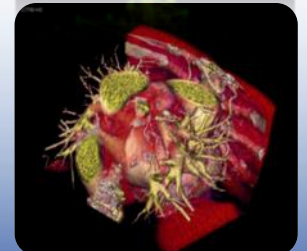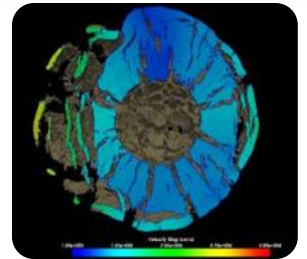# Big Data Analysis and Visualization with Open Source Tools

Julien Jomier, Jeff Baumes and Joachim Pouderoux
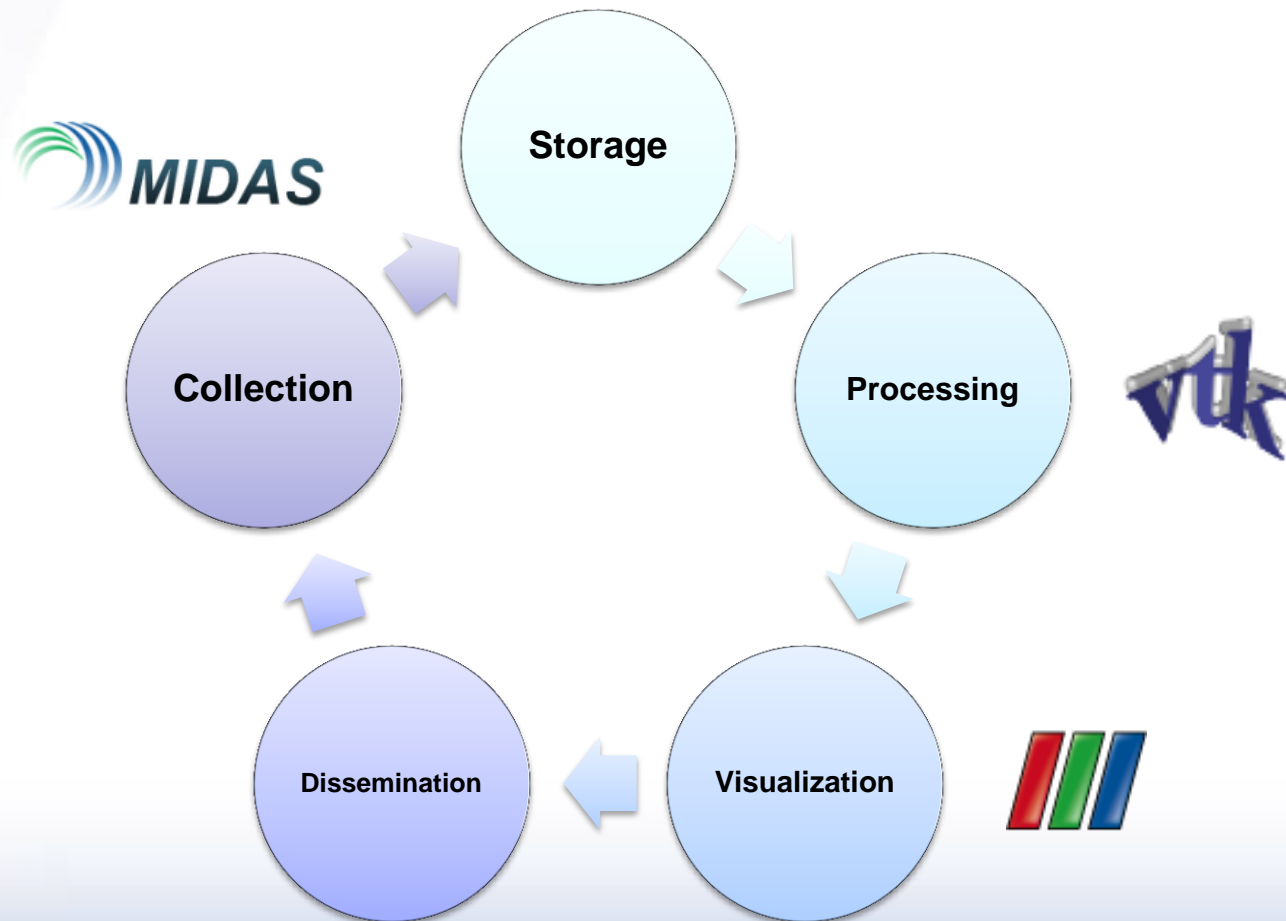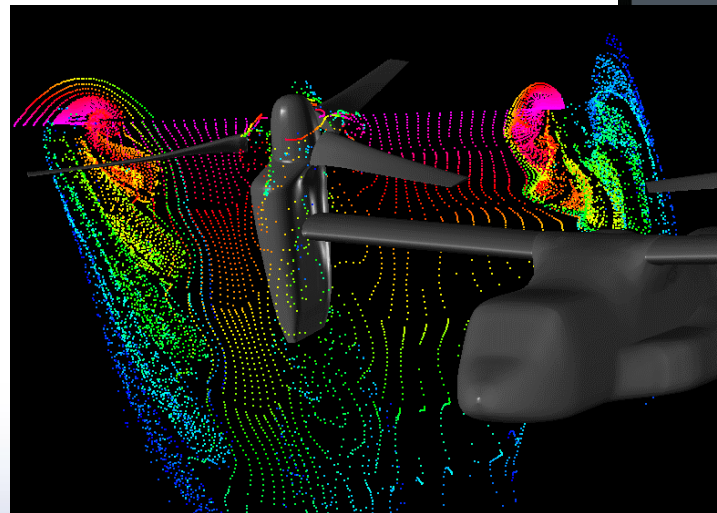julien.jomier@kitware.com
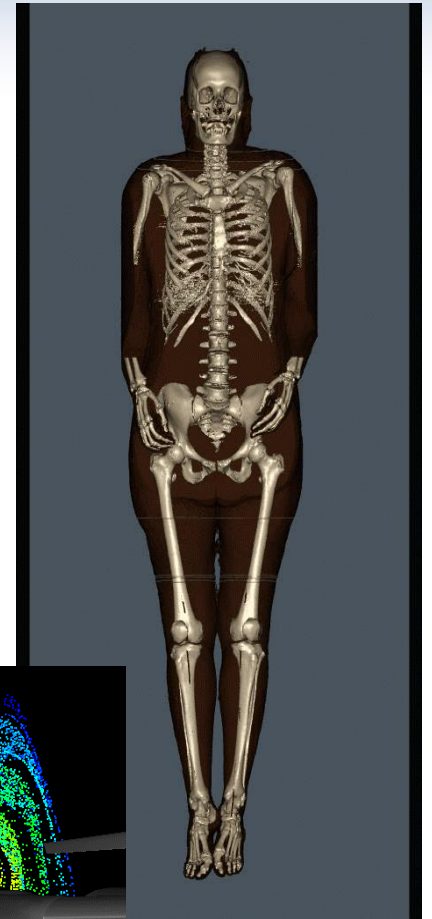
# Big data is everywhere

# Big Data Challenges

- **Parallel** (large data) vs. **Distributed** (many data)

# "Big" is relative

- Visible Human Project (1986)
  - Visible Woman CT Data
  - 870 MBytes 1734 Slices at 512x512x2

- Bell-Boeing V-2 2 tiltrotor
  - 140 GBytes

# Big Data Facts

- **Storage** is **cheap**
  - but fast disk access remains expensive
- **Moving data** across networks is the **bottleneck**
- **Most clusters** are **CPU-based** or **GPGPU** based
  - not necessarily optimal for visualization
- **Necessity for parallel/distributed computing**
  - processor speed is getting better

# Plan

- **Post-processing**
  - VTK/ParaView
- **Co-Processing/In-Situ**
  - Catalyst
- **Informatics**
  - Tangelo, Visomics
- **Computational Chemistry**
  - Avogadro, MongChem
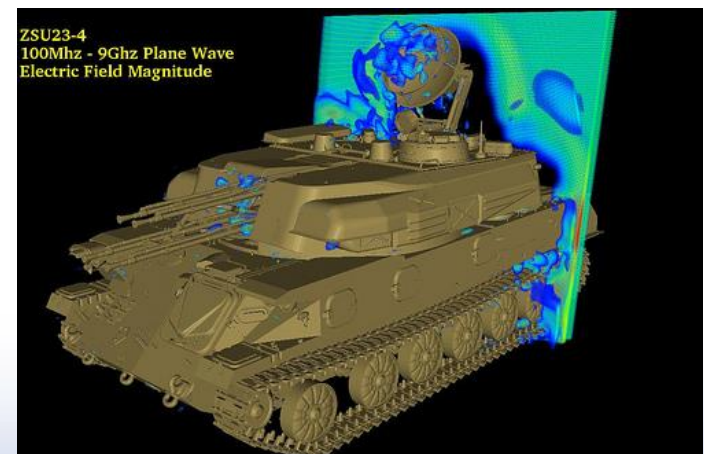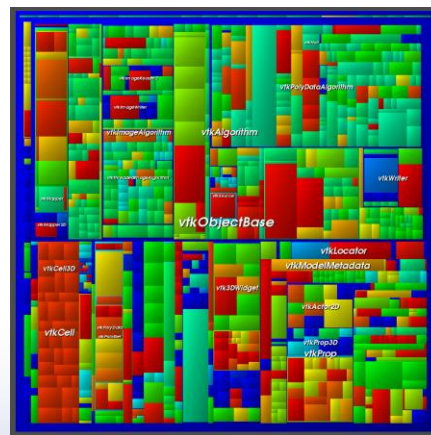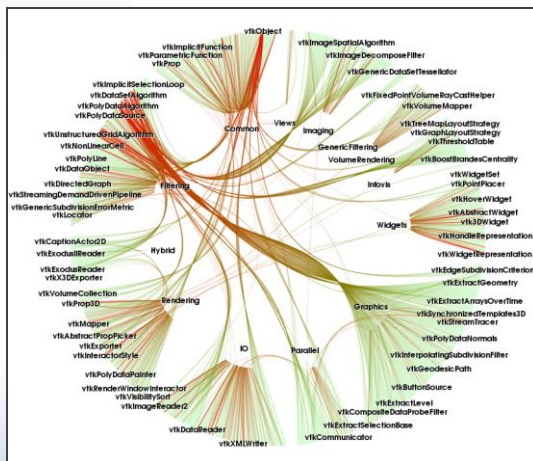- **Web Visualization**
  - vtkWeb/ParaViewWeb

