# Berkeley Lab Visualization Program
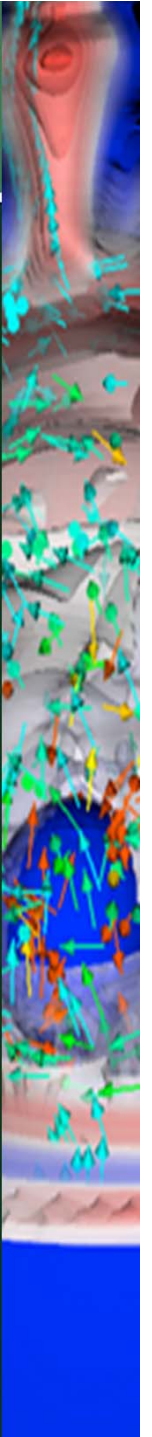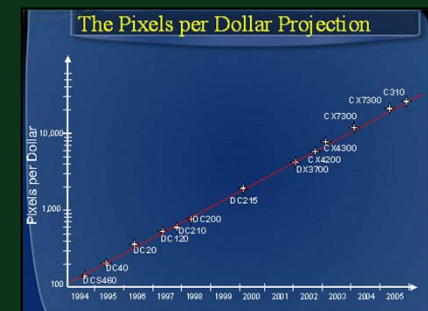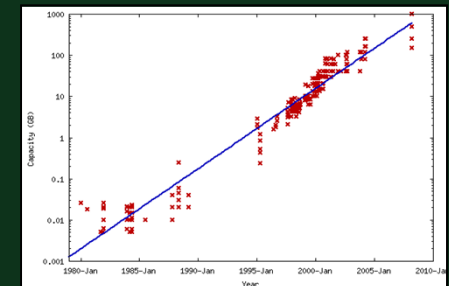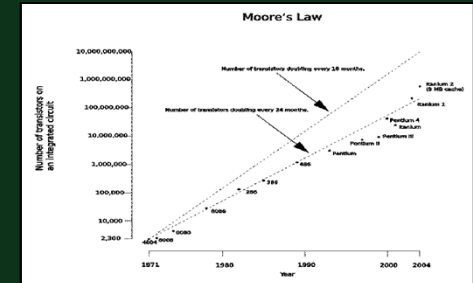
*E. Wes Bethel*

*Dec 15, 2010*

*Berkeley CA*

# Berkeley Lab Visualization Group Mission

- Enable scientific knowledge discovery through the research, development, deployment, and application of visual data analysis technologies in the modern regime of HPC and data intensive science.

- We accomplish this mission by:

  - Focusing R, D, & D efforts at all stages of the visualization pipeline.

  - Close collaborations with science stakeholders to maximize likelihood of science impact.

  - Tightly integrated and well coordinated interaction between research, development, and production deployment activities.

# A Big Problem: Too Much Data

- Our ability to create and store information exceeds our capacity to understand it.
  - Moore's law growth in transistor density, hard drive storage capacity, instrument resolution.
- Human cognitive capacity: flat, constant.
  - Shrinking? Information requires attention to process:
    - "A wealth of information creates a poverty of attention." – Hebert Simon, Nobel Prize, 1971.
- Major challenge: gain insight from data.
  - Visualization, visual data analysis are excellent tools for accomplishing this objective.
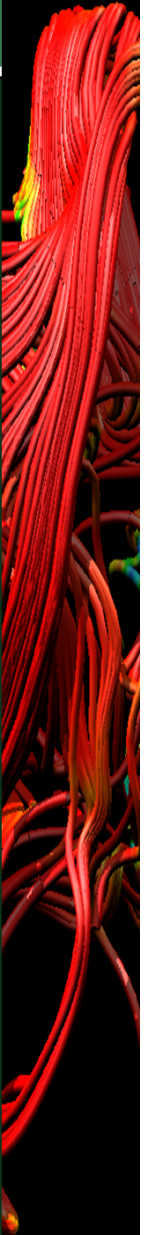
# Today's Vignettes

- **Query-driven visualization and high-energy physics accelerator design.**
  - *Accelerate scientific knowledge discovery.*

- **Topological analysis of combustion simulation results.**
  - *New algorithms enable new scientific insights.*

- **Hybrid-parallelism and extreme-scale visualization.**
  - *Advanced research to fully utilize tomorrow's architectures.*

- **Production-quality, petascale capable visualization software.**
  - *Delivering high-quality capabilities to the scientific community today.*
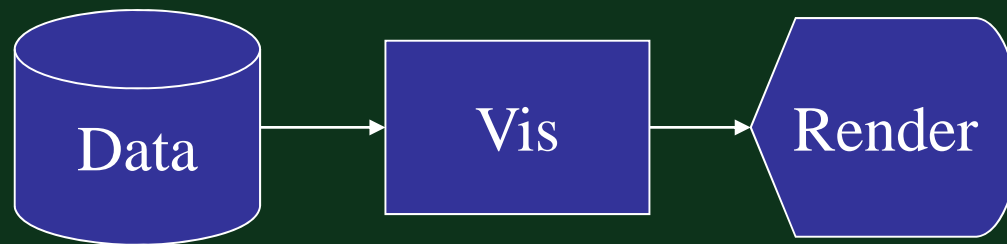
# Query-Driven Visualization

