# Occam's Razor and Petascale Visual Data Analysis

E. Wes Bethel

Lawrence Berkeley National Lab

26 Sept 2007

# Outline

- Occam's Razor and Petascale Visual Data Analysis?

- Challenges:
  - Limited cognitive bandwidth
  - Too much data
  - Not enough time
  - Production deployment

- Conclusion

# Occam's Razor

- Principle: "less is more"
  - Explanation of phenomena should make as few assumptions as possible, eliminate assumptions that make no difference in observable predictions of the explanatory hypothesis or theory.
- Attributed to 14[th] century English logician Franciscian friar William of Ockham
- Foundation of scientific method
  - "We are to admit no more causes of natural things than such as are both true and sufficient to explain their appearances." – *Sir Isaac Newton*
  - "Everything should be made as simple as possible, but not simpler." – *Albert Einstein*

# Occam's Razor Scooter

1. Supernatural being formed committee to design universal constants.

2. Committee of deities created gravity.

3. Gravity causes Razor scooter to roll down the hill.

1. Supernatural being created gravity.

2. Gravity causes Razor scooter to roll down the hill.

1. Gravity causes Razor scooter to roll down the hill.

# Occam's Razor

- **Physics**
  - Einstein's theory of special relativity vs. Lorentz's theory that rulers contract and clocks slow down when in motion through the Ether.
    - Equations are "the same."
    - Ether not detectable with Lorentz's equations.
  - Justification of Heisenberg's Uncertainty Principle in quantum mechanics
    - Impossible to know exact position and momentum of a particle at the same time.
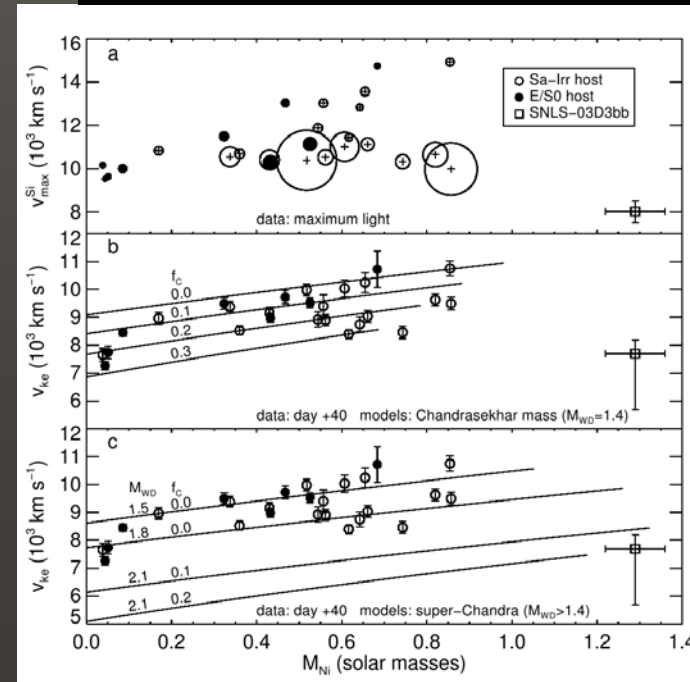    - Why?
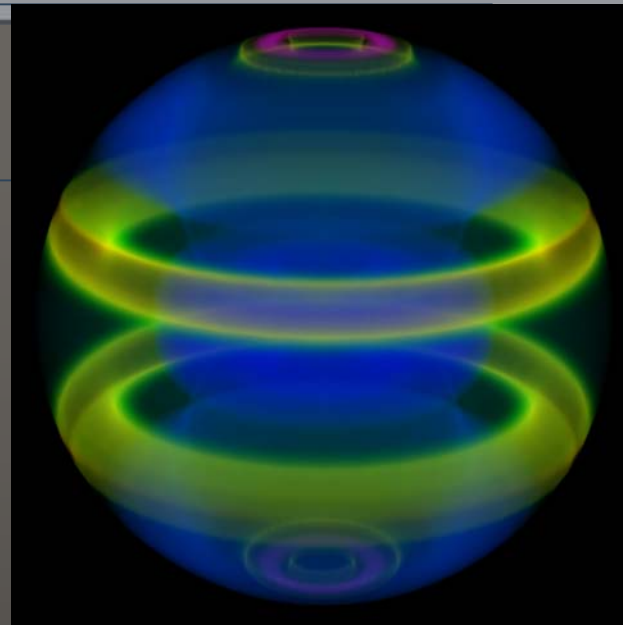    - The process of observation influences the observed.