# VTK/PARAVIEW

# The Visualization Toolkit (VTK)
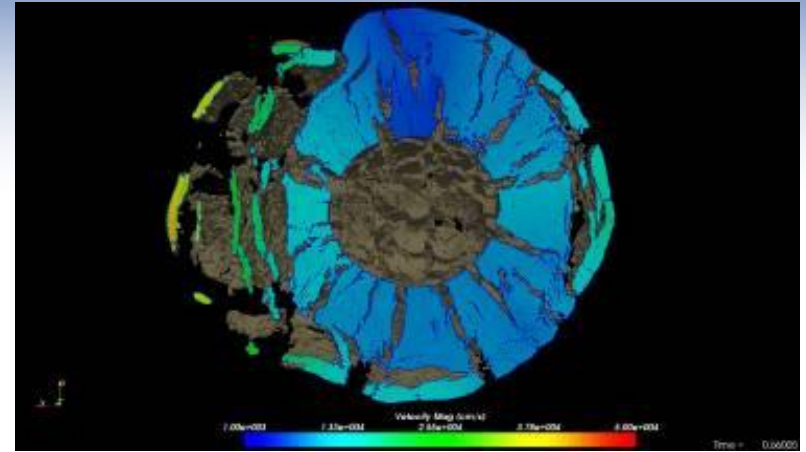
- www.vtk.org
- Started in 1993 at GE
- Visualization Library
  - Written in C++ (+5.5 million LOC)
  - Automatic binding for Java, TCL, Python
  - Portable by design: Linux, Windows, Mac OSX, Solaris…
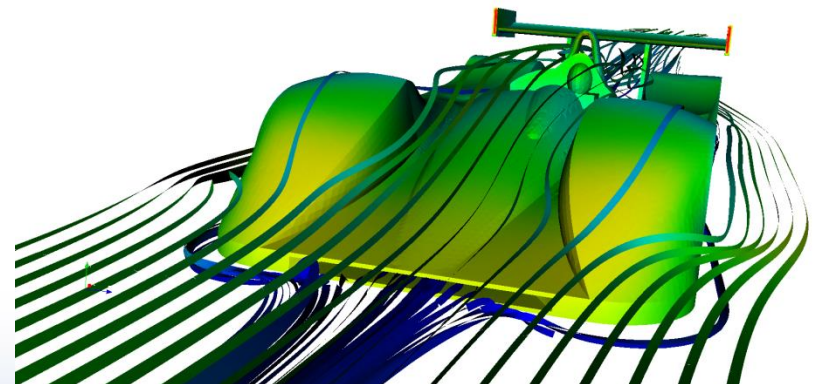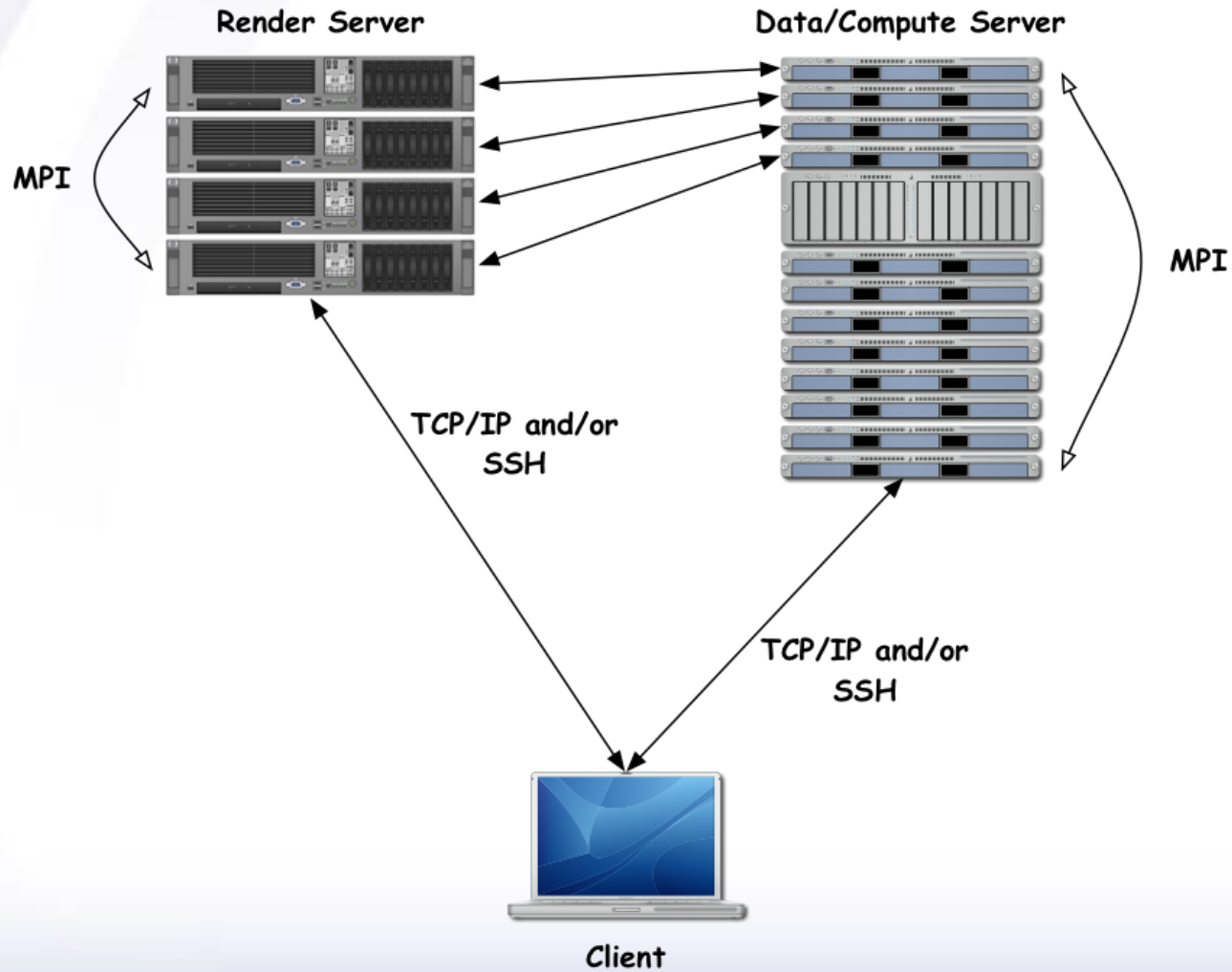- Very active community: 4000+ users

# ParaView



1 billion cell asteroid
detonation simulation

- www.paraview.org
- OpenSource (BSD)
- Based on VTK
- C++/Qt
- Python support
- Very active community (HPC wire award)
- Multi-core support (MPI)
- Co-Processing (in-situ)
- More than 50 news readers
- Visit plugins are supported
- User's guide online
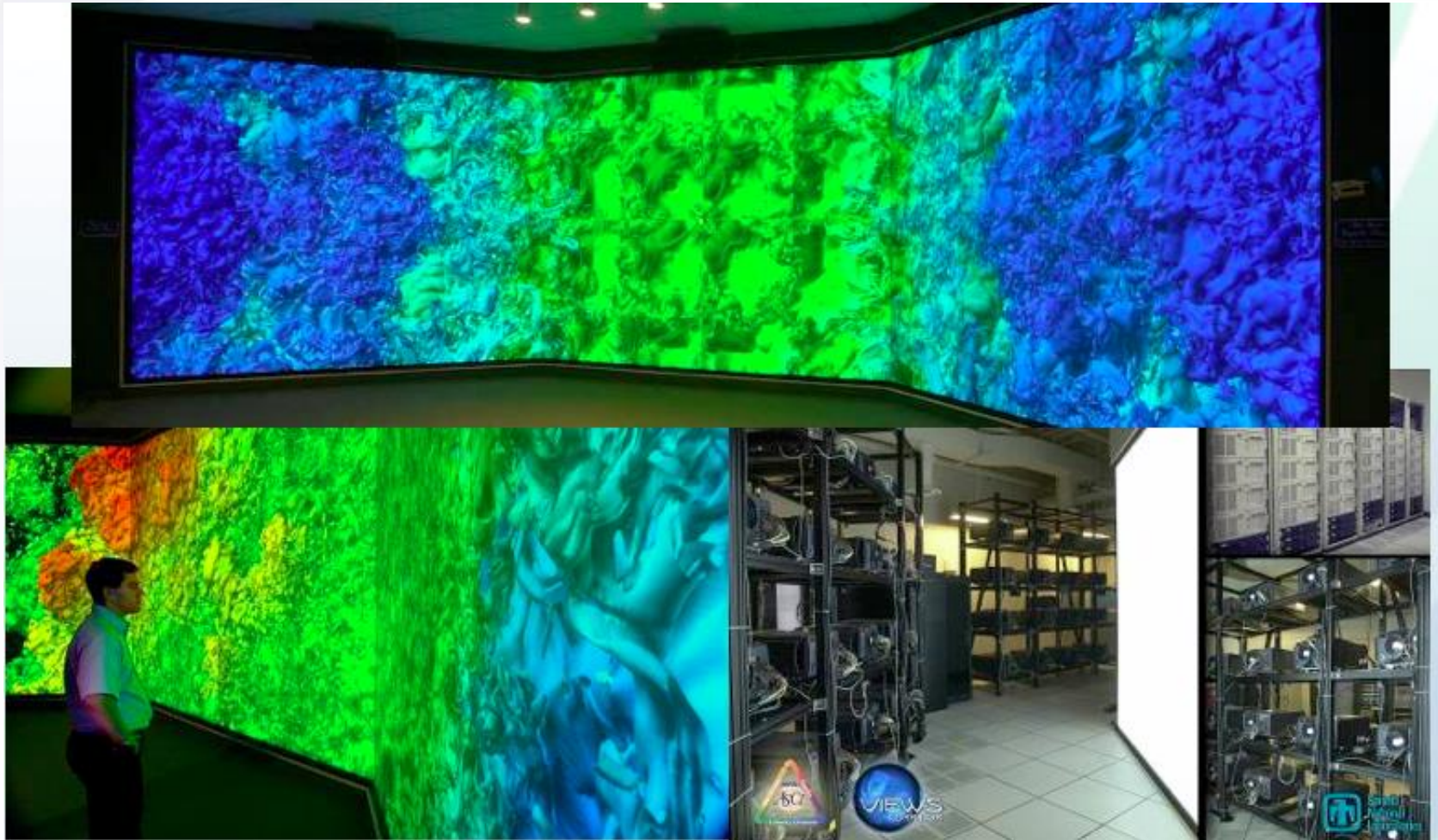
# ParaView Architecture

