

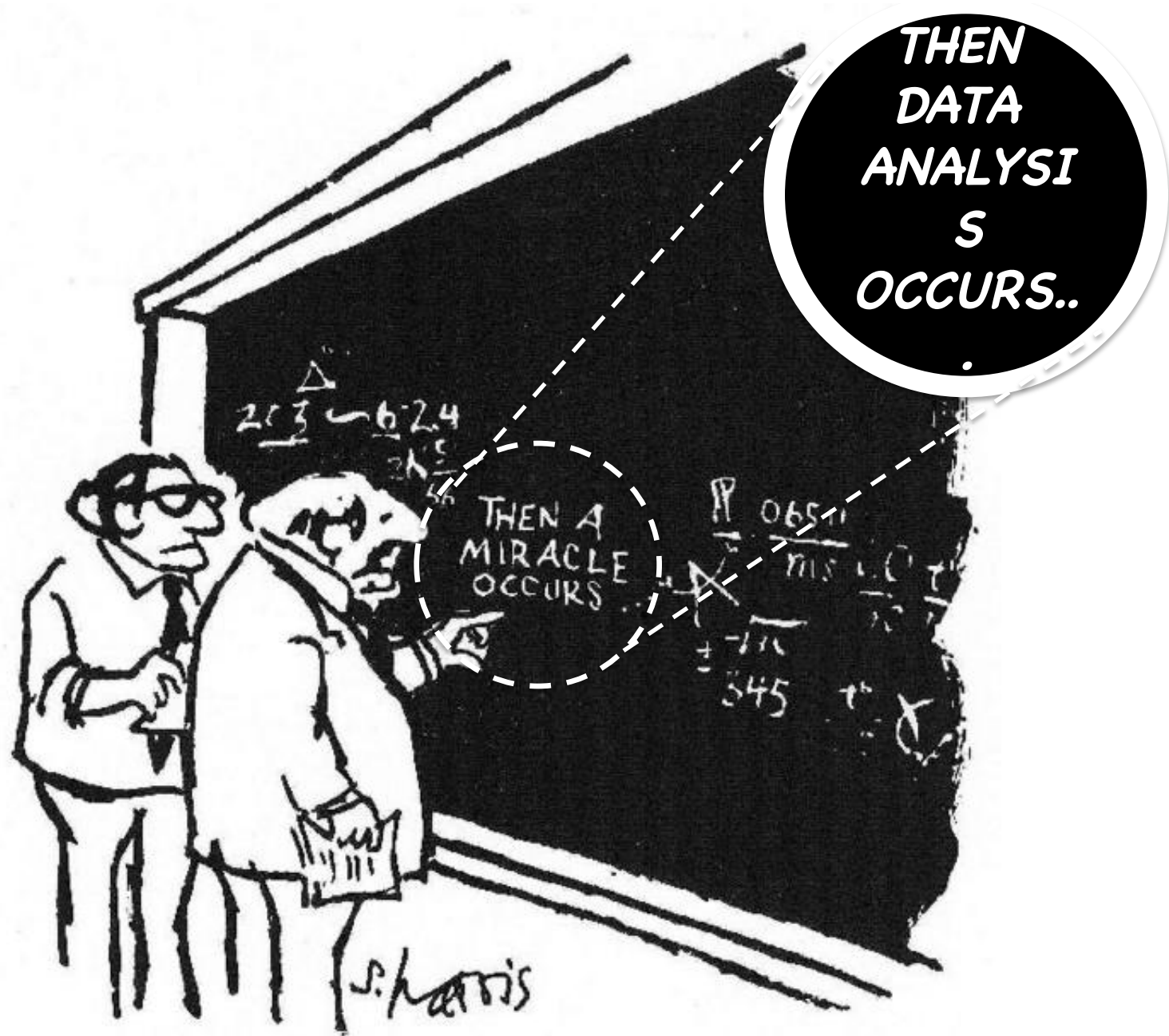


# **The Bigger the Data, the Harder it Falls: Current Problems in Large Data Analysis**

**David H. Rogers**

**Manager, Dept 1424  
Data Analysis and Visualization Department  
Sandia National Labs**

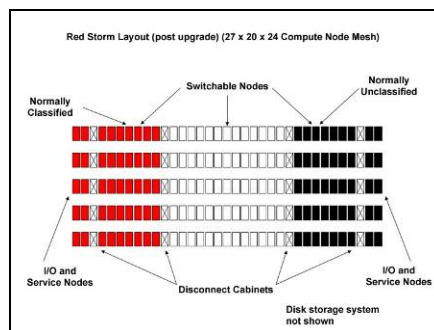
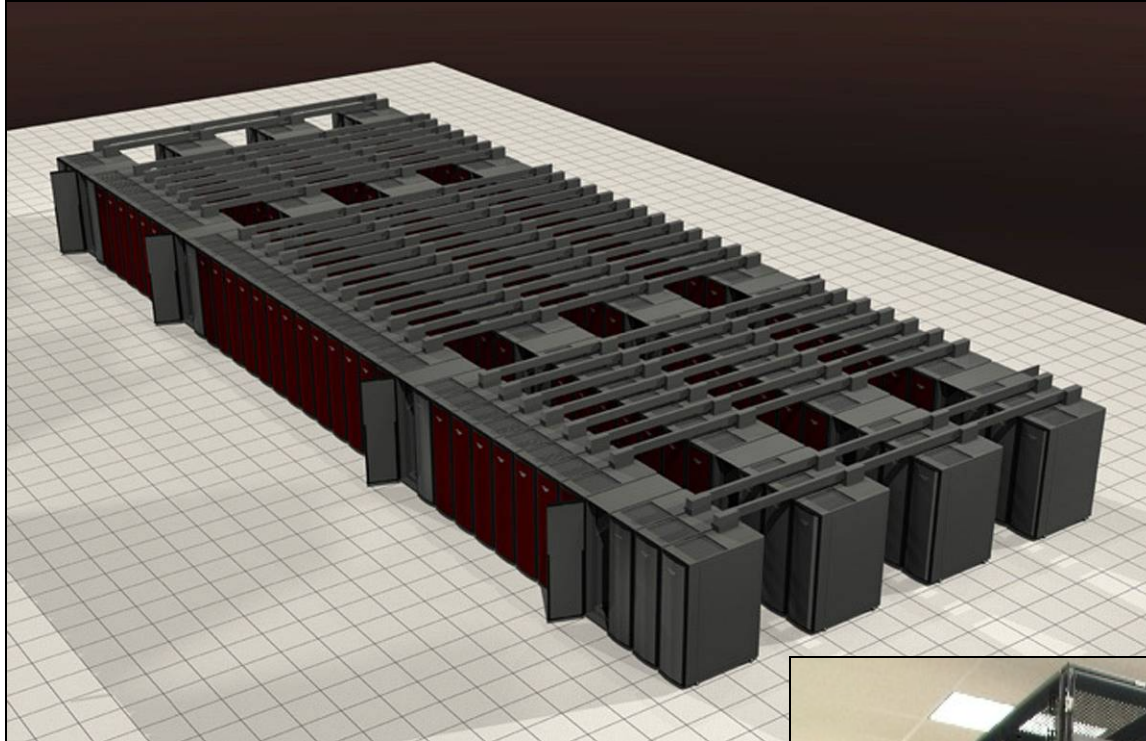




“I think you should be more explicit here in step two.”