- What is Query-Driven Visualization?
  - Find "interesting data" and limit visualization, analysis, machine and cognitive processing to that subset.
- One way to define "interesting" is with compound boolean range queries.
  - E.g., $(CH_4 > 0.1)$ AND $(T_1 < temp < T_2)$
- Quickly locate those data that are "interesting."
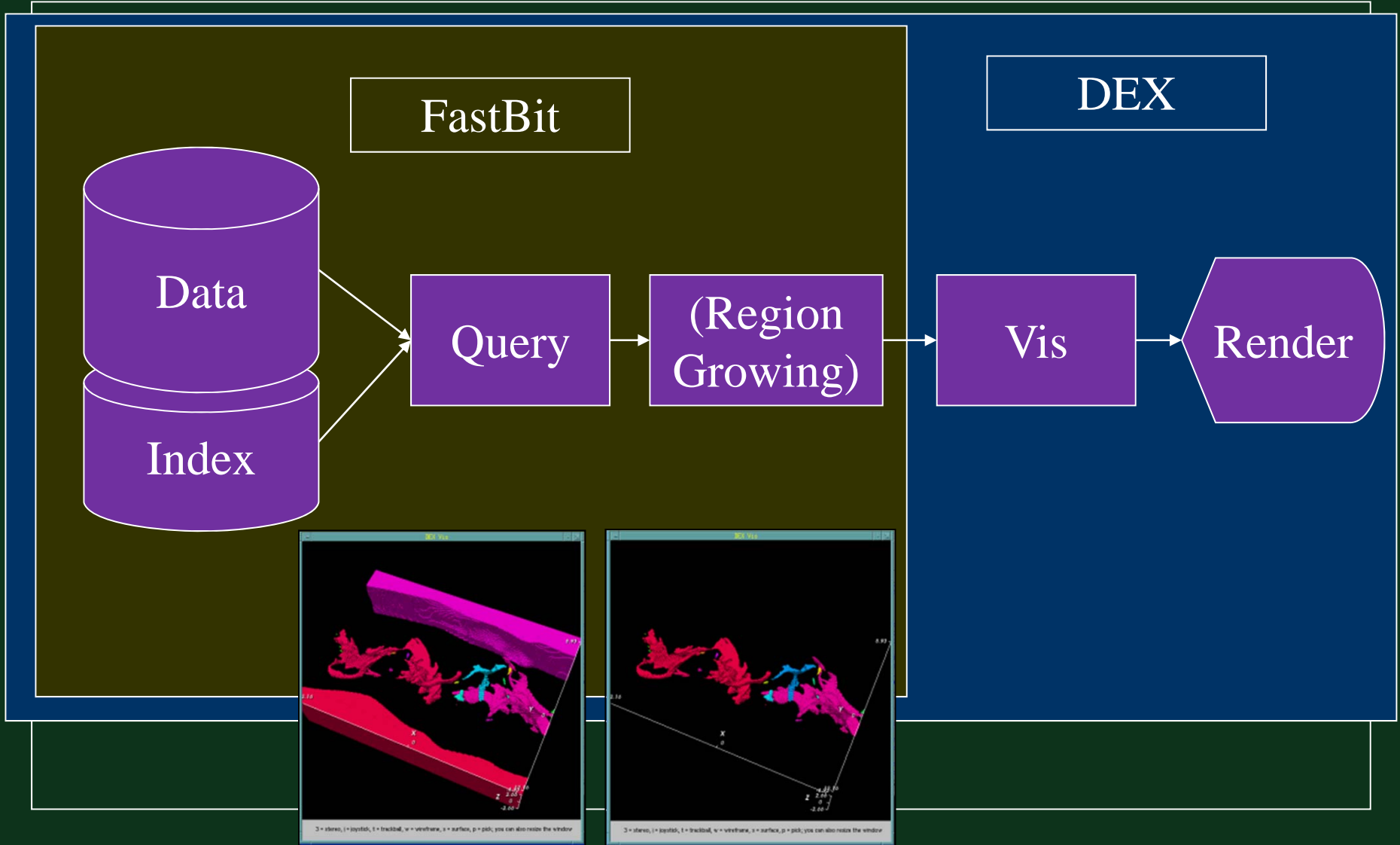- Pass results along to visualization and analysis pipeline.
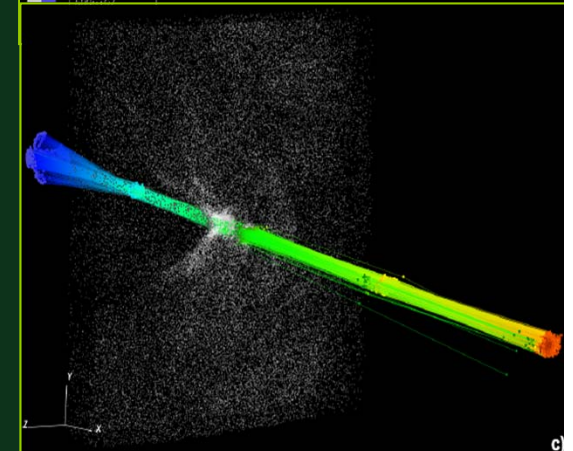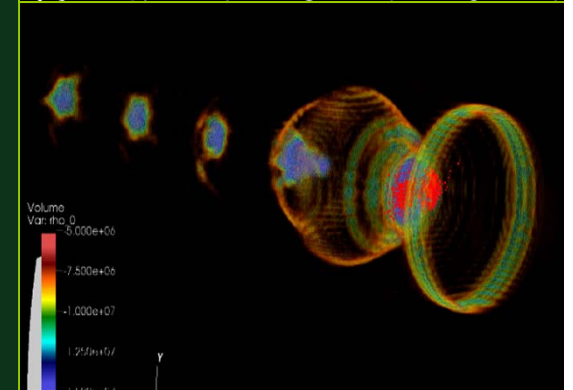
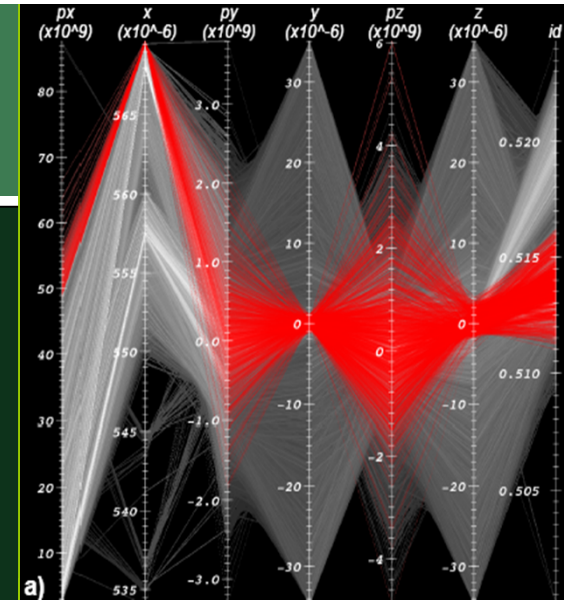# Query-Driven Visualization