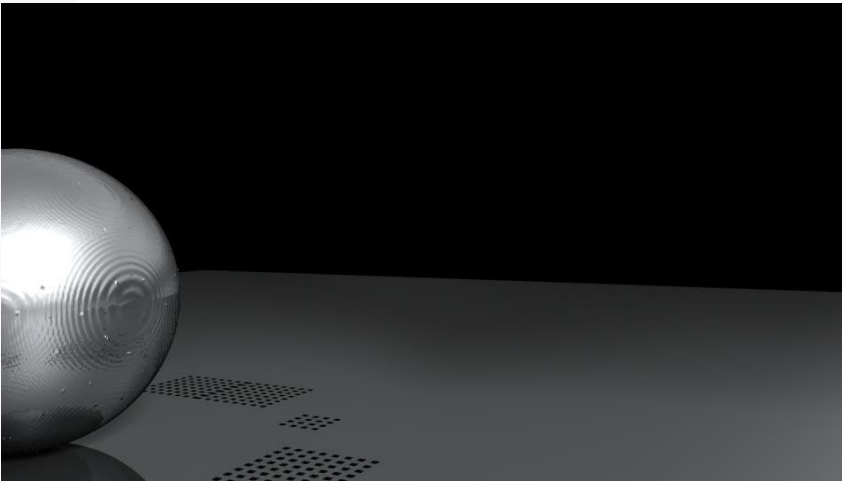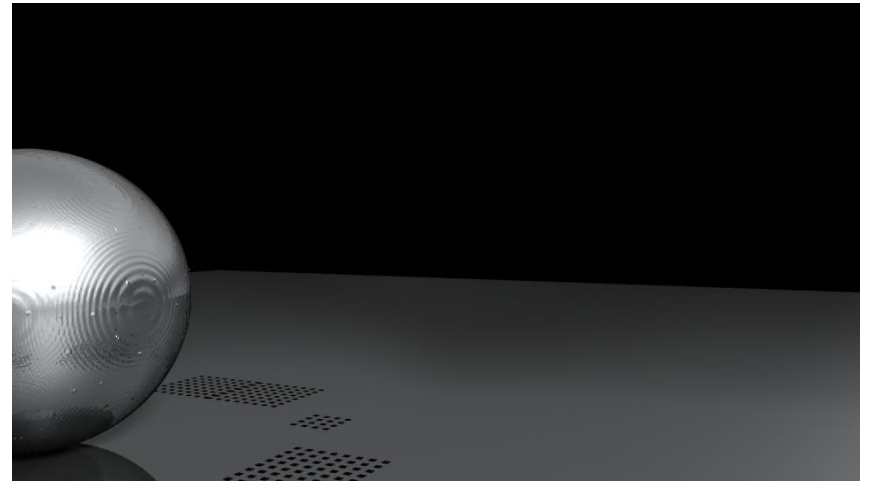# Support for Large Displays