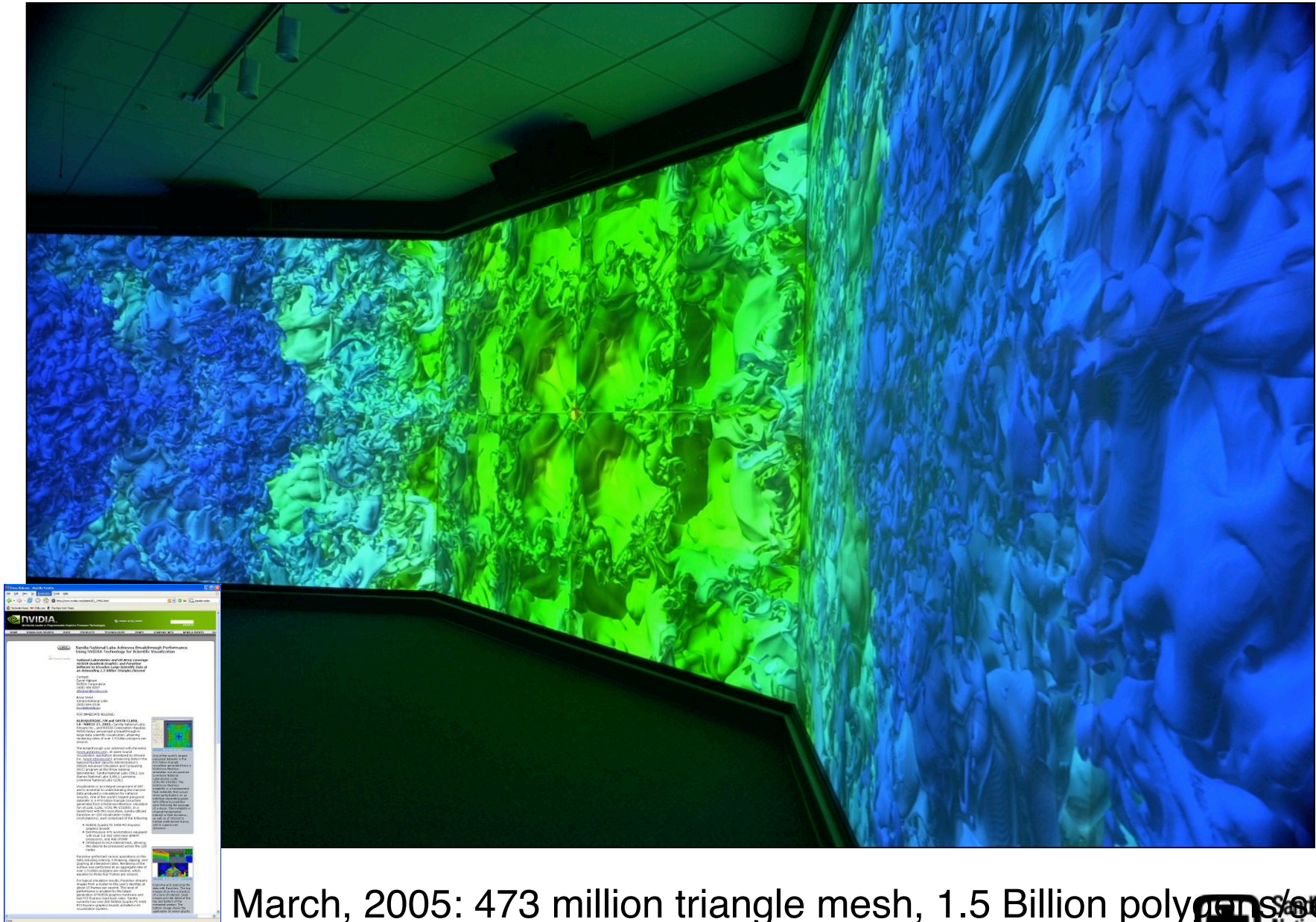
# Interactive Visualization of Extremely Large Data







# Interactive Visualization of Extremely Large Data



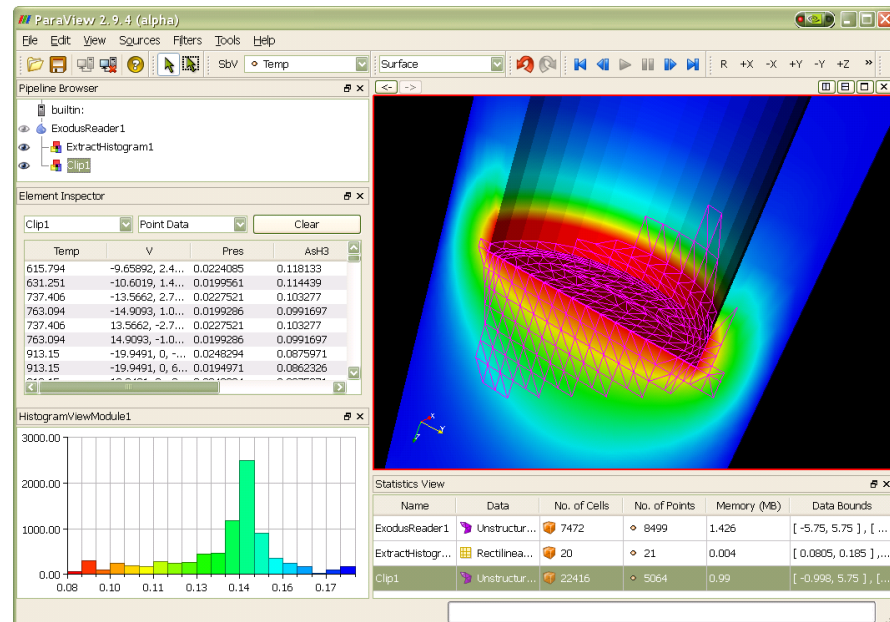
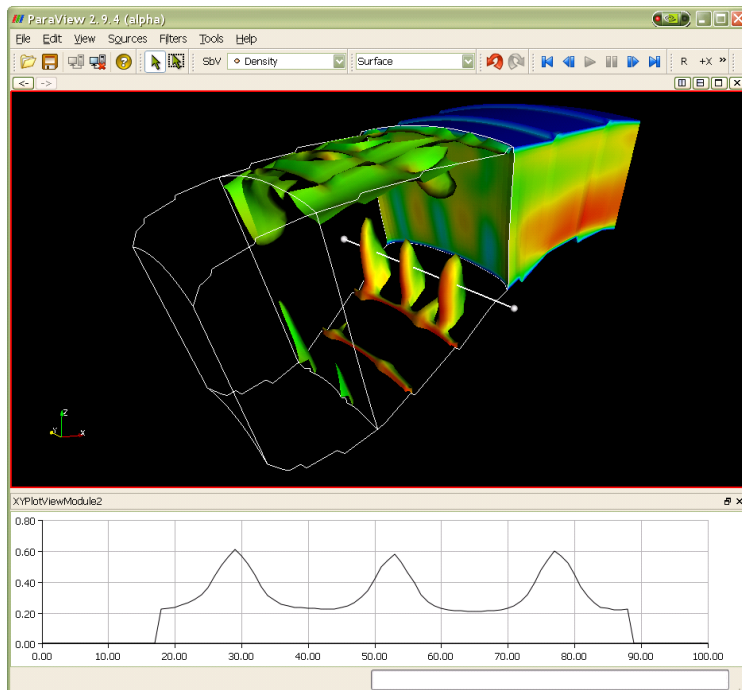
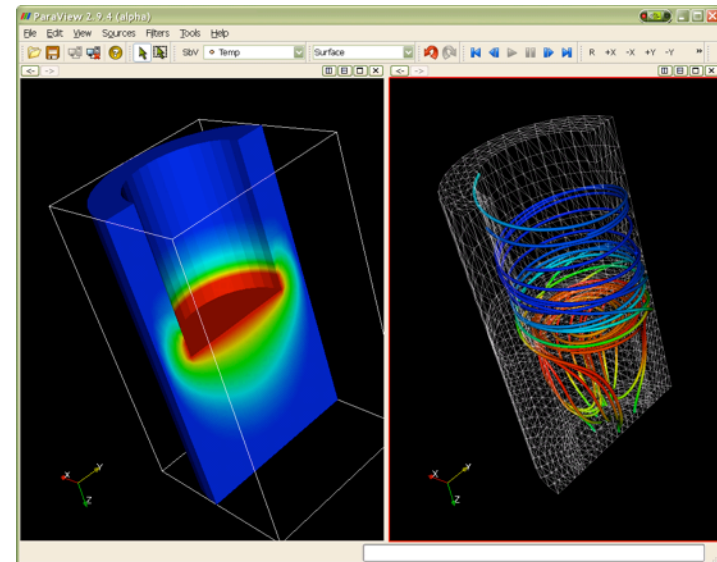
March, 2005: 473 million triangle mesh, 1.5 Billion polygons/sec.