The Canonical Visualization Pipeline
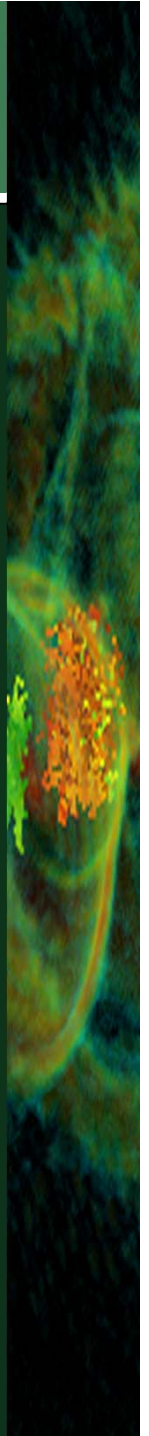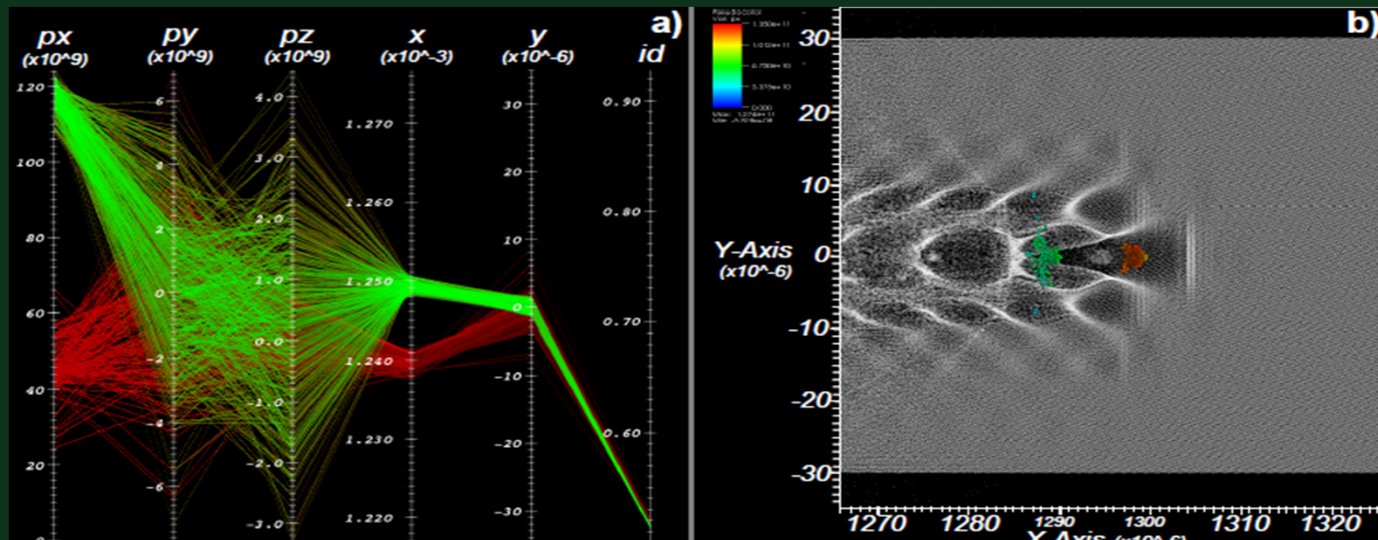
# Query-Driven Visualization

# QDV and Accelerator Modeling



- Problem: sheer size and complexity of data is a barrier to analysis. How to make the problem more tractable?

- Accomplishment:
  - Algorithms and production-quality s/w infrastructure to perform interactive visual data analysis (identify, track, analyze beam particles) in multi-TB simulation data.

- Science Impact:
  - Replace serial process that took hours with one that takes seconds.
  - New capability: rapid data exploration and analysis.

- Collaborators:
  - PI: C. Geddes (LBNL), part of SciDAC COMPASS project, Incite awardee.
  - SciDAC SDM Center (FastBit)
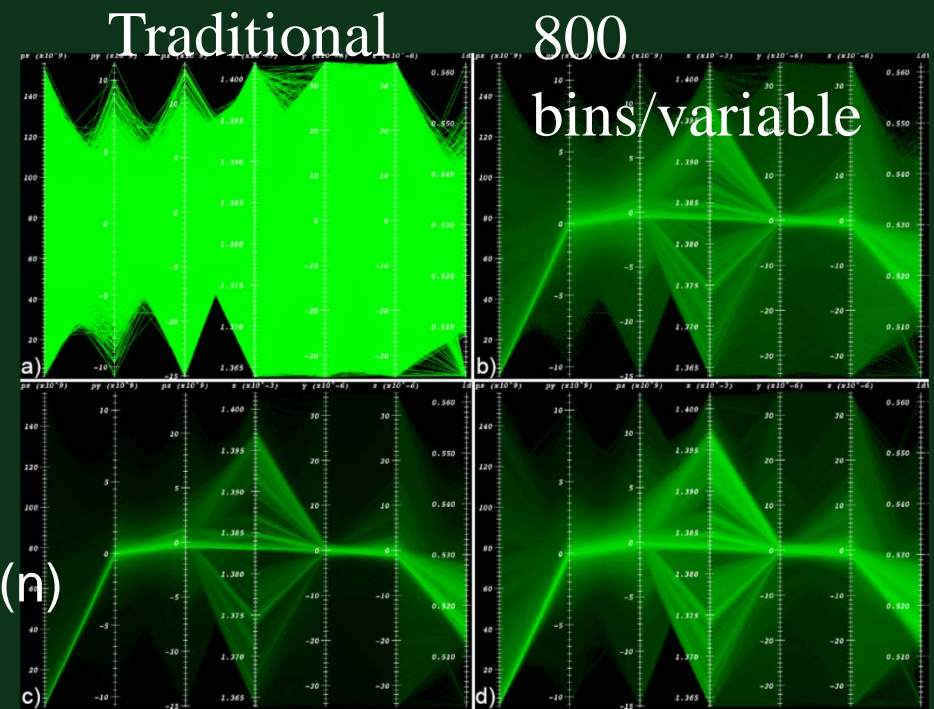  - Tech-X (Accelerator scientists)