# CATALYST

# In-Situ : Access to More/Richer Data

*Post-Processing*
*(every 100 time steps)*



*In-situ*
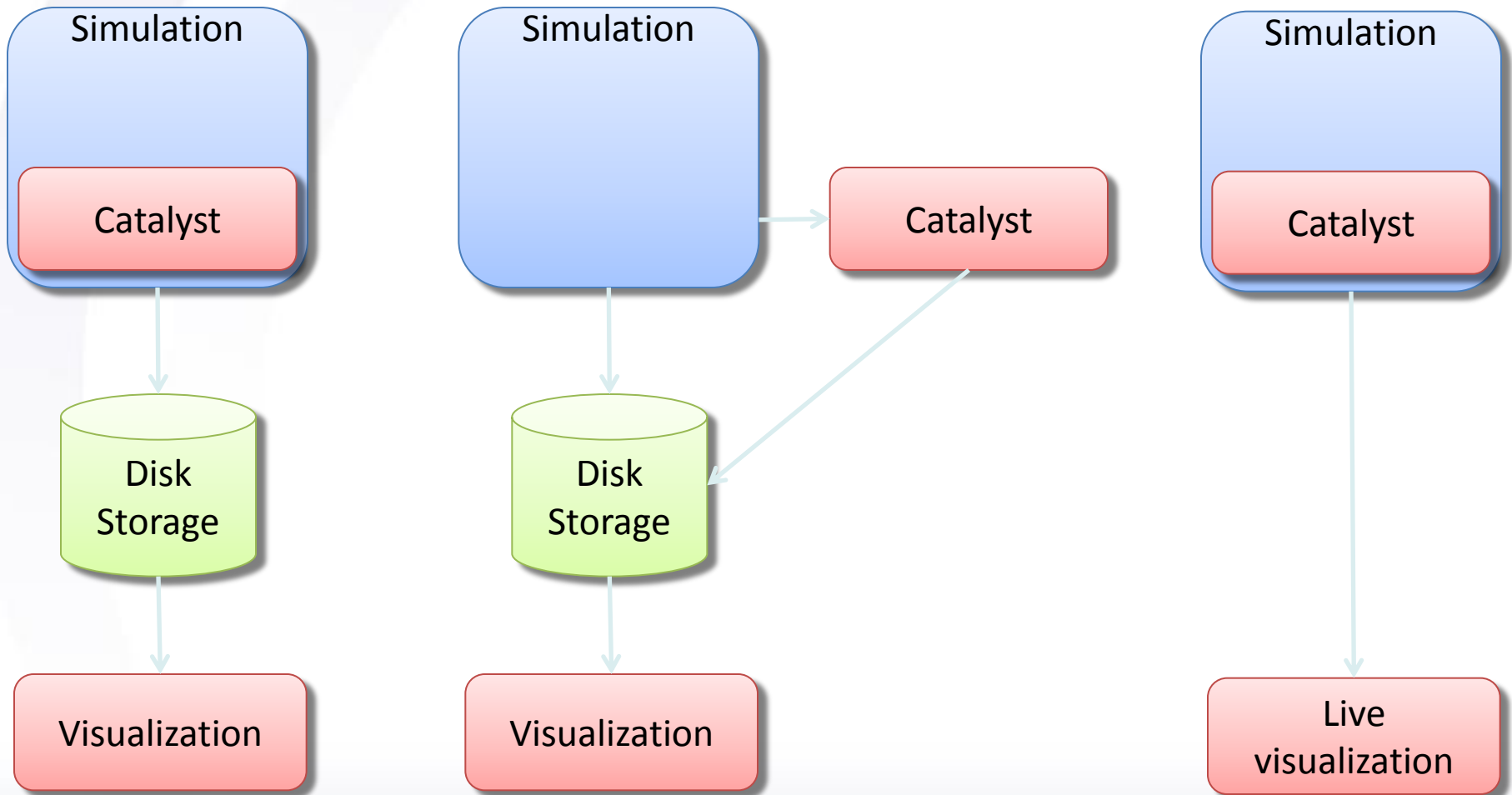*(every time step)*



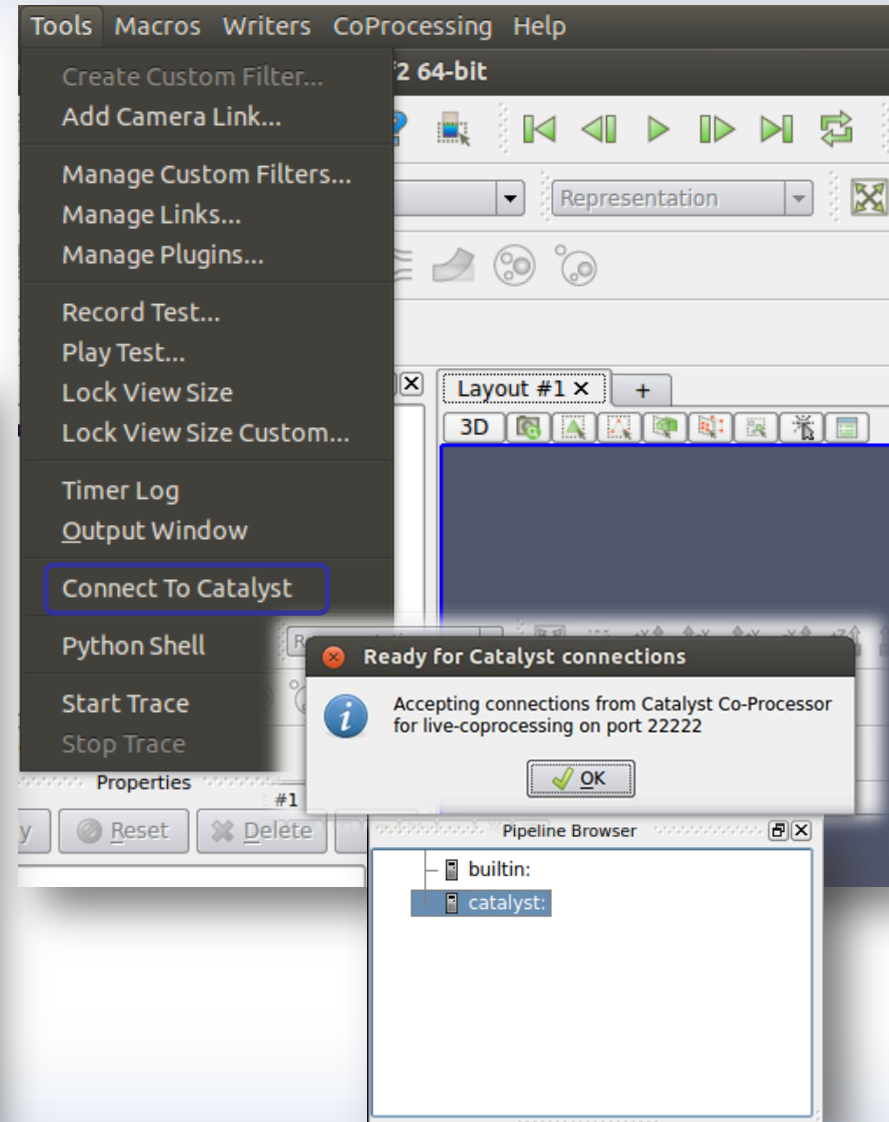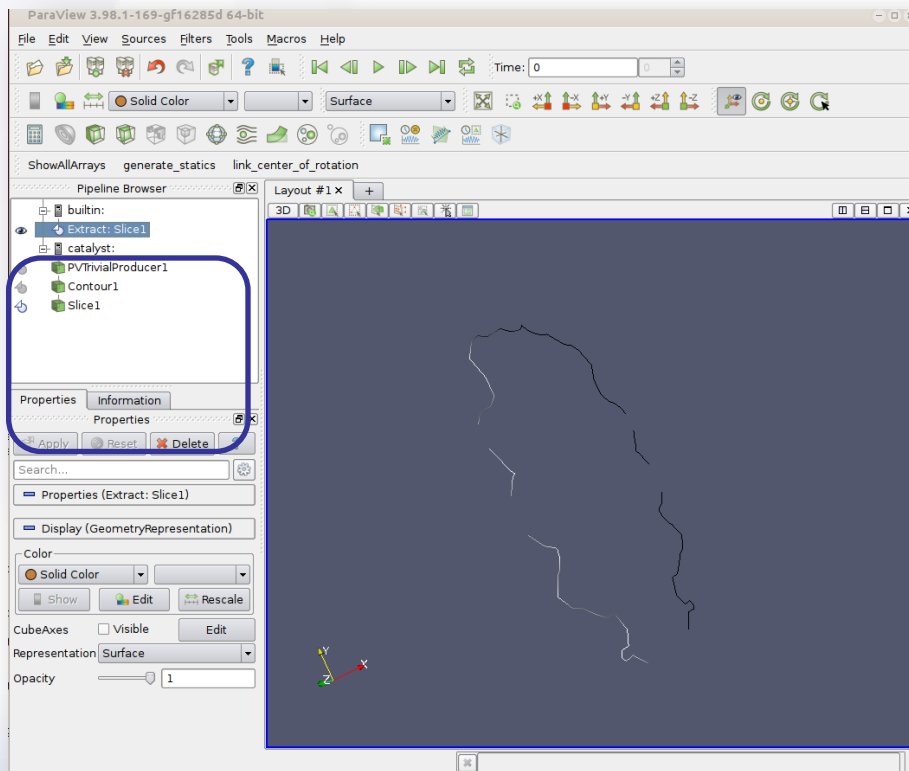*Note: Reflections and shadows added in post-processing*

# *In situ* Analysis & Visualization

# Connecting to Live InSitu Visualization

- ParaView client connects to Catalyst (acting as a server to receive live data)

# Catalyst Architecture



Catalyst

Python Wrappings

ParaView Server
Parallel Abstractions and Controls

VTK
Core Visualization Algorithms

# TANGELO/VISOMICS

# Tangelo

- Python web framework built on CherryPy
- Flexible HTML5 web server architecture
- Developed with a clean separation
  - Application in HTML, JavaScript, CSS
  - Service in pure Python (+ wrapped C/C++)
- Packages several other frameworks too
  - Bootstrap, D3, Vega,MongoDB
- Making web apps easier to develop/deploy
- htttp://tangelo.kitware.com

# Tangelo

- Python for server side, native web clients
- Easily add new services (single .py file)
  - Use RESTful API
  - JSON delivery of data
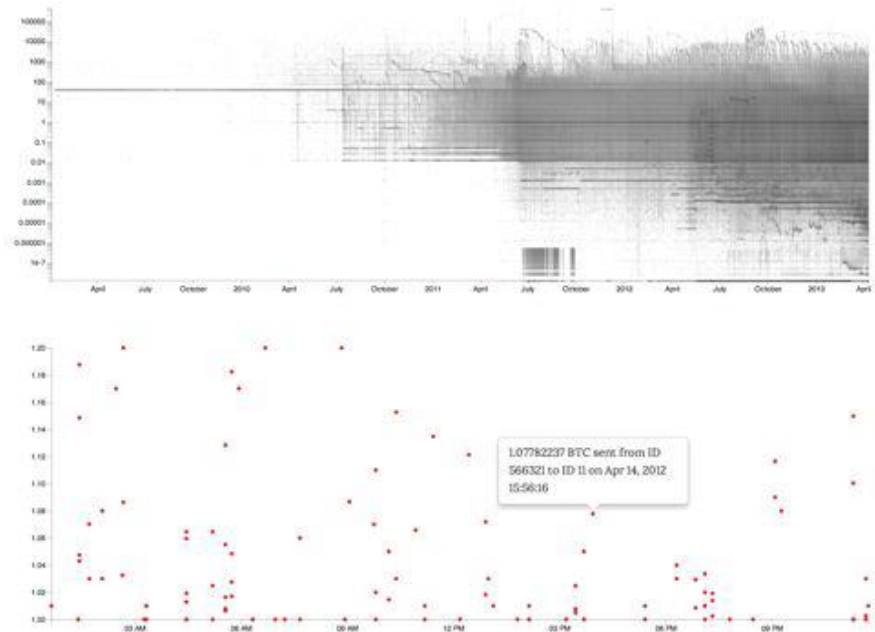  - Full power of Python
- Rapid prototyping

# Bitcoin Analysis

- Uses bitcoin blockchain
  – Individual transactions
- Intensity histogram with transaction volume in date/amount ranges
- Detail plot with individual transactions
- Anomaly search
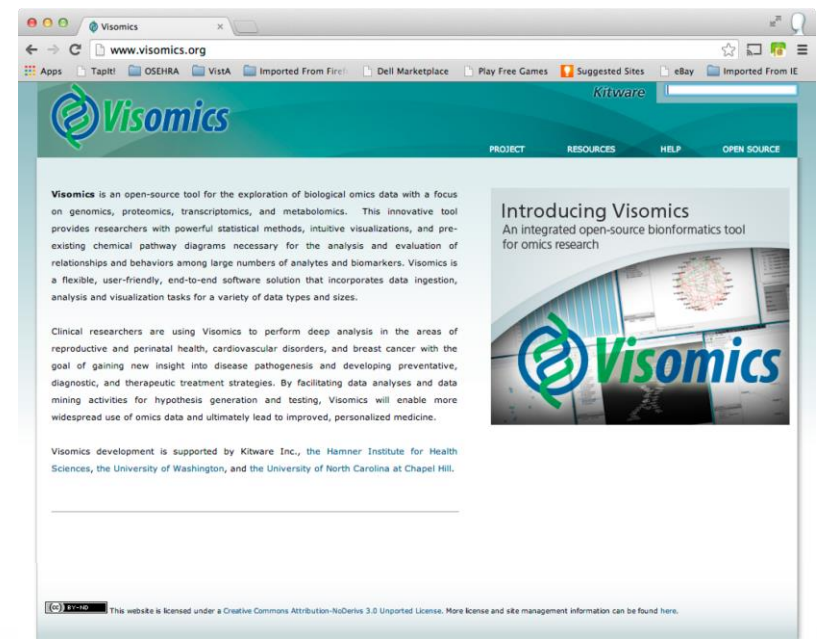  – Theft detection
- Study large scale behavior over time