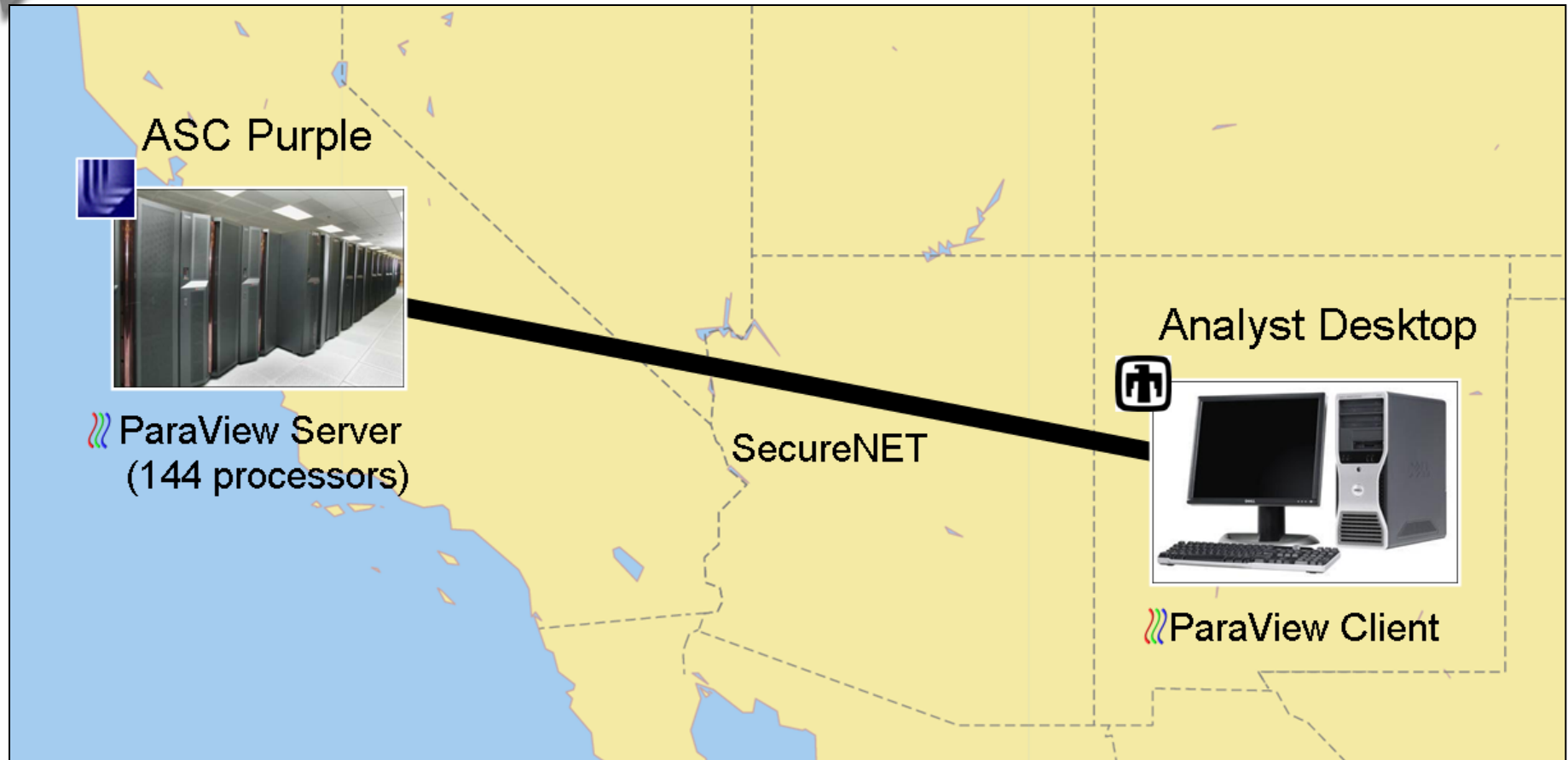
# Advancing Scalable Scientific Visualization with ParaView 3.0

```
Python Shell
Python 2.4.3 (#69, Mar 29 2006, 17:35:34) [MSC v.1310 32 bit (Intel)] on win32
>>> import paraview
>>> paraview.ActiveConnection = paraview.Connect()
>>> reader = paraview.CreateProxy("sources", "ExodusReader")
>>> print dir(reader)
['ListProperties', 'SMProxy', '__doc__', '__eq__', '__getattro__', '__init__', '__iter__',
 '__module__', '__ne__', '__pyProxy__AddProxy', '__pyProxy__AddToProperty',
 '__pyProxy__CreateDisplayProxy', '__pyProxy__GetProperty', '__pyProxy__RemoveFromProperty',
 '__pyProxy__SaveDefinition', '__pyProxy__SetProperty']
>>> print reader.ListProperties()
['ApplyDisplacements', 'BlockArrayInfo', 'BlockArrayStatus', 'CellArrayInfo', 'CellArrayStatus',
 'DisplacementMagnitude', 'DisplayType', 'ExodusModelMetadata', 'FileName', 'FilePattern',
 'FilePatternInfo', 'FilePrefix', 'FilePrefixInfo', 'FileRange', 'FileRangeInfo',
 'GenerateBlockIdCellArray', 'GenerateFileIdArray', 'GenerateGlobalElementIdArray',
 'GenerateGlobalNodeIdArray', 'HierarchyArrayInfo', 'HierarchyArrayStatus', 'MaterialArrayInfo',
 'MaterialArrayStatus', 'NodeSetArrayStatus', 'NodeSetInfo', 'PointArrayInfo', 'PointArrayStatus',
 'SideSetArrayStatus', 'SideSetInfo', 'TimeStep', 'TimeStepRangeInfo', 'TimeStepValues',
 'XMLFileName', 'XMLFileNameInfo']
>>>
```





## Distance Visualization of Terascale Data with ParaView



- Out-of-the-box solution (despite complexity of hardware)
- ParaView 2x faster than other solutions (27 million unstructured cells)
- ParaView only workable solution for large data (274 million unstructured cells)