# Analysis Task(s)

1.  ## Identify particles that form a beam
    *   Interactive visual data exploration
    *   Data subsetting: high energy, spatial coherency.
2.  ## Track them over time
    *   Given particle ID's from a given time step,
    *   Find all those particles in all time steps
    *   Subsequent visual data analysis.
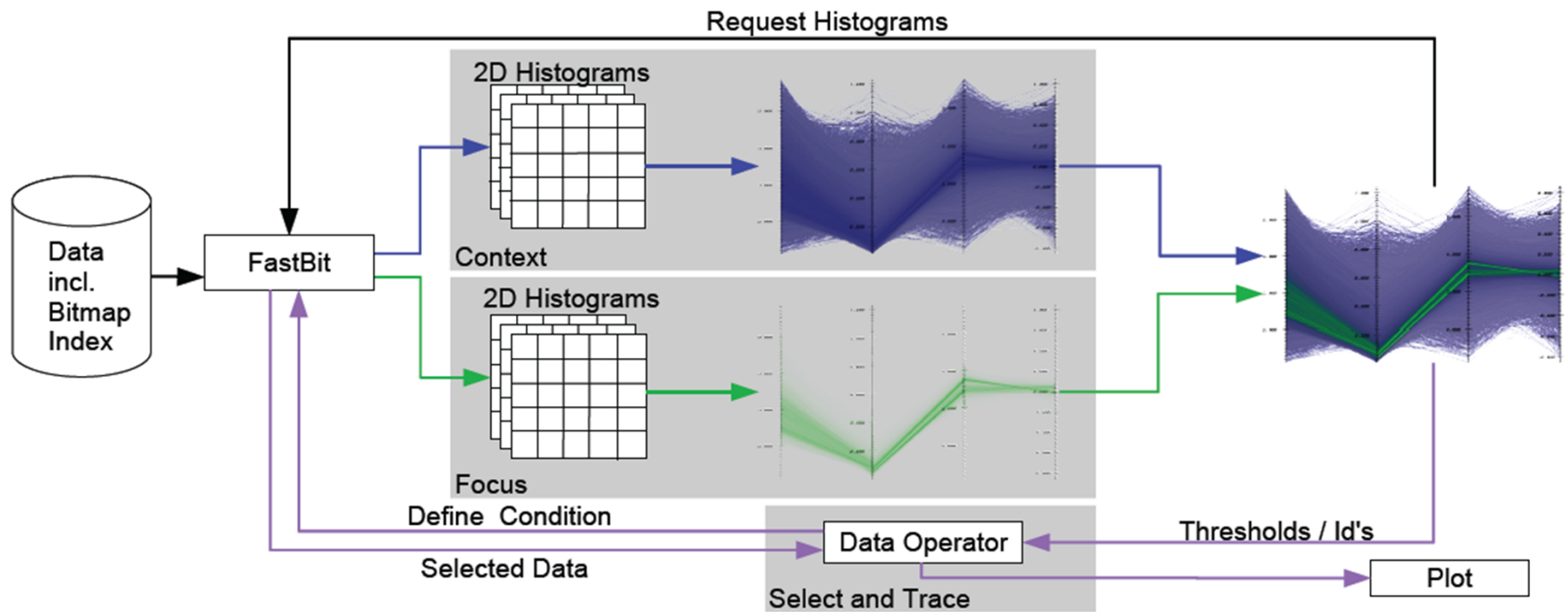
# Fundamental Problem #1 - Interface

- Solving this large-data problem required visualization R&D:
- Parallel coordinates
  - An interface for subset selection.
  - A mechanism for displaying multivariate data.
- Problems with large data:
  - Visual clutter
  - O(n) complexity
- Solution/Approach
  - Histogram-based p-coords
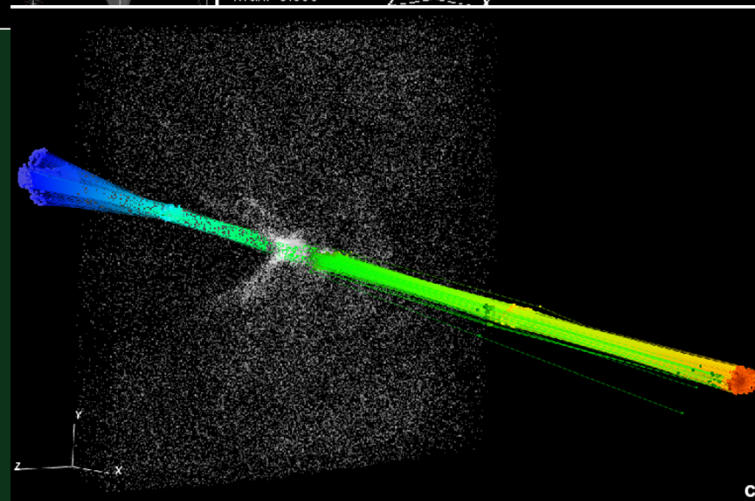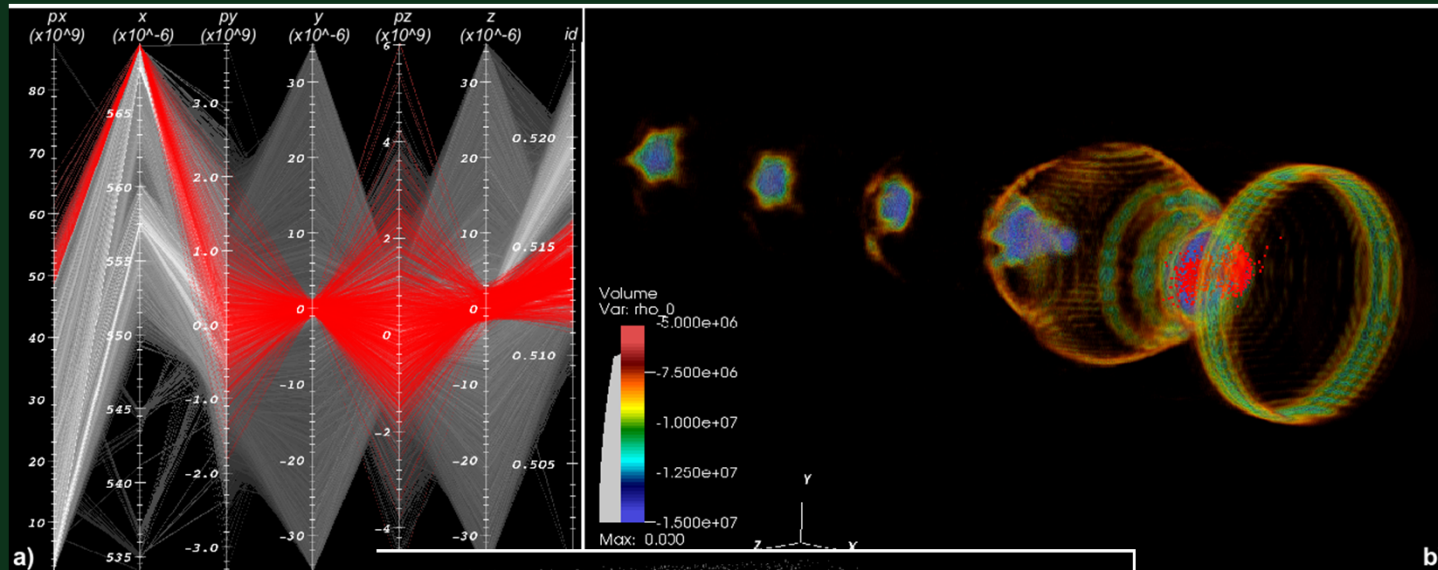  - Novel rendering technique
  - O(B) complexity, rather than O(n)