**Bitcoin Transaction Amounts over Time**

Hover to see the number of transactions in each bin consisting of a date and transaction amount range. Click to generate a detail plot of the individual transactions in that bin. Hover in the detail plot to show transactions with the same target user ID.
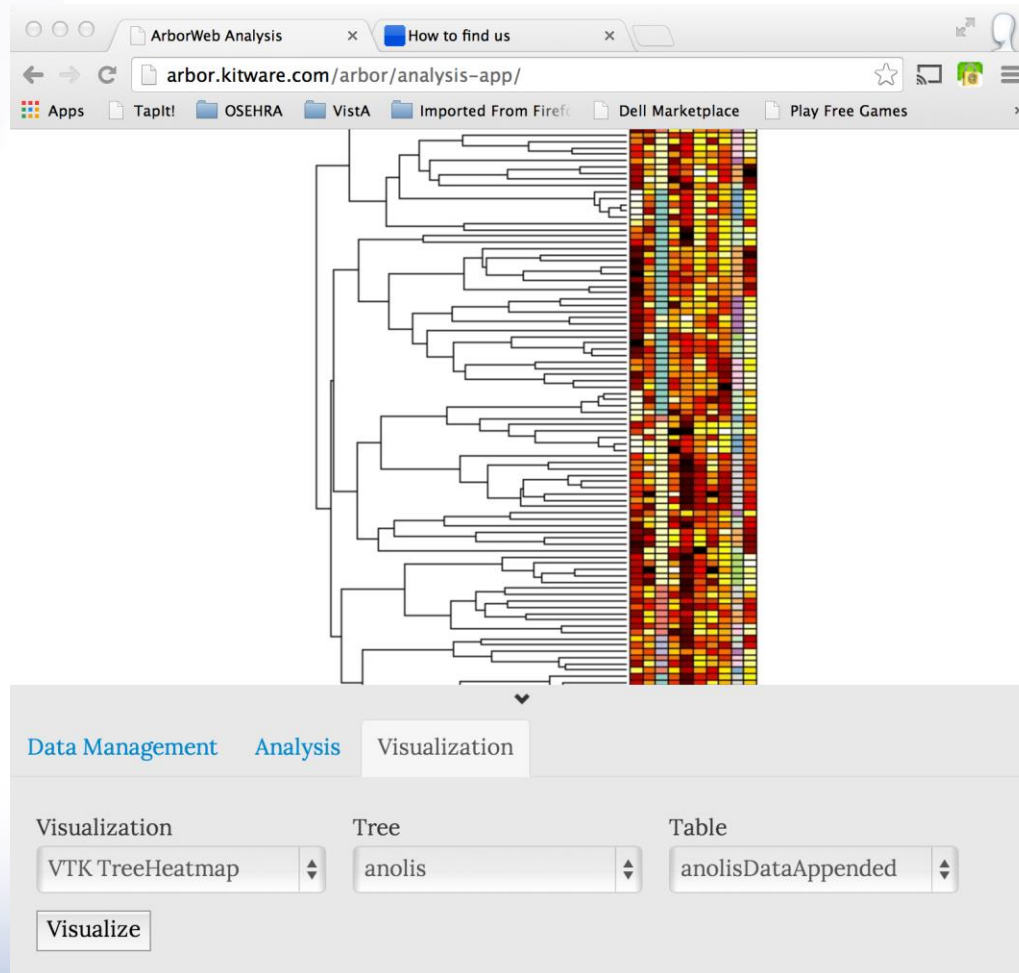
1.07782237 BTC sent from ID 566321 to ID 11 on Apr 14, 2012 15:56:16

# Visomics

- Visomics is an extensible, open source application for research in Omics applications.

- http://www.visomics.org/
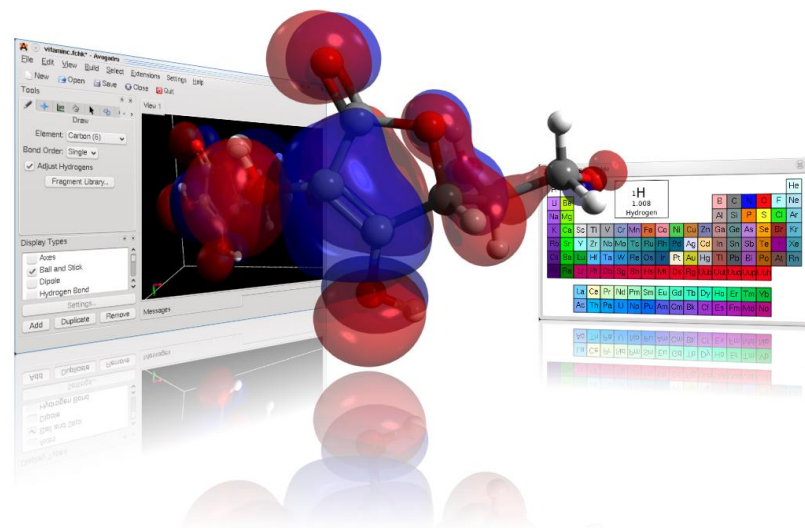
# Arbor

- Brining Visomics to the web

# OPEN CHEMISTRY

# Vision for Open Chemistry

- Advancing the state-of-the-art
- Tight integration is needed
  - Computational codes
  - Clusters/supercomputers
  - Data repositories
  - Reduce, reuse, and recycle!
- Facilitate sharing and searching of data
- Embracing data-centric workflows

# Open Chemistry

The **Open Chemistry** project is a collection of open source, cross platform libraries and applications for the exploration, analysis and generation of chemical data. The project builds upon various efforts by collaborators and innovators in open chemistry such as the Blue Obelisk, Quixote and the associated projects. We aim to improve the state of the art, and facilitate the open exchange of ideas and exchange of chemical data leveraging the best technologies ranging from quantum chemistry codes, molecular dynamics, informatics and visualization.