# Occam's Razor and Petascale Visual Data Analysis?

- What is the minimum needed to do "the job?"
- What is "the job?"
  - Visualization: a flexible, powerful knowledge discovery technology/set of methodologies.
  - Can do basically "anything"
  - "Any sufficiently advanced technology is indistinguishable from magic." – *Arthur C. Clarke*
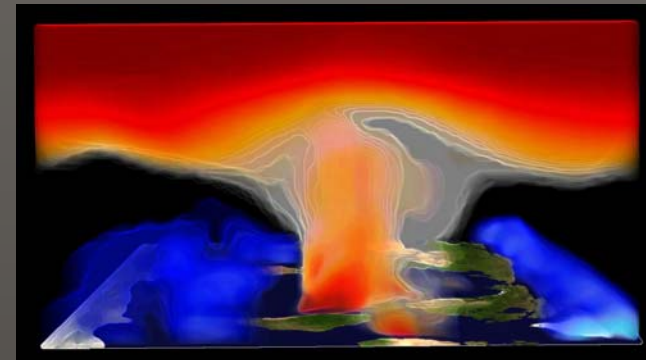- With infinite flexibility, must choose carefully

# Visualization Use Models

- ## Presentation visualization
  - You know what's there and want to show it to someone else

- ## Analytical Visualization
  - You know what you are looking for

- ## Discovery Visualization
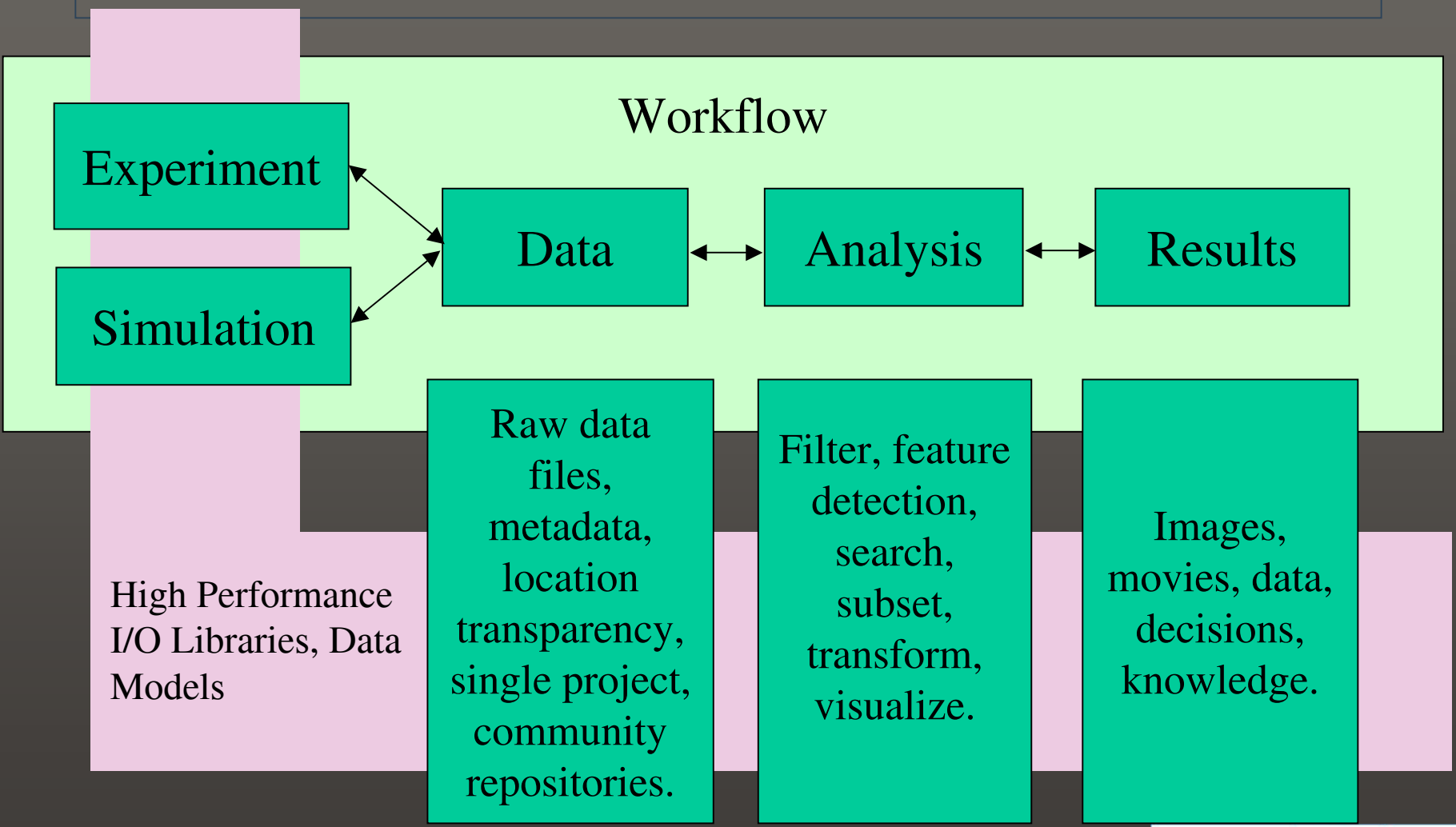  - You have no idea what you're looking for