Traditional        800 bins/variable

Lower gamma        80 bins/variable

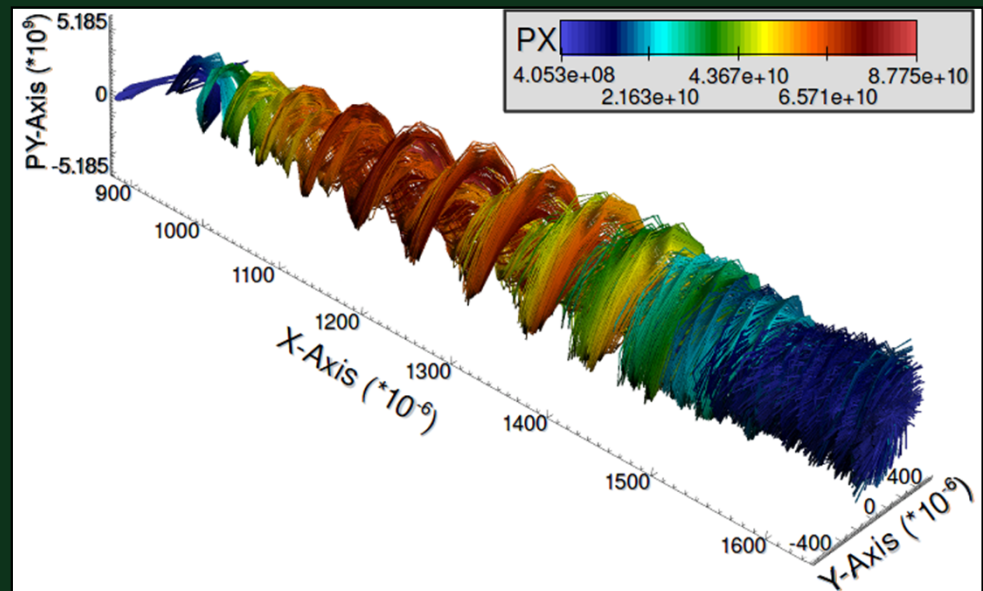# 3D Example
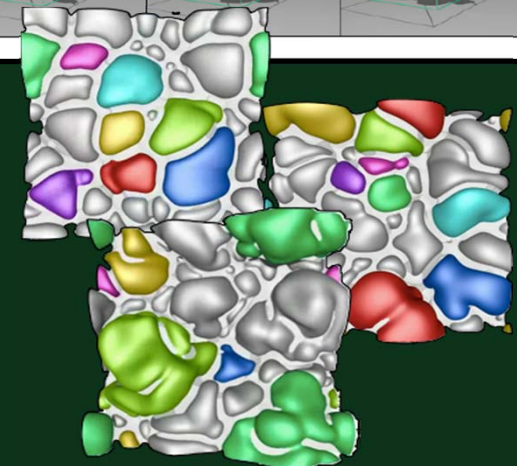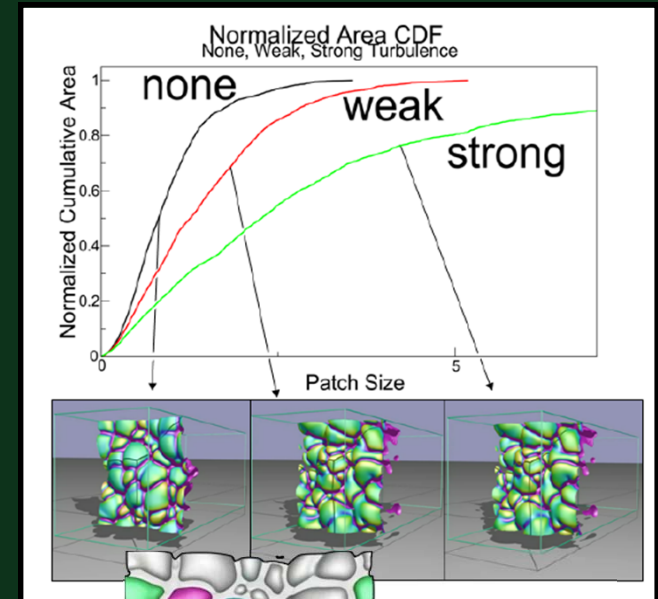
- Understanding particle behavior over time:
  - After finding interesting particles and tracing them through time,
  - Particles start out slow (blue, left), undergo acceleration (reds), then slow again as the plasma wave outruns them (blue, right).
  - Spiral structure shows particles oscillating transversely in the focusing field (new science).

# Analysis of Combustion Simulation Data

- ■ **Problem:** Data of increasing size and complexity increasingly difficult to analyze.

- ■ Accomplishments:
  - • New approaches based upon topological methods offer the means to discover relationships, features, and characteristics in today's largest datasets.

- ■ Science Impact:
  - • First-ever quantitative analysis of large, time-varying combustion simulation data to study influence of turbulence on size/shape of combustion regions in lean, premixed hydrogen flames.

- ■ PI: John Bell (LBNL), SciDAC Community Astrophysics Consortium Partnership, Incite Awardee.

# Hybrid-parallelism: proof at the petascale holds promise for the exascale.

- Existing programming models may not work well at the exascale: multi- and many-core processors.

- Early studies show promise: hybrid-parallel approach outperforms MPI-based approaches on largest-ever visualization runs on DOE supercomputers.

- These results suggest hybrid-parallelism likely a good approach for exascale class machines.