## Open Chemistry
Explore, analyze and generate chemical data

## News                                                    More News >

**05.29.2013** New Open Access Article on a Collaborative Project Between NWChe...

**08.17.2012** Avogadro Featured in Journal of Cheminformatics

**01.24.2012** Kitware Receives Phase II Funding for the Development of a Comput...

## Events
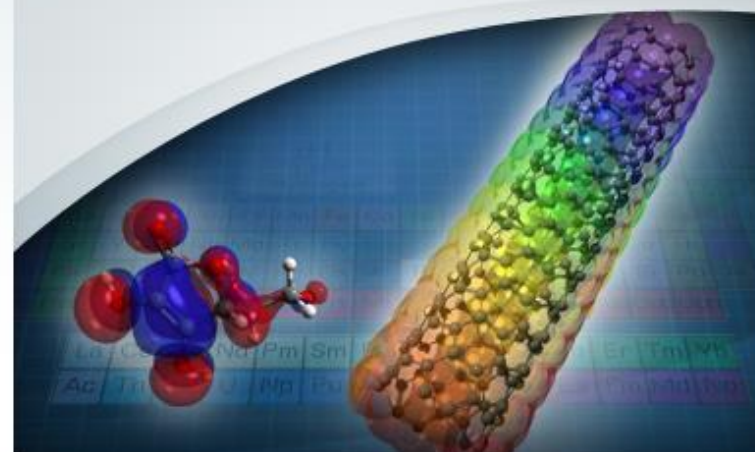
**06.07.2013** LA-SIGMA

## Blog Posts

**05.22.2013** New Input Generator Framework in Avogadro 2

**05.01.2013** Using VTK's Image Regression Tests in Avogadro 2

**04.11.2013** First Open Chemistry Beta Release

# Applications Being Developed

- Three independent applications
- Communication handled with local sockets
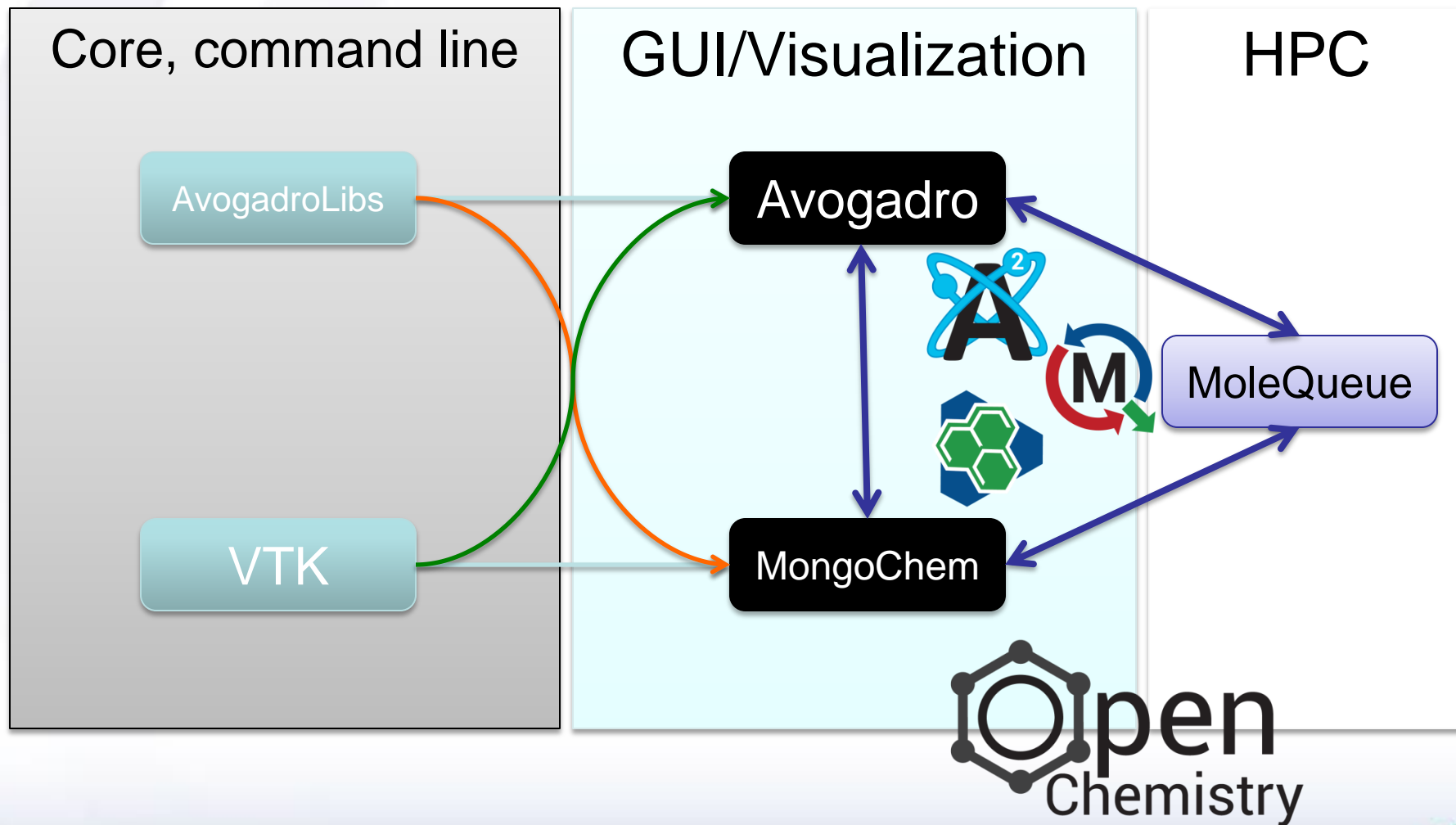- **Avogadro 2:** Structure editing, input generation, output viewing, and analysis
- **MoleQueue:** Running local and remote jobs in standalone programs, and management
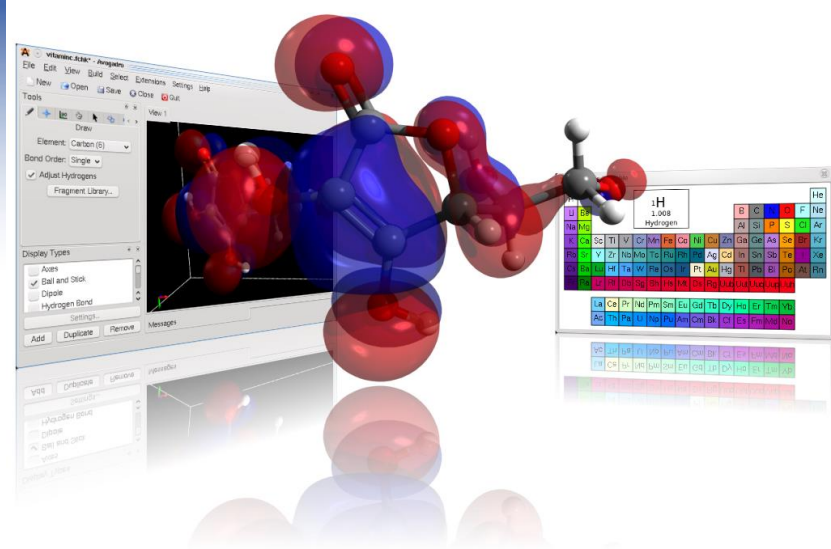- **MongoChem:** Storage of data, searching, entry, and annotation

# Project Diagram: Libraries/Apps

# Avogadro$^2$



- Rewrite of Avogadro
- Split into libraries & application (plugin-based)
- Still one of very few open source **editors**
- Still using Qt, C++, Eigen, OpenGL, and CMake
- Use AvogadroLibs for core data
- Introduces client-server dataflow/patterns
- New, efficient rendering code
- More liberally-licensed: from GPL to BSD

# MoleQueue: Job Management

- Tight integration with remote queues
- Integration with databases
    - Retains full log of computational jobs
    - Triggers actions on completion
- Plugin-based system
    - Easy addition of new codes
    - Easy addition of new queue systems
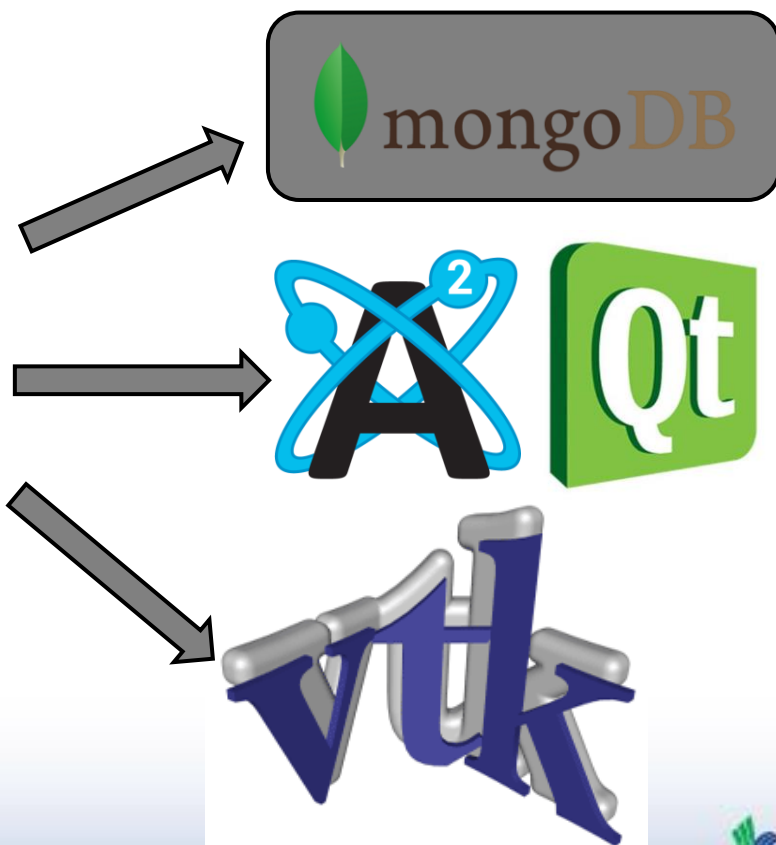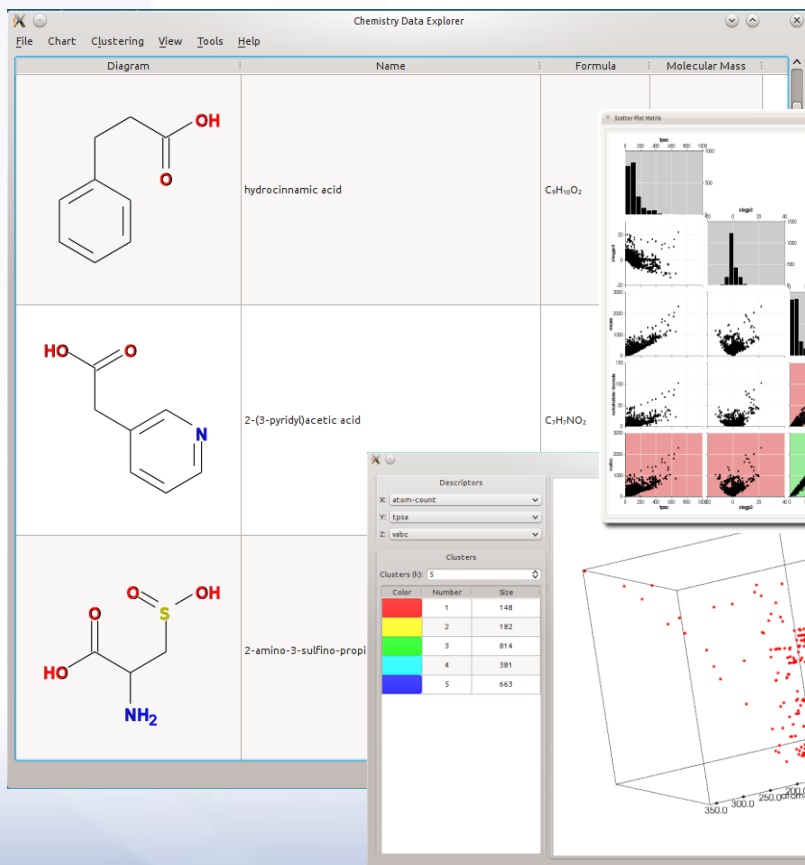- Provides client API for applications

# MongoChem

- A desktop cheminformatics tool
  - Chemical data exploration and analysis
  - Interactive, editable, and searchable database
- Leverages several open-source projects
  - Qt, VTK, MongoDB, Avogadro 2, Open Babel
- Designed to look at many molecules
- Spots patterns, outliers; runs many jobs
- Scales to studies with ~3 million structures

# MongoChem Architecture

- Native, cross-platform C++ application built with Qt and Avogadro 2
- Stores chemical data in a NoSQL MongoDB database
- Uses VTK for 2D and 3D dataset visualization