# Brief Digression – Theory of Education

- **Three stages of learning:**
  - Romance
    - "Play and wonder," which are initial steps towards
  - Precision
    - Deeper understanding leading to
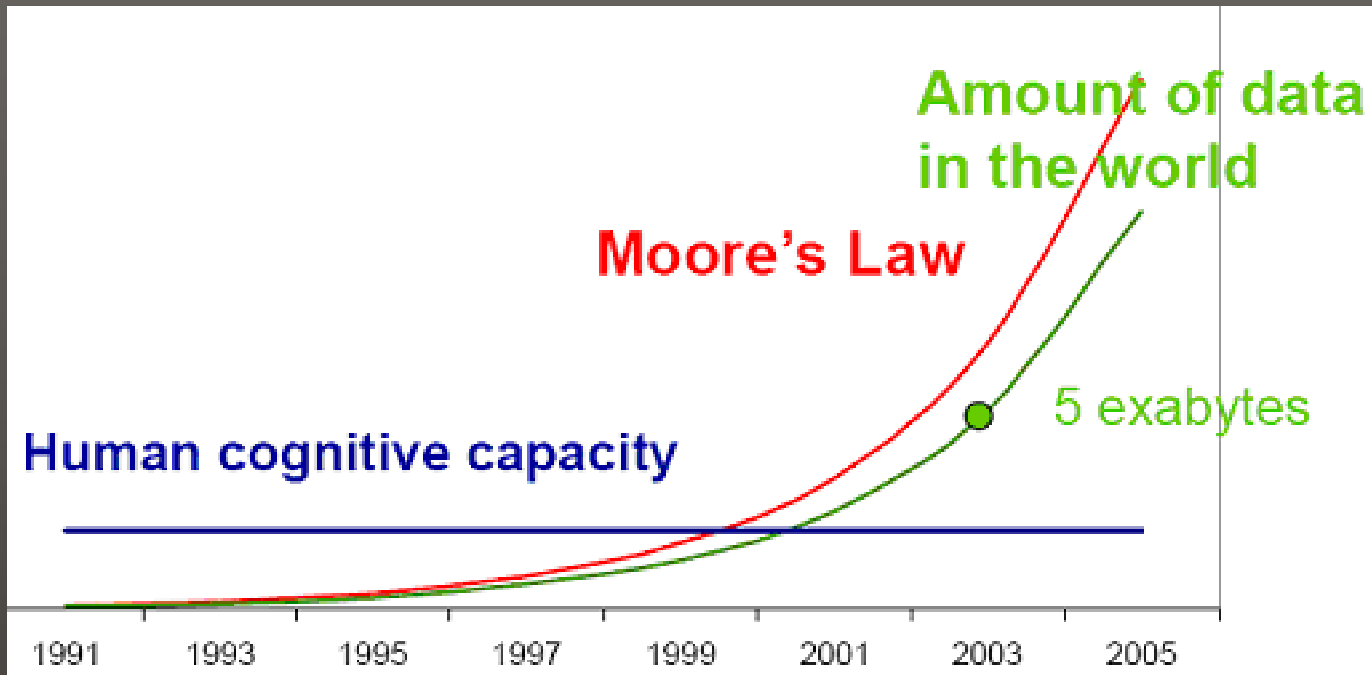  - Generalization
    - Broad insights

# Map of the Problem Space

Workflow

Experiment

Simulation

Data

Analysis

Results

High Performance I/O Libraries, Data Models

Raw data files, metadata, location transparency, single project, community repositories.