Hybrid-parallel volume rendering of 64-billion zones from combustion simulation on 216,000 cores of JaguarPF at ORNL.
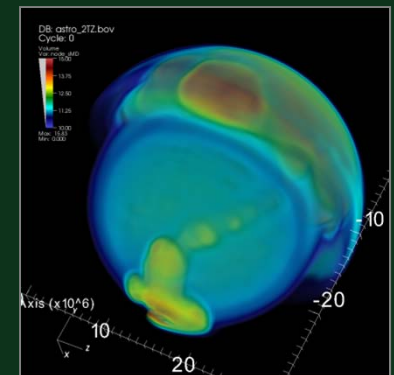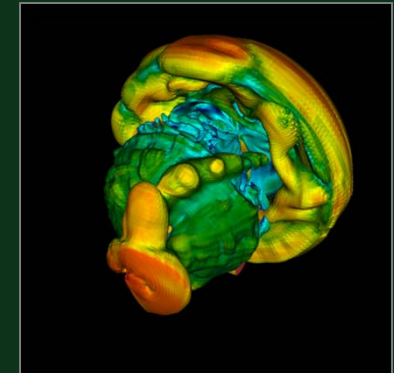
# Summary of Results

- Absolute runtime: -hybrid twice as fast as –only at 216K-way parallel.
- Memory footprint: -only requires 12x more memory for MPI initialization then –hybrid
  - Factor of 6x due to 6x more MPI PEs.
  - Additional factor of 2x at high concurrency, likely a vendor MPI implementation (an $N^2$ effect).
- Communication traffic:
  - -hybrid performs 40% less communication than -only for ghost zone setup.
  - -only requires 6x the number of messages for compositing.
- Image: $4608^2$ image of a ~$4500^3$ dataset generated using 216,000 cores on JaguarPF in ~0.5s (not counting I/O time).

# Production Visualization at the Petascale

- Petascale machines are unique, need visual data analysis tools capable of leveraging the entire resource to ingest and process today's largest scientific datasets.

- SciDAC Visualization and Analytics Center for Enabling Technologies produces such software, proves its effectiveness on all major DOE computational platforms, and distributes it at no charge to the science community.

- Investments in software infrastructure pay off by producing visualization software that can effectively harness the power of today's largest supercomputers for scientific data analysis.





Visualization of supernova simulation results, conducted at 32,000-way parallel on JaguarPF (ORNL) and Franklin (NERSC).
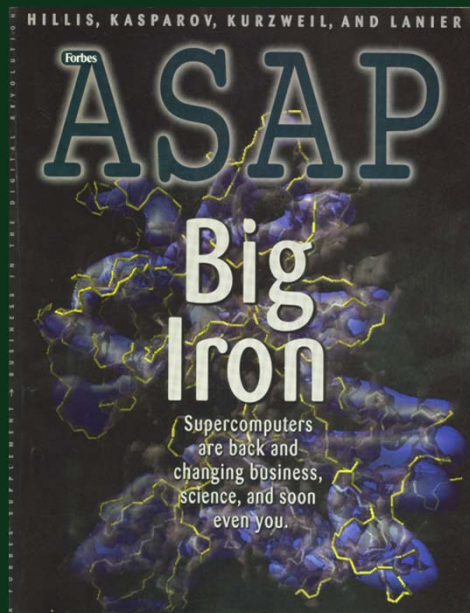
# Glimpse of Current Work (partial)

- **Climate data analysis**
  - Increasingly refined simulations produce data too large for legacy visual data analysis and exploration software.

- **High Performance I/O**
  - You'd be surprised how much time is consumed doing I/O on supercomputers!

- **Topological Data Analysis**
  - New analysis methodology applied to multiple science domains.

- **Carbon Sequestration**
  - Machine learning, computer vision, multivariate analysis, geometric analysis, and visualization help provide traction on understanding how $CO_2$ interacts with porous storage media.

# Summary

- **Berkeley Lab visualization program are field-leading researchers in high performance visualization**
  - Vibrant and productive R&D program
  - Achieving scientific impact through collaborative endeavors
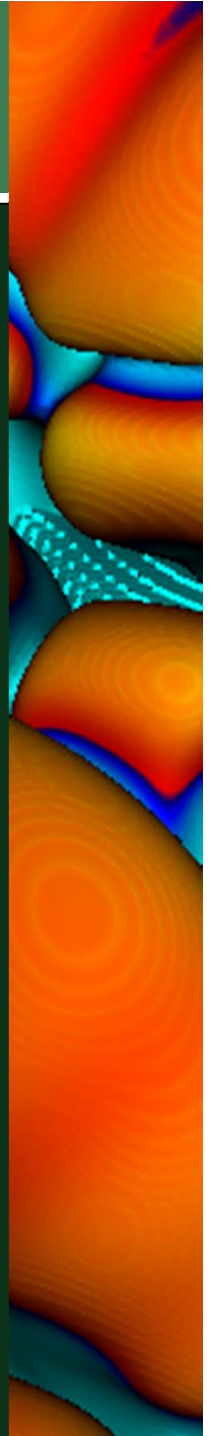- **More information at http://vis.lbl.gov/**

1998
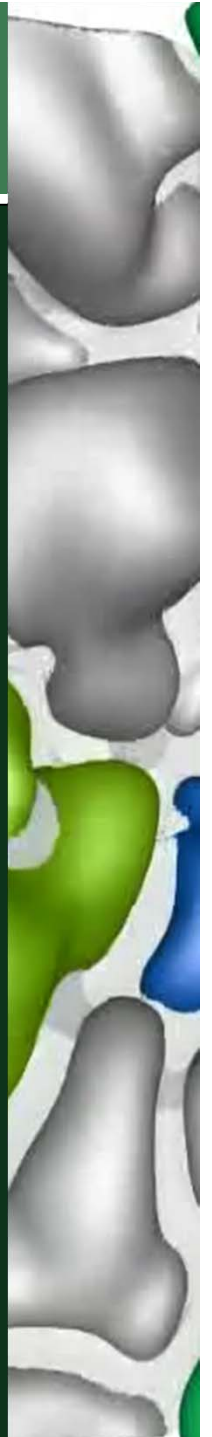
2009

# QDV Cybersecurity Case Study

- The next sequence of slides discusses application of the work to a cybersecurity application – evidence that the idea is generally applicable to large data visualization.