31

# WEBVIZ

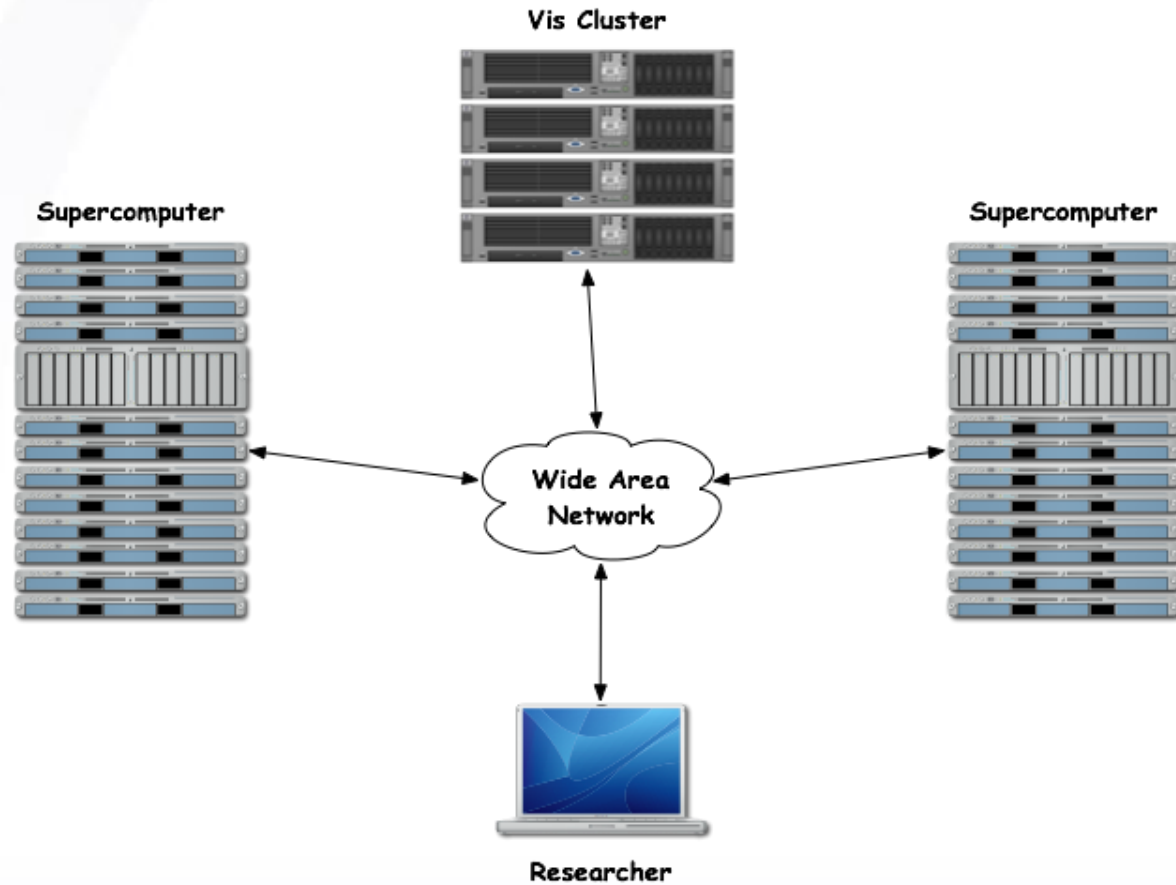# Online Visualization – vtkWeb/ParaViewWeb

- No plugin
- Works on all devices and browsers
- Instant Visualization (fast loading)
- Fully interactive visualization
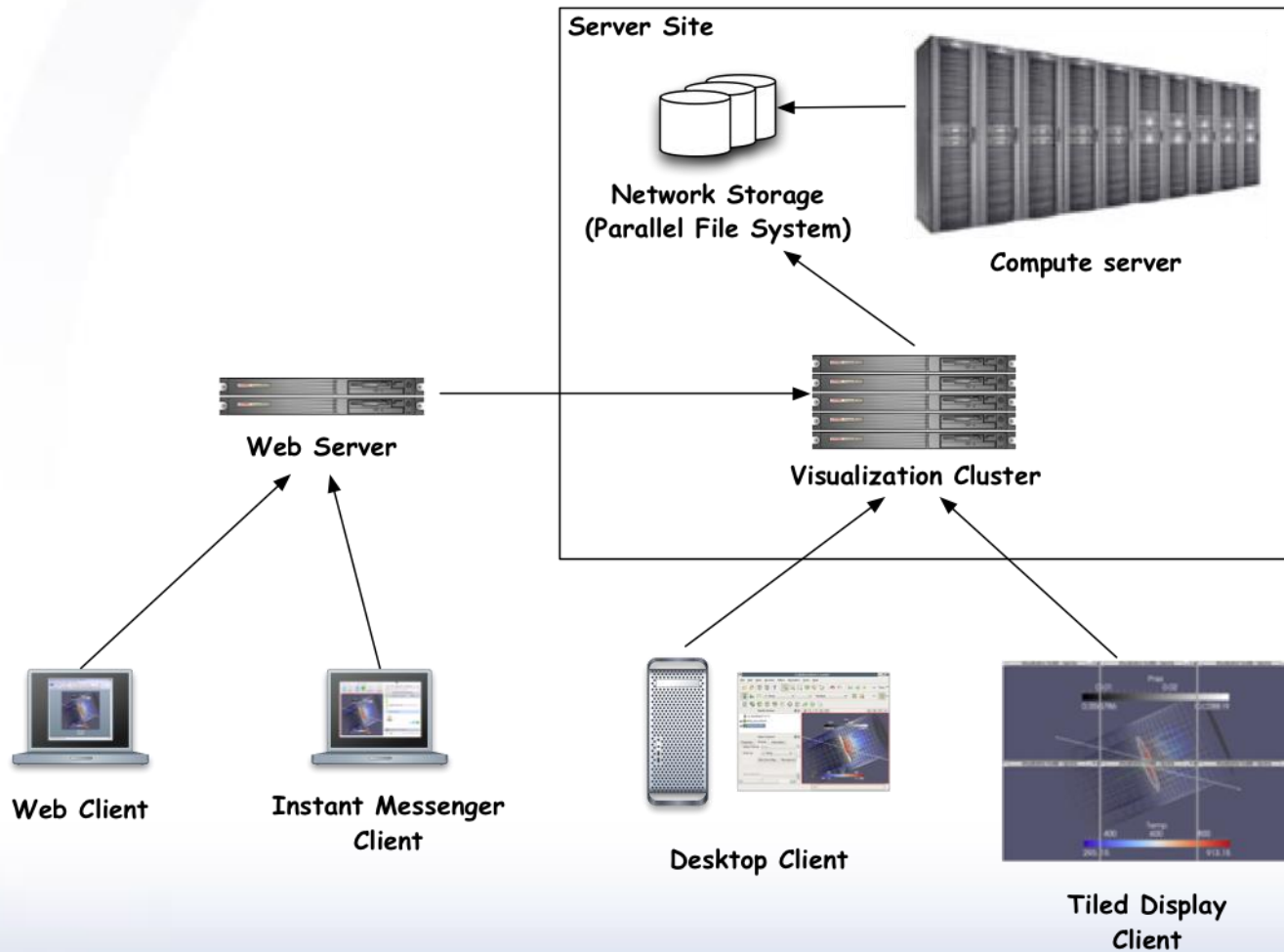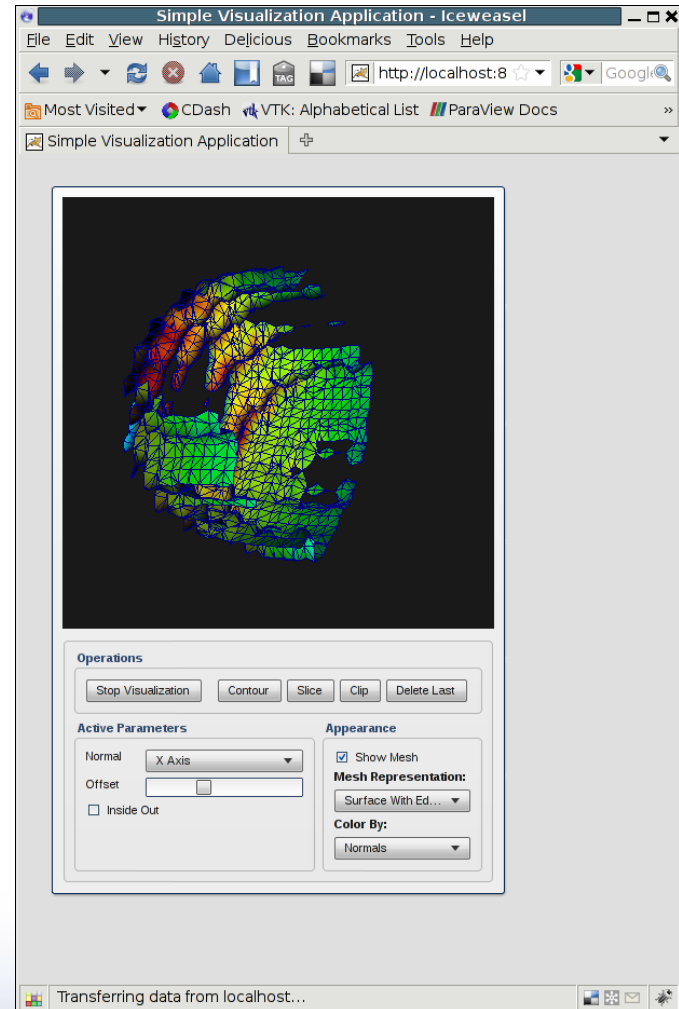- **http://www.webviz.org**
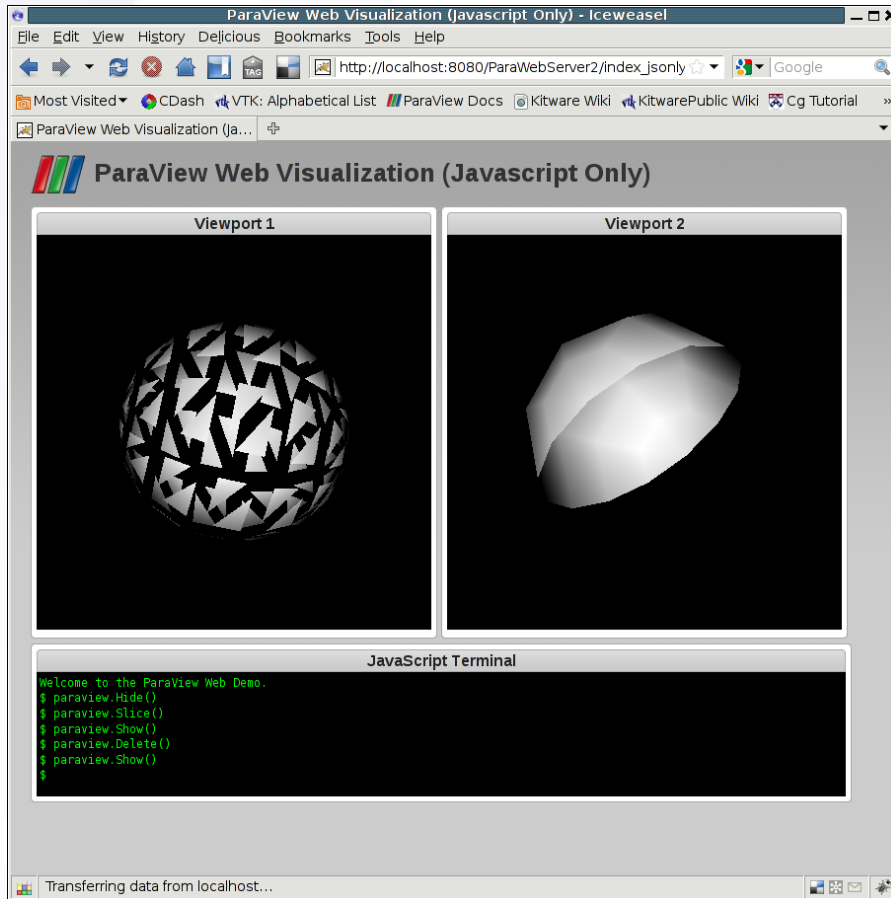
# ParaViewWeb  - Remote Access

