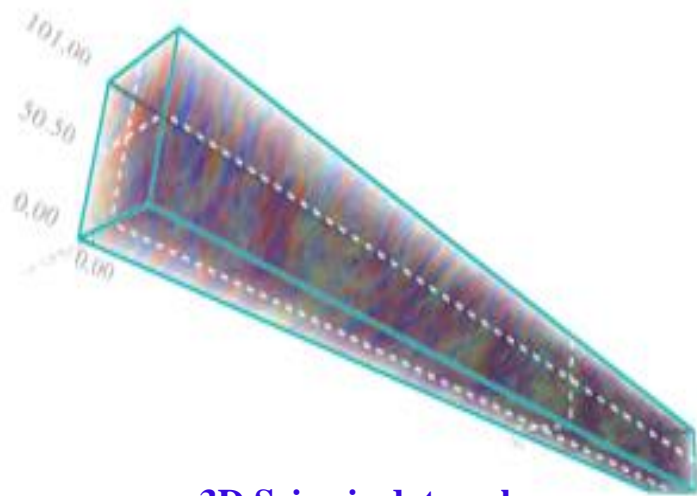
**Electrical exploration**

**Magnetic exploration**
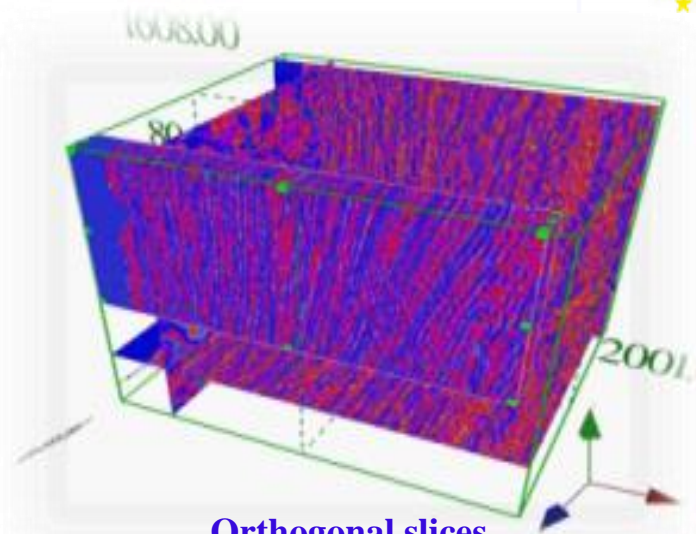
**Gravity Exploration**

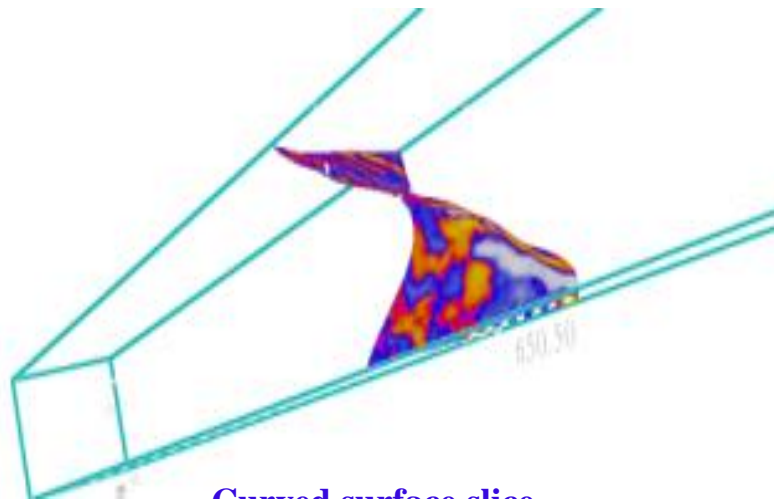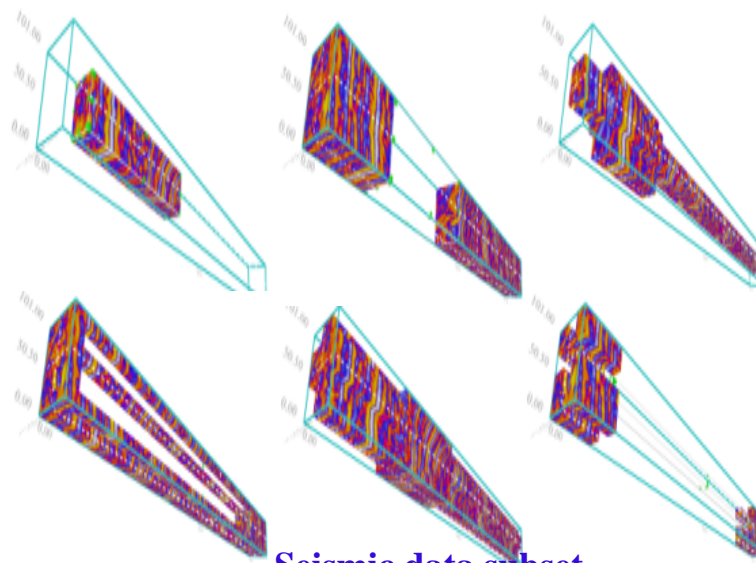**Well logging**