- The team:

  - NERSC Network Security

  - ESnet Network Engineers

  - Scientific Data Management Research

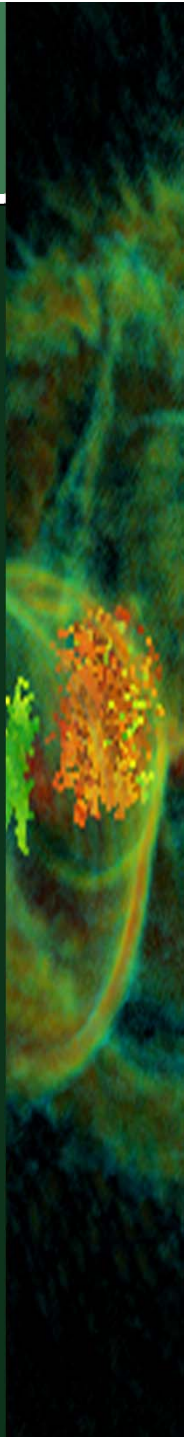  - Visualization Research

# QDV – Detecting Distributed Scans

- The problem:
  - One day's worth of traffic consists of tens of millions of individual connections.
  - Traffic increasing by an order of magnitude every 48 months.
    - ESnet monthly traffic levels now exceed 1 PB.
  - The Internet is a hostile environment, and it will get worse.
  - Objective: enable rapid forensic data analysis (network flow records).

# QDV – Detecting Distributed Scans

- ## The data:
  - 42 weeks' of connection records from Bro (NERSC).
  - 281GB for raw data, 78GB for compressed bitmap indices.

- ## "Hero-sized problem"
  - No previous network analysis work has ever attempted to perform interactive visual analytics on data of this scale (ca. 2006).
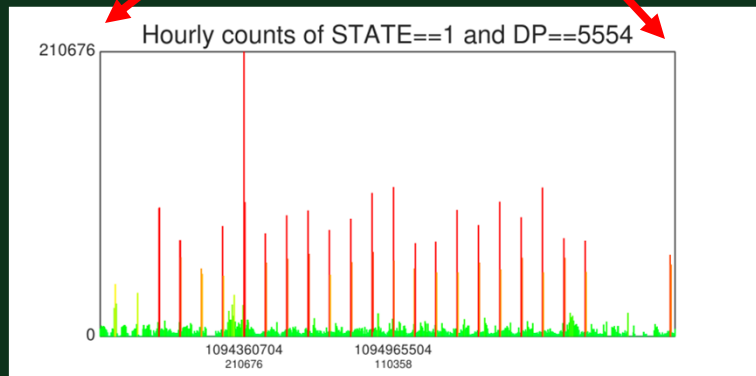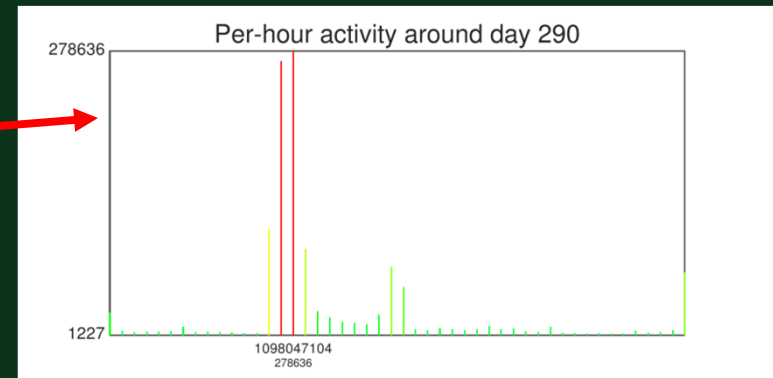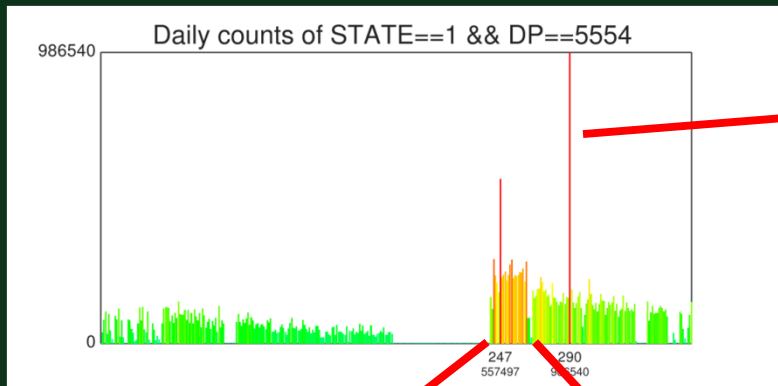  - Result: what once took days (if at all possible) now takes seconds.

# QDV – Detecting Distributed Scans

- ## The starting point:
  - You are a network security analyst
  - Your beeper goes off [at lunch, in the shower…]
  - You receive an alert that "something odd is happening with the network…IDS showing unusual levels of activity on port 5554"
  - Your job – answer questions:
    - What's going on now?
    - How long has this been happening?
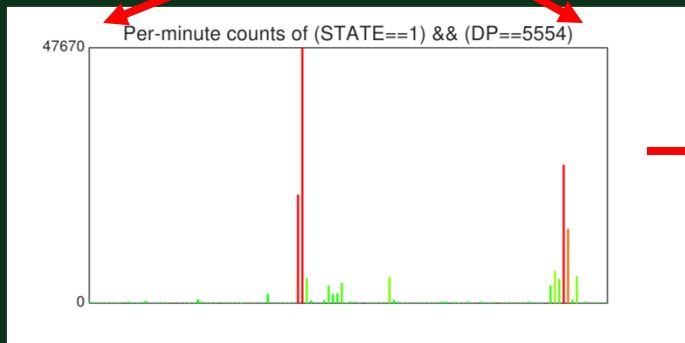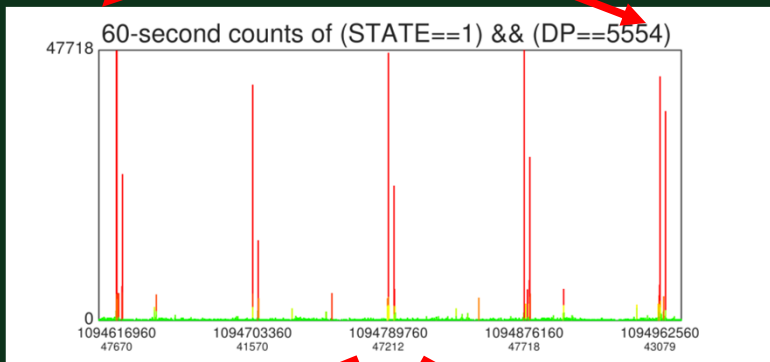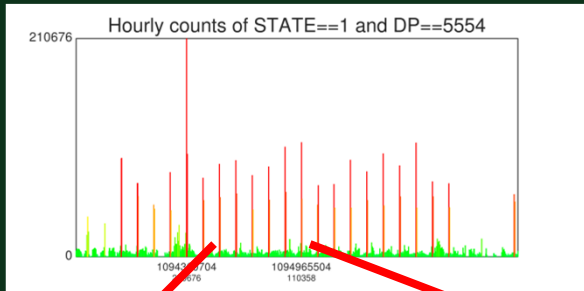    - Implications?