Filter, feature detection, search, subset, transform, visualize.

Images, movies, data, decisions, knowledge.

# The Big Challenge

- How to effectively enable knowledge discovery at the petascale given:
  - Limited cognitive bandwidth
  - Too much data
  - Not enough time
  - Deploy _production quality_ tools for scientific communities
- Rest of discussion examines these topics in more detail.

# Limited Cognitive Bandwidth

# The Cognitive Bandwidth Challenge



Experiment

Simulation

Data

Analysis, Filter

Vis

Render

# Cognitive Bandwidth Challenge

- "Picture worth a 1000 words." F. Barnard, 1921, N. Bonaparte c. 1800.
  - At a glance, humans are great at perceiving trends, anomalies, deriving meaning and understanding from seemingly "random" data.
- "One equation is worth 1000 pictures." J. Bell, 2006.

# Cognitive Bandwidth Challenge

- Can't "see" a TB, much less a PB
  - Mesh cells map to subpixels
  - Depth complexity
  - Spatial complexity

- Occam might say:
  - Show only what is important to more closely match cognitive capacity.

- Counterexample:
  - Visual simulation and ultra high resolution displays

# Photorealistic Rendering of Molecules

# Visual Simulation – UCLA's UST

# Data Tsunami Challenge

- Science bottleneck: management of and knowledge discovery from large collections of scientific data

# Reducing I/O (and Compute) Load

- Adaptive Mesh Refinement
- Compression
- Subsampling, statistical models
- Other issues: storage, access, high performance I/O, movement, sharing, provenance, ontology, …

# Another Approach – In-place Analysis

- Idea: do visualization and analysis in the simulation, save the analysis results rather than the simulation data.
  - Pro:
    - Avoids I/O and storage issue by having the simulation perform analysis and save only analysis results.
  - Cons:
    - No possibility of unexpected discovery
    - Assumes image is final analysis result
    - Simulation code "feature creep"

# Another Approach – Focus on Interesting

- Restrict visual analysis and associated I/O, processing to "interesting data."
  - Pros:
    - Applicable to many problem domains, including both experimental and computational data
    - Applicable to all visualization use modalities
  - Cons:
    - Still have the I/O and storage problem
- Comment
  - Good fit for human learning

# Query-Driven Visualization

- Focus visualization processing on subsets of data deemed to be "interesting."
  - "Interesting" is something the user needs to define.

- Challenges
  - How to define "interesting."
    - Formulation of definition (domain-specific).
    - Expression of definition (semantic).
  - Find interesting data quickly (data management).
  - Effective visual presentation of "interesting data" (visualization).
  - Architectures/deployment that complements existing visualization algorithms and applications (computer science).

# QDV and Fusion

- ## GTC is a PIC code for modeling microturbulence

  - Top: all particles from a single timestep

  - Bottom: particles that undergo "trapping" at least 20 times.

  - *New VACET project: statistical downsampling to reduce I/O load.*



www.vacet.org

# QDV and Fusion

- Application: fusion microturbulence

- Objective: interactive multidimensional filtering to locate and analyze interesting phenomena.

- Result: capability in production software (VisIt).

# QDV at the Petascale

- Problem: want to do better than O($n$) for data analysis (where $n$ is the size of the data).
- Solution(s):
  - Use state-of-the-art index/query. VACET and SDM Center integrating such technology at the data I/O layer in a way that is transparent to simulation code developers.
    - Note: tree-based index/query systems suffer from "*The Curse of Dimensionality*"; compressed bitmap indices do not:
    - **O(N\*\*D) vs. O(N\*D): Trees vs. CBI**
  - Leverage this capability in visual analysis tools, can provide assistance for use in other types of tools.