# Principles Resulting from Large Scientific Data Analysis

## Scalability

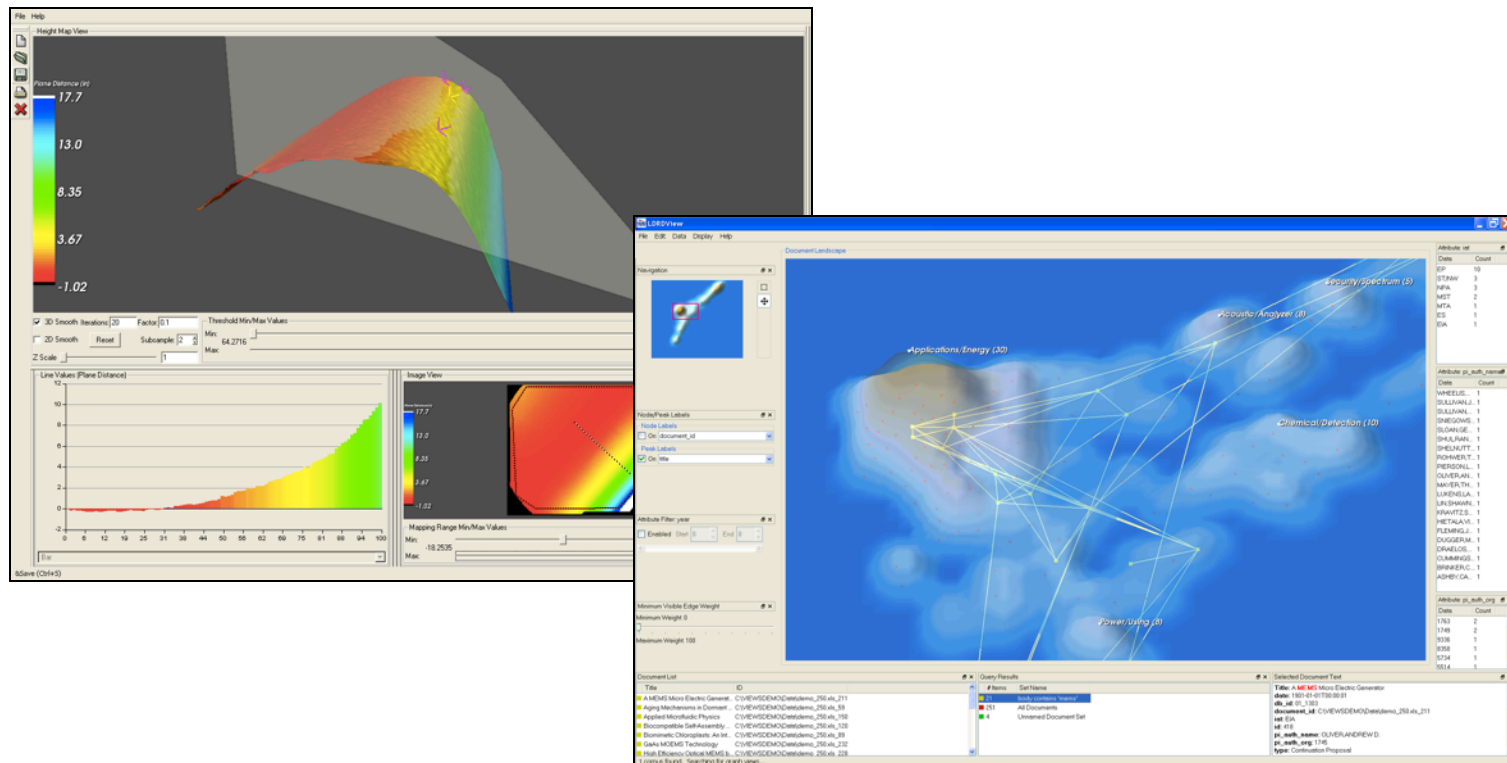
- Run the same software on everything from laptop to cluster
  - Data size shouldn't matter to the analyst
- Algorithms, data structures and visualization must adapt to data size
- Client (user-side) software can connect to appropriate remote hardware
- Tools must be cross-platform
  - Linux, Windows, Mac
- Responsive innovation leads to success
  - Balance between research, long term planning and immediate problems



# Targeted Tool Development

Rapid Prototyping, with flexible technologies (Qt, VTK, etc.)

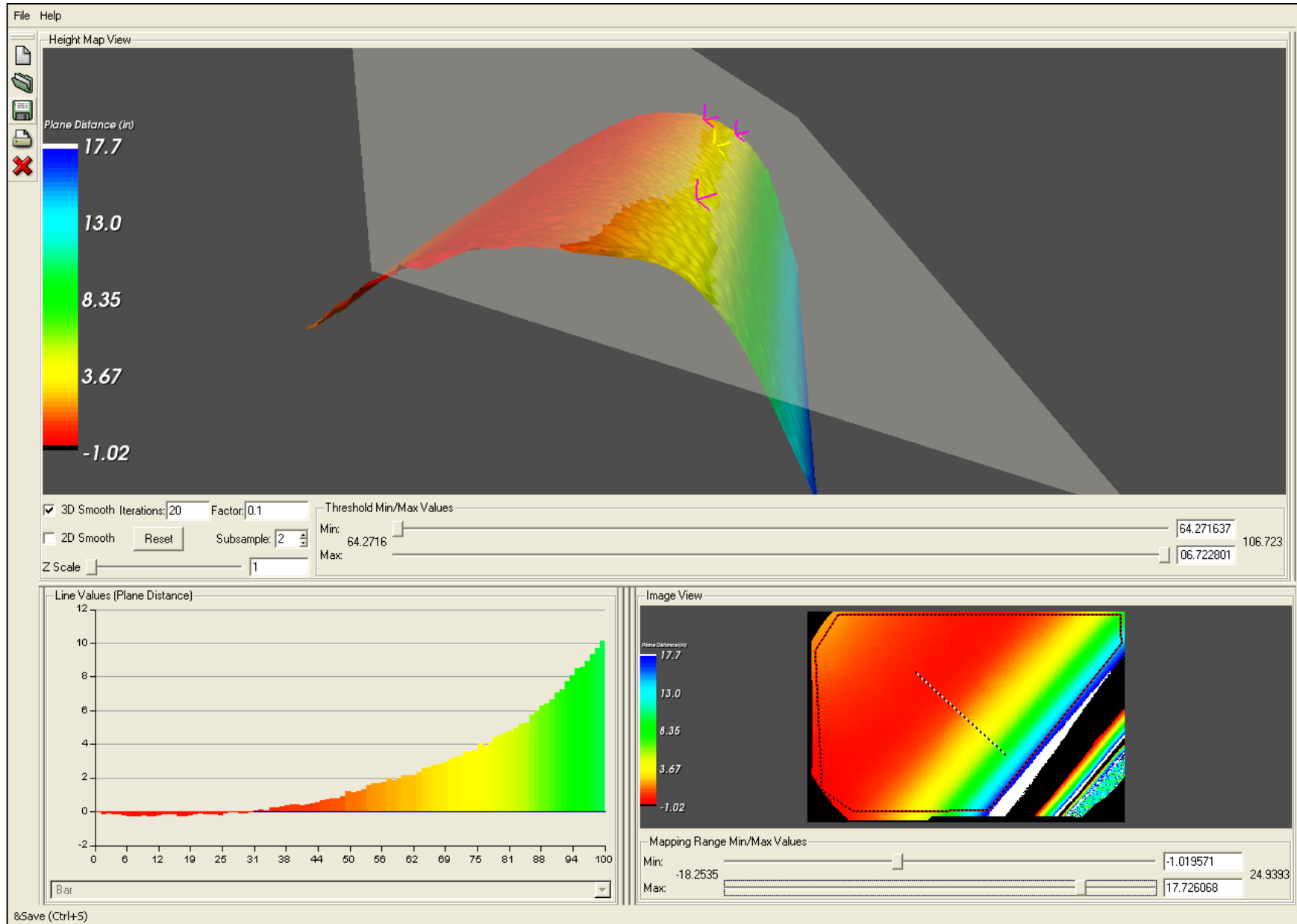
- Promotes interactive, cross-platform tool development
- Attacks the ‘real problems’
- Addresses customer concerns up front
- Promotes creative, collaborative development