# QDV – Detecting Distributed Scans



Daily counts of STATE==1 && DP==5554



Per-hour activity around day 290



Hourly counts of STATE==1 and DP==5554

1. Query to produce a histogram of unsuccessful connection attempts over a 42-week period at one-day temporal resolution (upper left).

2. Drill into the data, query to produce a new histogram covering a four-week period at one-hour temporal resolution (lower left).

3. Generate a histogram of one-hour resolution over a two-day period around day 290 (upper right).
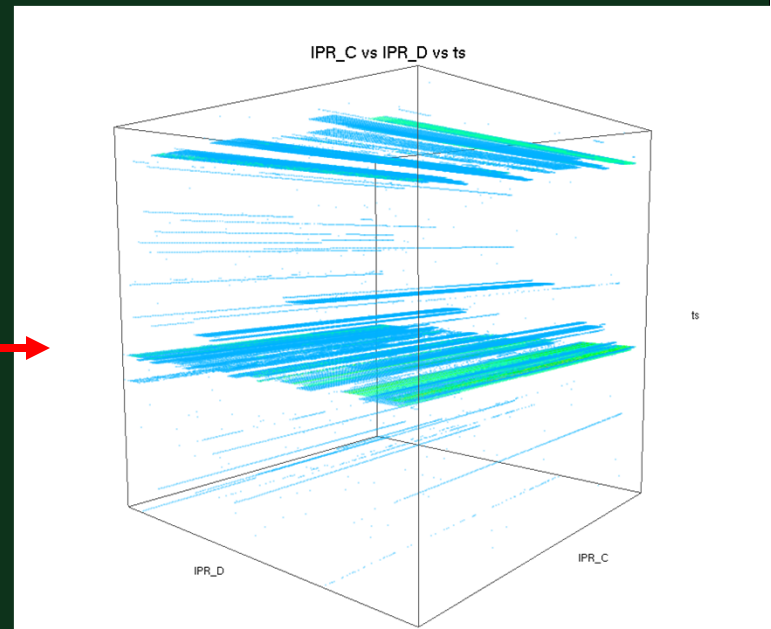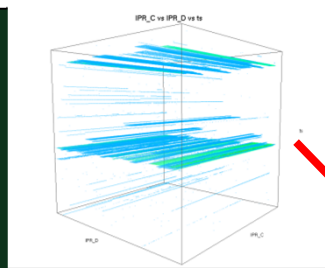
# QDV – Detecting Distributed Scans



5. Query to generate a histogram of unsuccessful connection attempts over a five-day period sampled at one-minute temporal resolution (middle, left). Regular attacks occur at 21:15L, followed by a second wave 50 minutes later.

6. Query to generate histogram over a two-hour period at one-minute temporal resolution (lower left).

7. Query to generate a 3D histogram showing the coverage of attacks in destination address space (lower right).
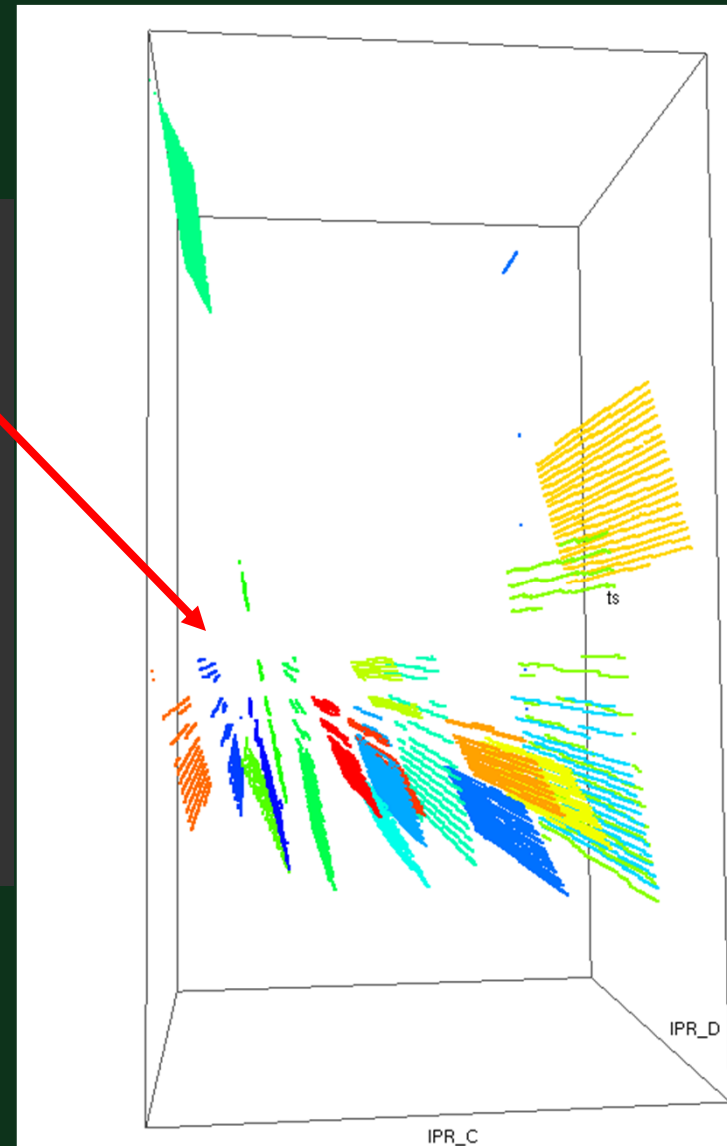
# QDV – Detecting Distributed Scans



After establishing that (1) a temporally regular activity is occurring, and (2) that it is in fact a systematic probe (scan) of entire blocks of network addresses, the next task is to determine the set of remote hosts participating in the attack.

Working backwards, we isolate the A, B, C and D address octets of the hosts participating in the attack.

This image shows a 3D histogram of the destination address space being attacked by each of 20 different hosts. The vertical axis is time – a seven-minute window at one-second temporal resolution.

# QDV – Detecting Distributed Scans

- Our analysis was performed in statistical space only.

  - We never accessed the raw data.

  - Our processing and visualization used only the index data.

  - Performance study focuses on parallel algorithms for multidimensional histogram computation from compressed bitmap indices.