**Seismic Exploration**
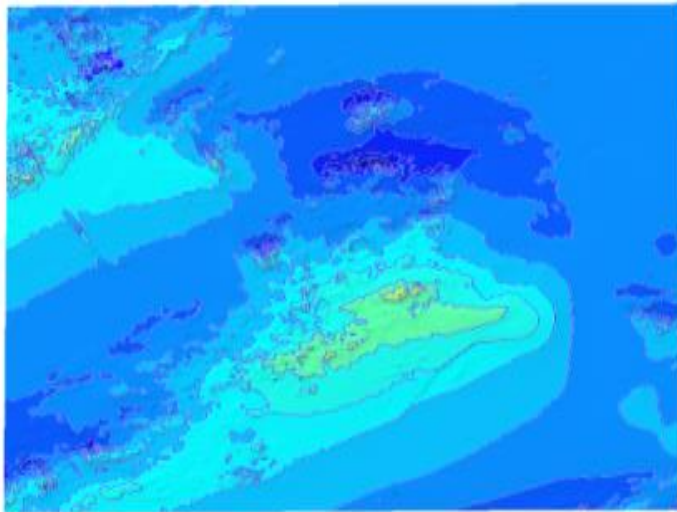
# VolumeViz Module



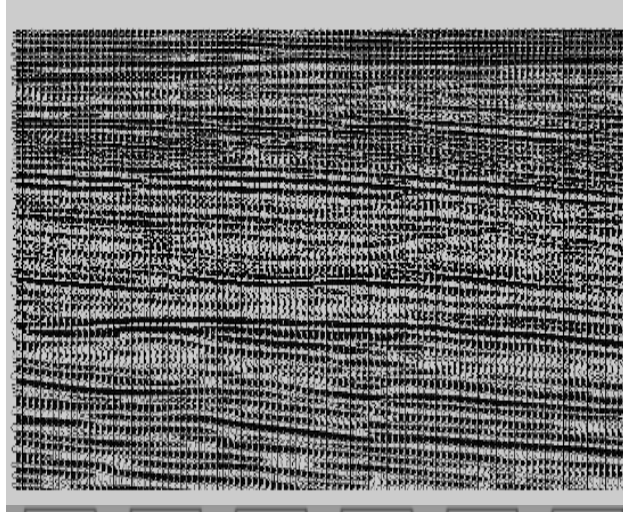3D Seismic data volume



Orthogonal slices
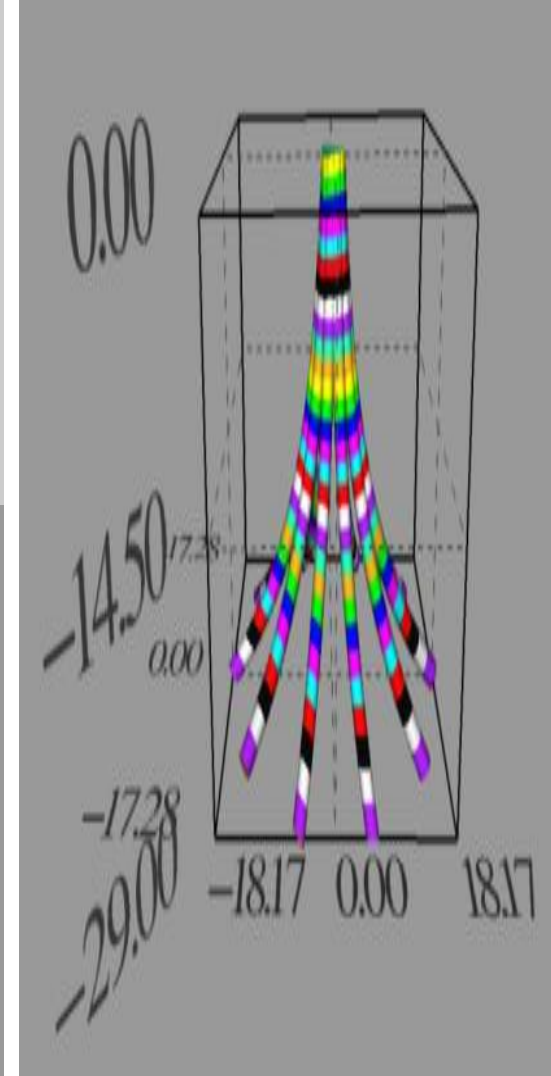


Curved surface slice



Seismic data subset

# MeshViz  Module



2D Grid



3D Grid



2D curve (seismic & logging)



Curve hose (logging tracks)

# LDM: High efficiency management module of large scale seismic data

- Able to handle ultra-large-scale data, the data size unlimited.
- Real-time rendering of massive data.
- Supports multi-resolution seismic data rendering.

Data size：26Gb

Data size：100Gb

- Data mapping algorithm
- Data Fast Load
- tile, surface domain lookup
- Data Logic Coding
- Data storage organization

Computing Cluster

RAM

Storage Cluster

Extended Octree

LDM Core technology architecture diagram

# Outline

- **TH-1A system and its application**

- **Large-Scale data Visualization**

  - ➢ Large-Scale Flow Visualization

  - ➢ Multi source geological data visualization graphics engine——OpenProbe

- **Summary**

# Summary

- Propose novel approaches to integrate flow field analysis with multivariate and ensemble analysis.

- Propose a parallel sparse data management to support advection-based flow analysis.

- Develop a set of common, object oriented 3D graphics toolkit, supporting geoscience application software secondary development.

- OpenProbe offers a range of application modules for geoscience data features, and uses LDM to achieve massive geoscience data 3D stereoscopic dynamic visualization.

# **Thanks**