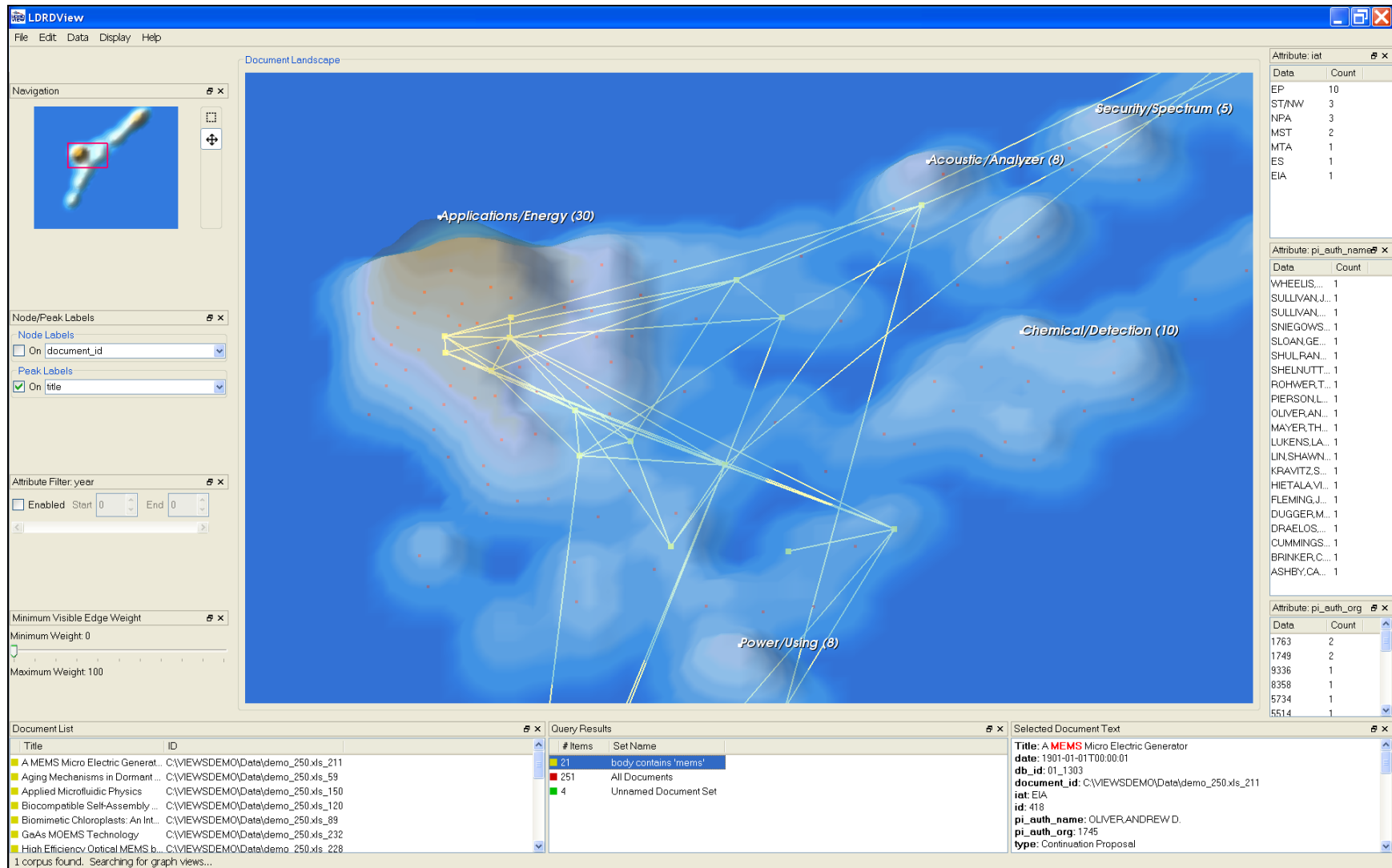


# Analysis of NASA Telemetry Data (Leading Edge of Orbiter Wings)



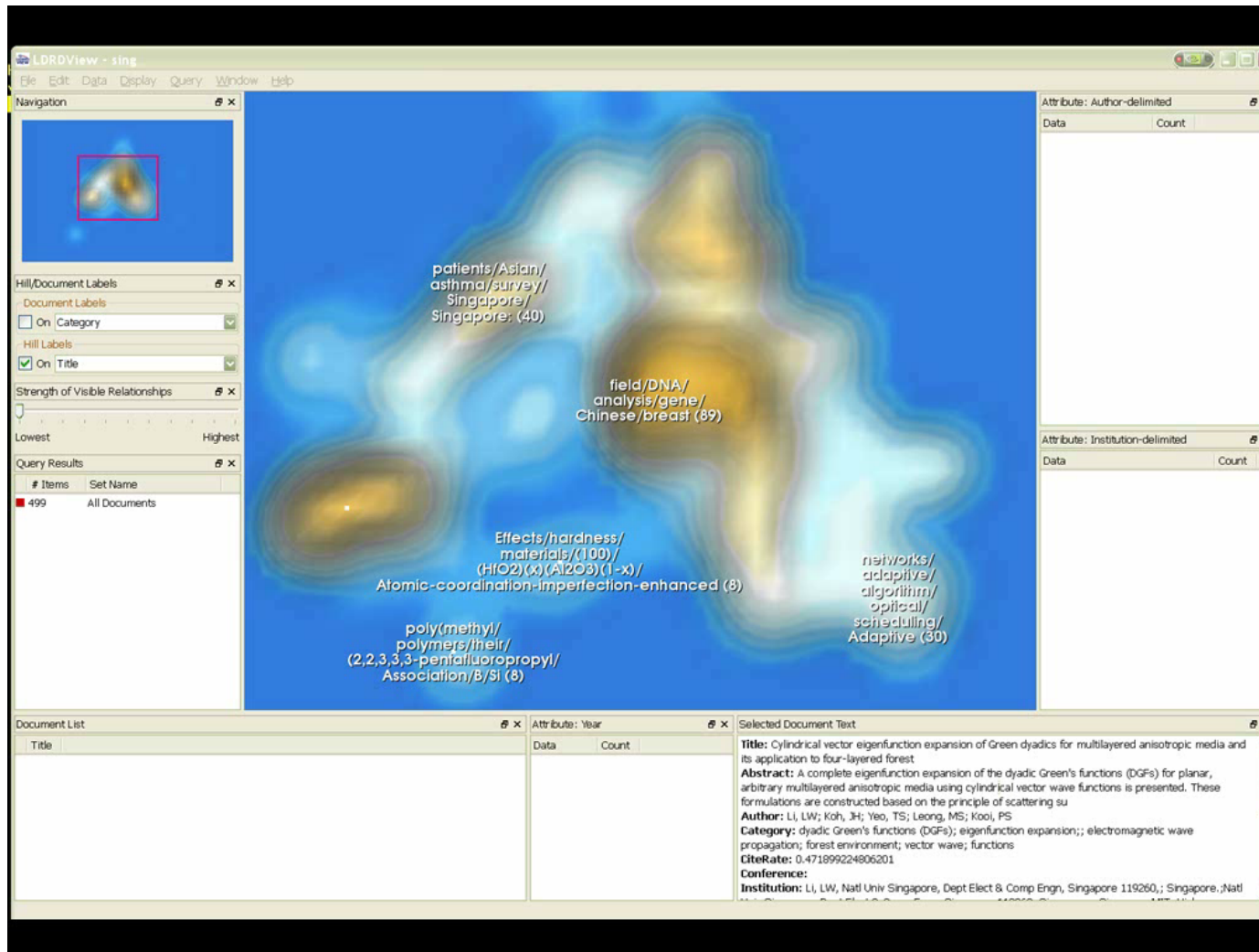


# Accelerating Insight into Complex Data with LDRDView





# Accelerating Insight into Complex Data with LDRDView



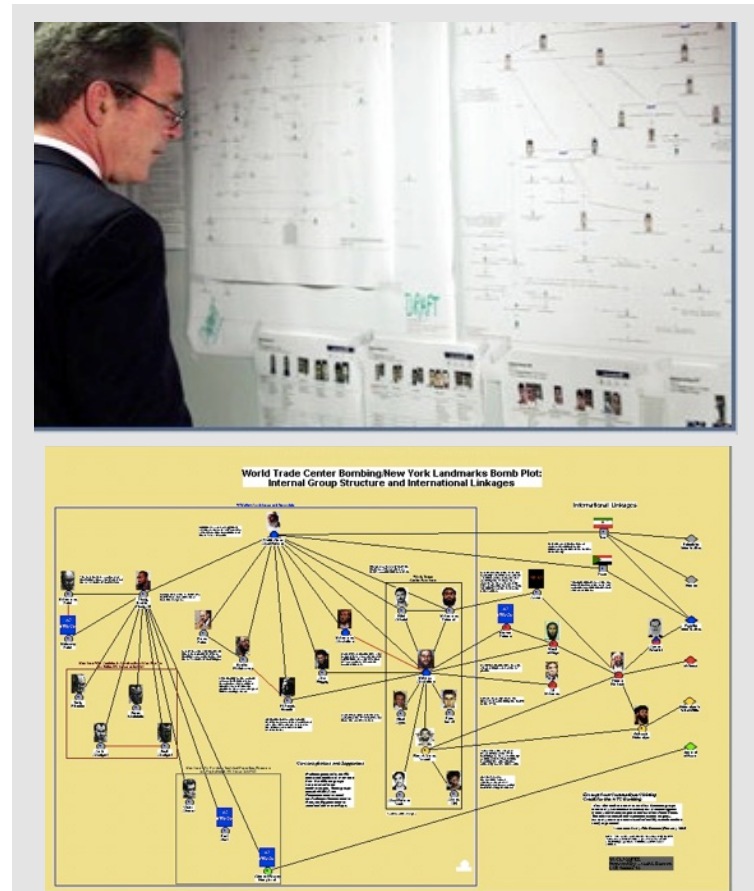


# **The Networks Grand Challenge?**



## ***“our real adversaries are networks...”***

- Many national security threats come from loose, dynamic networks of people & organizations.
  - Facilitated by networks of finance, shipment, recruiting, smuggling, etc.
  - E.g., terrorism, proliferation, cyber, drug trafficking.
- Apprehending individuals or preempting events doesn't remove the threat.
  - Need means to discover and defeat the network.
- Individual tidbits of data look benign.
  - Only recognizable in larger context of related activities
- R&D gap in issues around scale and automation:
  - Lacking scalable methods for processing very large network graphs.
  - Need to find very faint signatures (e.g., 1013 bytes within 1013 data).
  - Batch-processing unacceptable (analysts need answers within seconds).