# Not Enough Time Challenge

- Topics
  - Make it go faster
    - Software architecture for production quality high performance visualization on DOE's leading computational platforms (Hank's talk).
  - Make it do more, make it do what I want.
    - Analysis and visualization
    - Production analytics pipelines
    - Community-centric

# Visualization and Analysis

- Objective: convey deeper meaning than possible with "data exploration"
  - Relationships between variables
  - Characteristics of data
  - Compute, display and analyze "features"

# Analysis and Visualization

- Premise: find and analyze "features" in data.
- Methods of location and analysis vary:
  - PCA, ICA, Machine Learning, Support Vector Machines, etc.
  - Topological analysis.
- Impact:
  - Quantitative analysis
  - Traction on "big data problems"
  - One potential basis for comparative analysis (rather than "chi-by-eye")
- Examples

# Analysis and Visualization – Climate

**Extracted features can be used as templates for finding similar features.**

*Tropical storm visible in sea level pressure simulations at multiple time steps.*

**In this case, the features were variations on rotating low-pressure systems. This was not assumed a priori.**

C1  C2  C3  C4

C5  C6  C7  C8

C9  C10

*Images of extracted features: top ten independent components were extracted from set of all 8x8 subimages.*

# Analysis and Visualization – Topology



- Channel structures in porous media : green isosurface separates solid material and empty space; red curves, connecting maxima and 2-saddles, represent channels.

- Combustion kernel feature identification, tracking and analysis. Kernels appear, merge, and extinguish over time. Why? How many?





(a)

# Discovering Relationships

- How are variables related to one another?

- $O_2$ and $CO_2$ concentration (left and middle), correlation field (right)

- Same idea, but in 3D and mapped onto varying isotherms.





(a)    (b)    (c)    (d)    (e)    (f)

# Production Analytics

- ## Visual Data Cartography
  - ### Google Earth, Google Sky, Sloan Digital Sky Survey, etc.

[slide by S. Bailey]

# SNFactory Pipeline circa 2005

NEAT
(Near Earth Asteroid Tracking images)

**Supernova search**

Email, cut-and-paste

**Scanning**

Scan, view websites, wait for data (HPSS, web)

**Vetting**

Manual scheduling

**SNIFS data taking**

Manual data download, custom Perl scripts

**Postmortem spectra processing**

Cut-and-paste

**Twiki page**

= needs human intervention

# SNFactory Pipeline – Present Day

# Sunfall: Supernova Warehouse

- A comprehensive supernova data management, workflow visualization, collaborative scientific analysis tool

# Sunfall: Supernova Warehouse

- Improved situational awareness, decreased repetitive labor, more science!



www.vace

# Production Analytics Impact on Science

- Significant, measurable effect on a major science project, with impact on physics grand challenge (dark energy).
    - Reduced number of false positive candidate supernovae by 80%
    - Up to 90% reduction in labor costs in areas of SNfactory pipeline
    - This freed up more time for science
        - SNF pubs 2005: 0 refereed
        - SNF pubs 2006 & 1st 3 mos 2007: 3 refereed (1 submitted)

# Production Analytics Impact on Science

- From the researchers:
  - "We're awash in Supernova"
  - More SN discovered in past year than in all previous history
- Observation
  - This system the product of a multidisciplinary team of dedicated scientific staff
  - Mission focus guides R&D

# Production Analytics - Challenges

- Coupled code projects:
  - Fusion: FACETS
  - Accelerator: COMPASS
  - Climate: Earth System Model
- Data from one code is input to another
- V&V at each stage
- Model and result is multimodal, multivariate (and multiscale)

# Challenge: Software Architecture and Engineering

- Objective: production-quality, petascale-capable visual data analysis software (Hank's talk)

- Challenges:
  - Who is the user/stakeholder?
  - What platforms?
  - Use models?
  - What is effective balance between R&D?

# Conclusion

- Petascale visual data analysis is a problem-rich environment.

- Given limited bandwidth (cognitive, network, processing, etc.) and time, let Occam's Razor help as a guide to eliminate unnecessary motion.

- We've discussed some, but not nearly all, of the challenges in this space.

# Next Up

- ## Hank Childs, LLNL
  - *Large-scale viz (GNEP connection)*

- ## Kwan-Liu Ma, UC Davis
  - *Parallel viz. pipeline and insitu Feature detection/tracking*