# ParaViewWeb - Collaboration

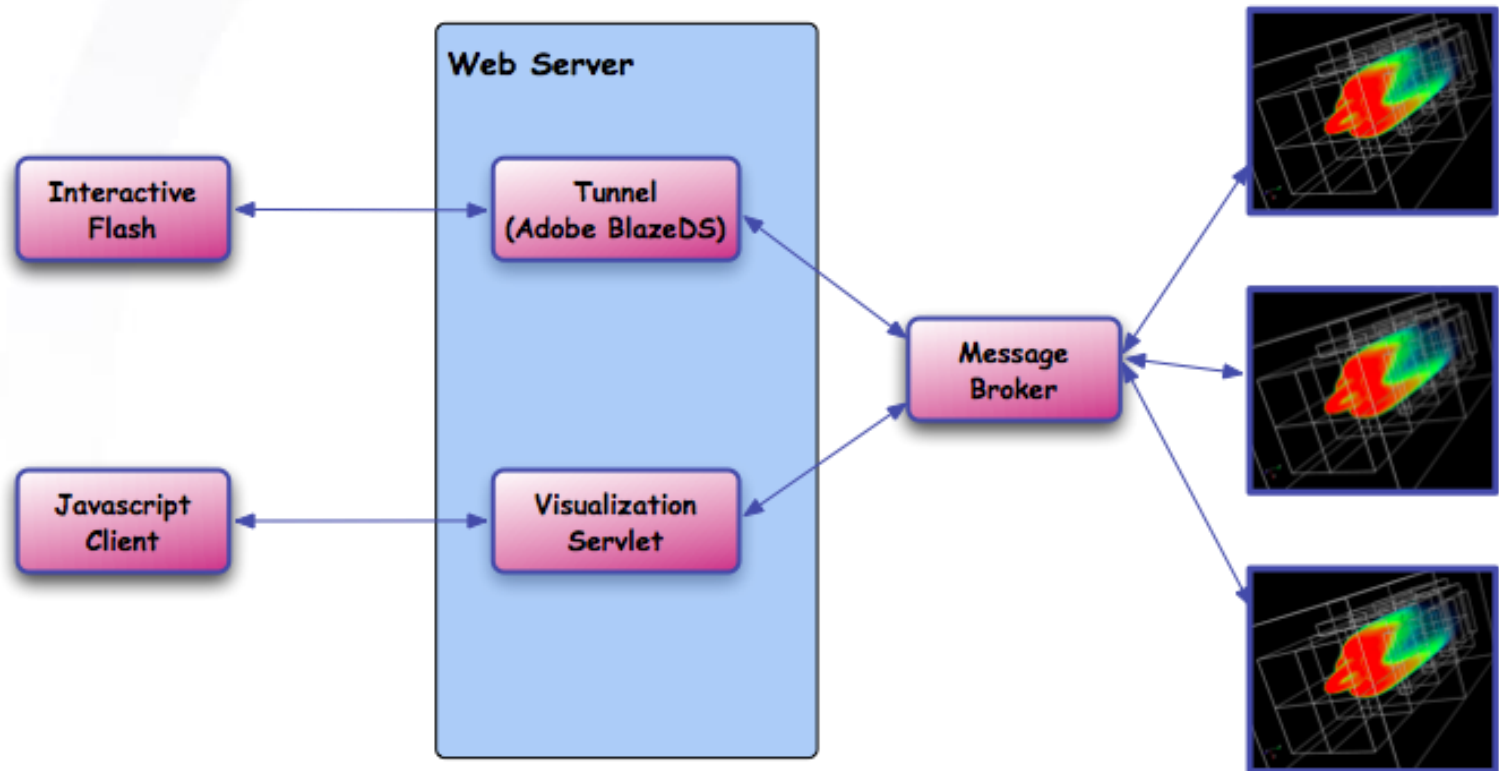# ParaViewWeb - Visualization

# ParaViewWeb - Architecture

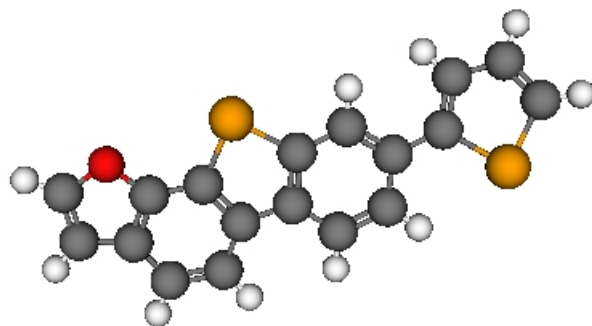# Digital Pathology

- https://slide-atlas.org/

# VTKWeb and MongoChem

- Uses VTK's web architecture

- Performs interactive 3D rendering

- Runs in any modern web browser

- Same MongoDB server as MongoChem

- Moves more to the client JavaScript code

- Using a simple, Python-based server

    - Easy to add new APIs

    - Easy to deploy/integrate into other solutions

# VTKWeb and Open Chemistry

# VES

- www.kiwiviewer.org

- Mobile visualization toolkit based on VTK

- Native applications

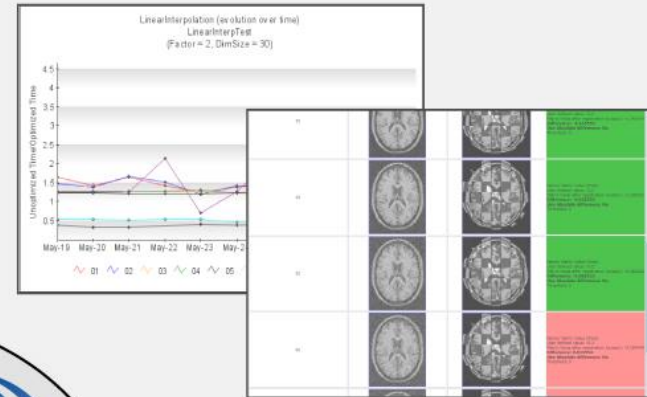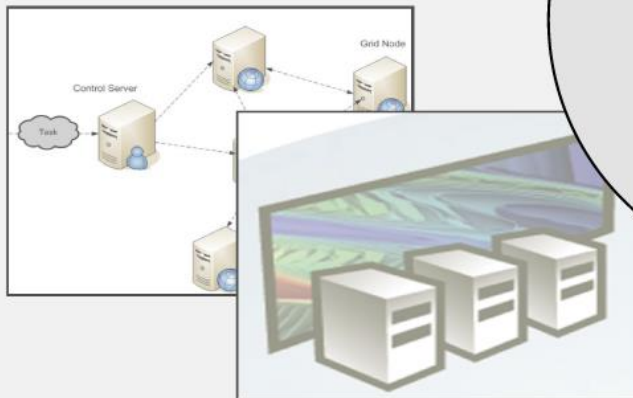- Open Source

- Works on iOS and Android

- High performance
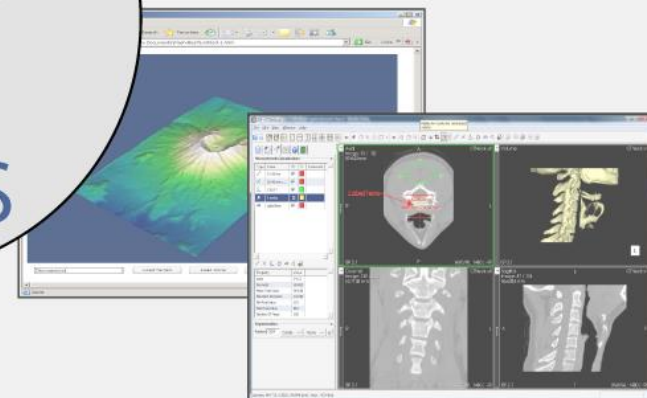
# Complex Data Management



Digital Storage

Online Reporting

Server-Side Processing

Interactive Visualization

MIDAS

# Where are we heading?

- **Computing infrastructure** is available
  - Cloud computer
  - Large clusters
- **Data management**
  - Interoperability
- **Pre-processing tools**
- Coupling **computing** with **visualization**
  - In-situ
- **Visualization nodes** on large cluster is still an issue
- **Easier/Better interfaces** for the users

# **Thank You!**

Julien Jomier, Jeff Baumes and Joachim Pouderoux
julien.jomier@kitware.com