President Bush looks over a chart depicting Osama bin Laden's financial network; an excerpt from the chart (created using one of the leading tools, i2 Analyst's Notebook)

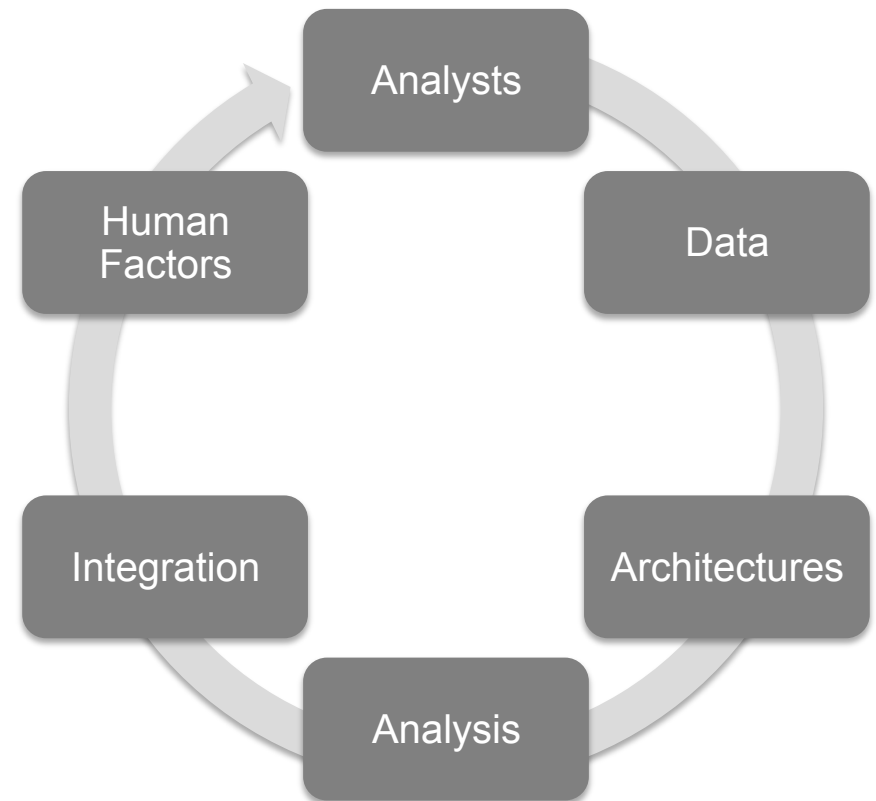
**....and so, SNL is conducting R&D to yield a radical improvement on analytical methods and tools.**





# NGC Gathers Expertise from across the Labs, working on an End-to-end Approach

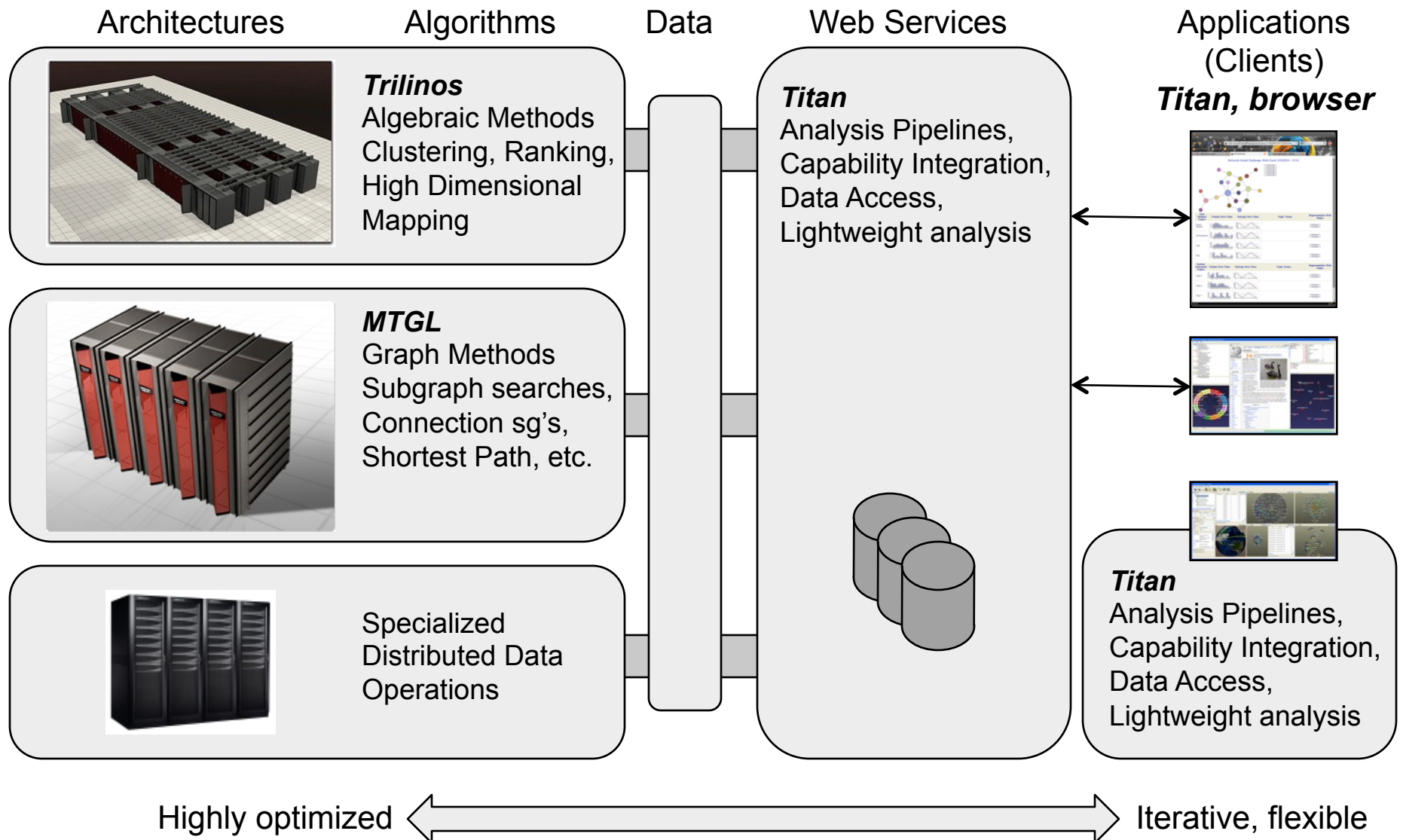
- **Data**
  - Provides data sets to the project, creates and supports the data infrastructure required by the project.
- **Architectures**
  - Develops software to efficiently integrate specialized hardware devices required by the research teams
- **Analysis (Discovery and Forecasting Teams)**
  - Develops capabilities relevant to the needle-in-a-haystack kinds of problems that concern intelligence analysts.
- **Integration**
  - Integrates NGC technologies and techniques into tools that are usable by analysts
- **Human Factors**
  - Responsible for all elicitation and knowledge representations, as well as software evaluation and assessment of technology impact on analyst performance.
  - Team of social and computer scientists interested in the relationship between human beings and computer technologies.





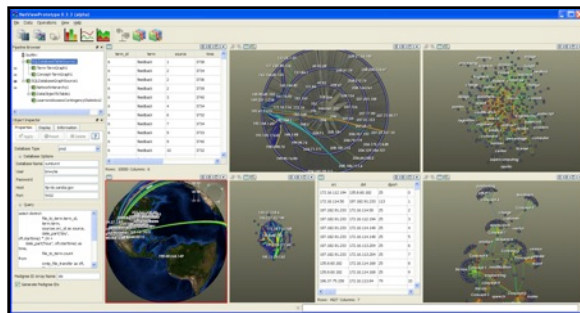
# NGC System Diagram

*"This project seeks to bring these two strengths – a solid reputation for excellence in computing, and our niche expertise in specific classes of intelligence analysis – to bear on a thorny problem: developing advanced informatics capabilities that are both usable and useful to analysts who are drowning in data." NGC project proposal*



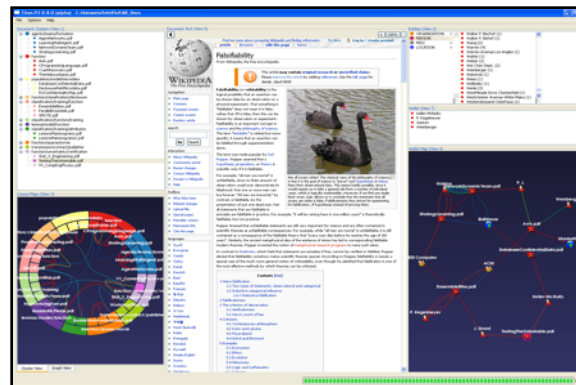


# NGC's Commitment to Prototypes Promotes End-to-end Integration



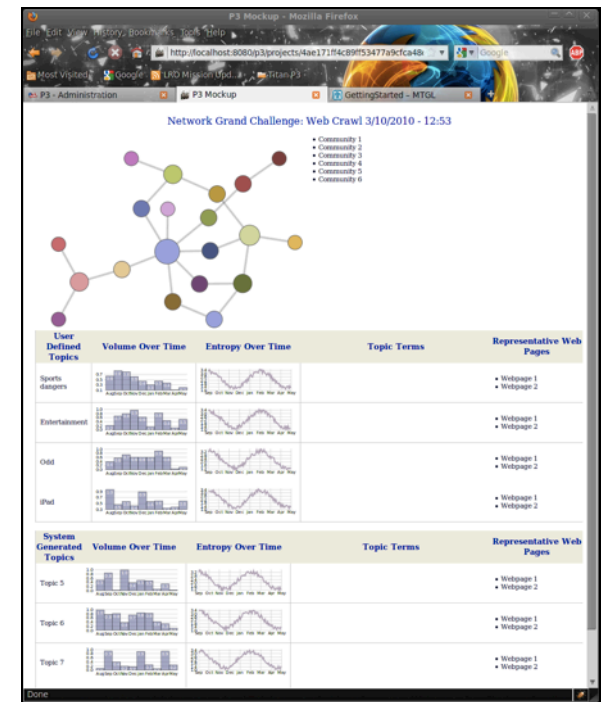
## PI: Cyber Application

- Capability integration and demonstration



## PII: Document Analysis

- Targeted development with analysts
- Iterative development
- Currently under consideration for funding



## PIII: Web 'prediction'

- New approach: web services architecture, lightweight application in browser



# There are Specific Reasons to Use HPC

- Iterative questioning in *Analyst Time*
- “Firehose, Stopwatch, or Dump truck” problems
  - Data Constraints
  - Time Constraints
  - Complexity of the query





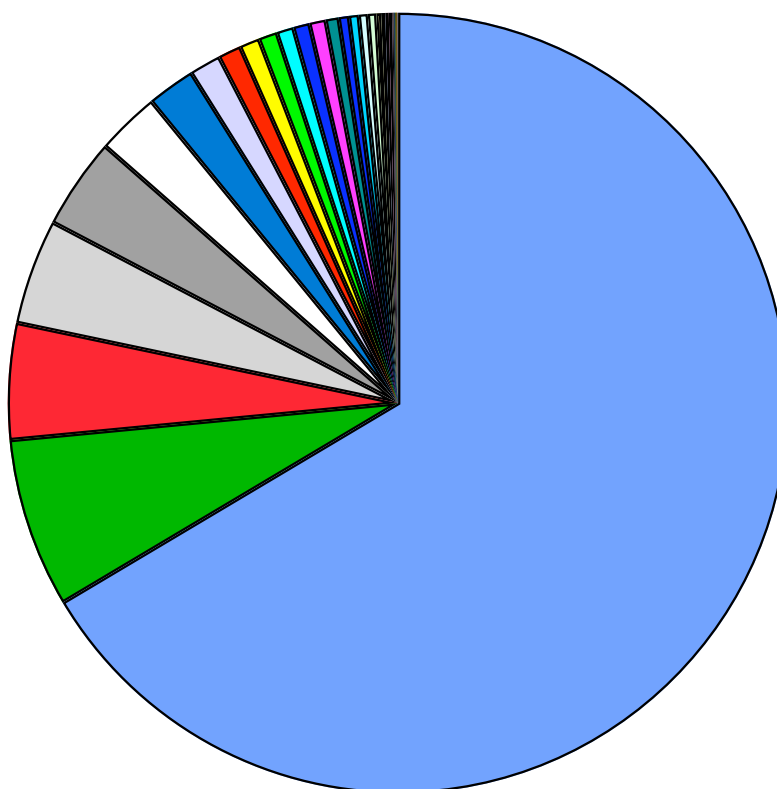
# **NGC research thrust: Multilingual Text Analysis**



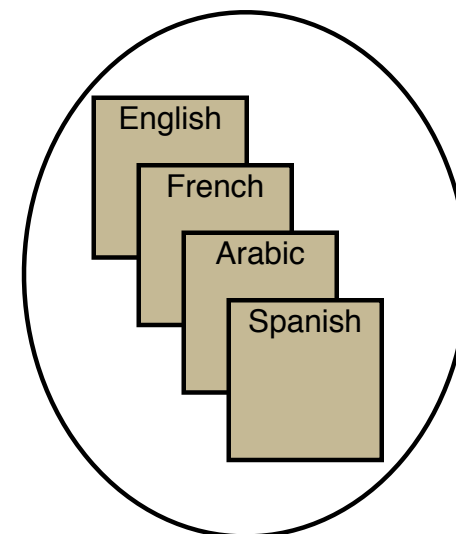
# Cross-language Information Retrieval (CLIR)

Documents could be in any language

Example: languages on the web



Goal: Cluster documents by topic regardless of language



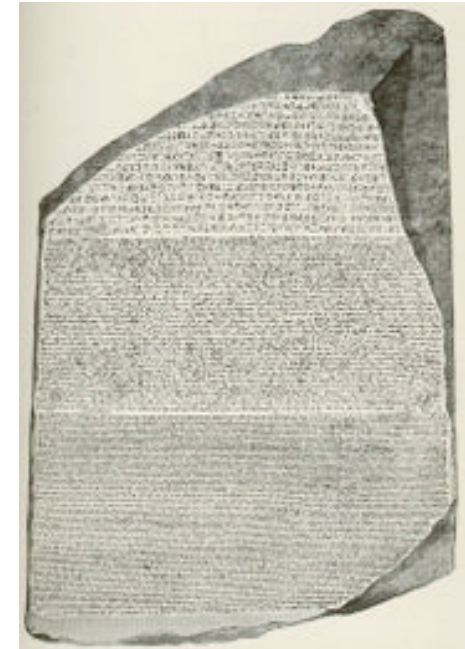
- Translation triage
- Multilingual sentiment analysis
- Ideological classification

# Bible as a 'Rosetta Stone'

- The Bible has been translated carefully and widely
  - 451 complete & 2479 partial translations
- Verse aligned

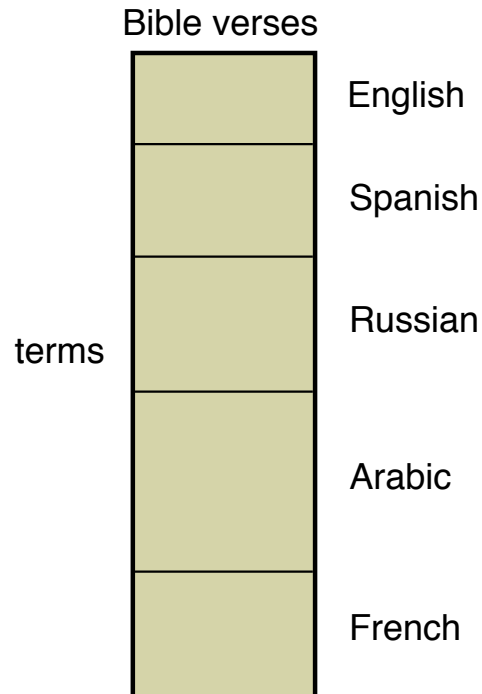
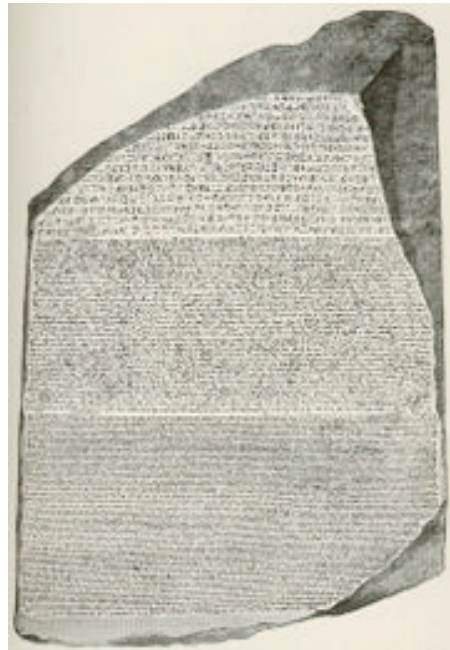
Sandia's database: 54 languages: 99.76 % coverage of web

Afrikaans	Estonian	Norwegian
Albanian	Finnish	Persian (Farsi)
Amharic	French	Polish
Arabic	German	Portuguese
Aramaic	Greek (New Testament)	Romani
Armenian Eastern	Greek (Modern)	Romanian
Armenian Western	Hebrew (Old Testament)	Russian
Basque	Hebrew (Modern)	Scots Gaelic
Breton	Hungarian	Spanish
Chamorro	Indonesian	Swahili
Chinese (Simplified)	Italian	Swedish
Chinese (Traditional)	Japanese	Tagalog
Croatian	Korean	Thai
Czech	Latin	Turkish
Danish	Latvian	Ukrainian
Dutch	Lithuanian	Vietnamese
English	Manx Gaelic	Wolof
Esperanto	Maori	Xhosa



# Term-Doc Matrix

Term-by-verse matrix  
for all languages



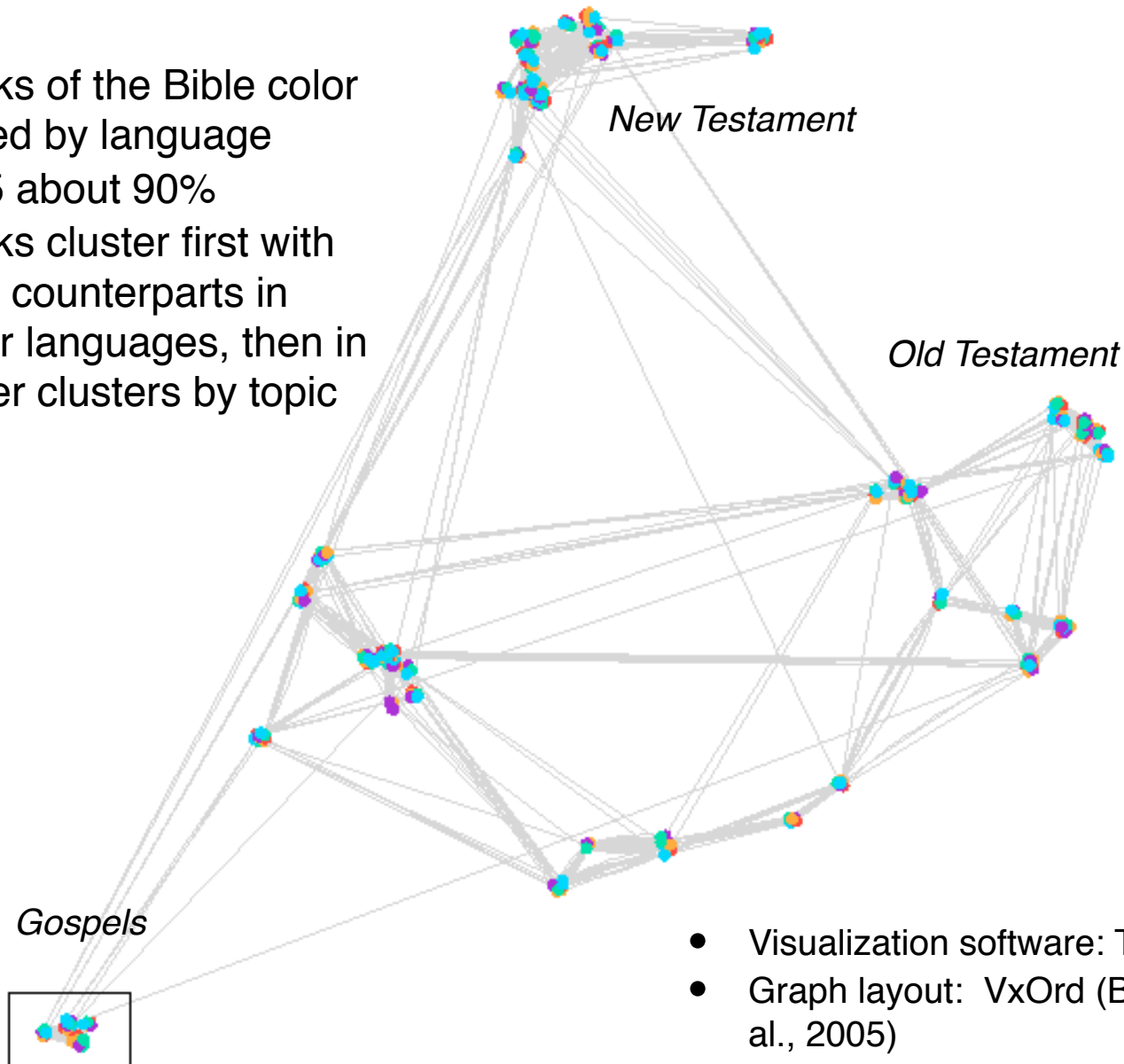
163,745 x 31,230

Look for co-occurrence of  
terms in the same verses  
and across languages to  
capture latent concepts

- Approach is not new: pairs of languages in Latent Semantic Analysis (LSA)
  - English and French (Landauer & Littman, 1990)
  - English and Greek (Young, 1994)
- *Multi-parallel* corpus is new

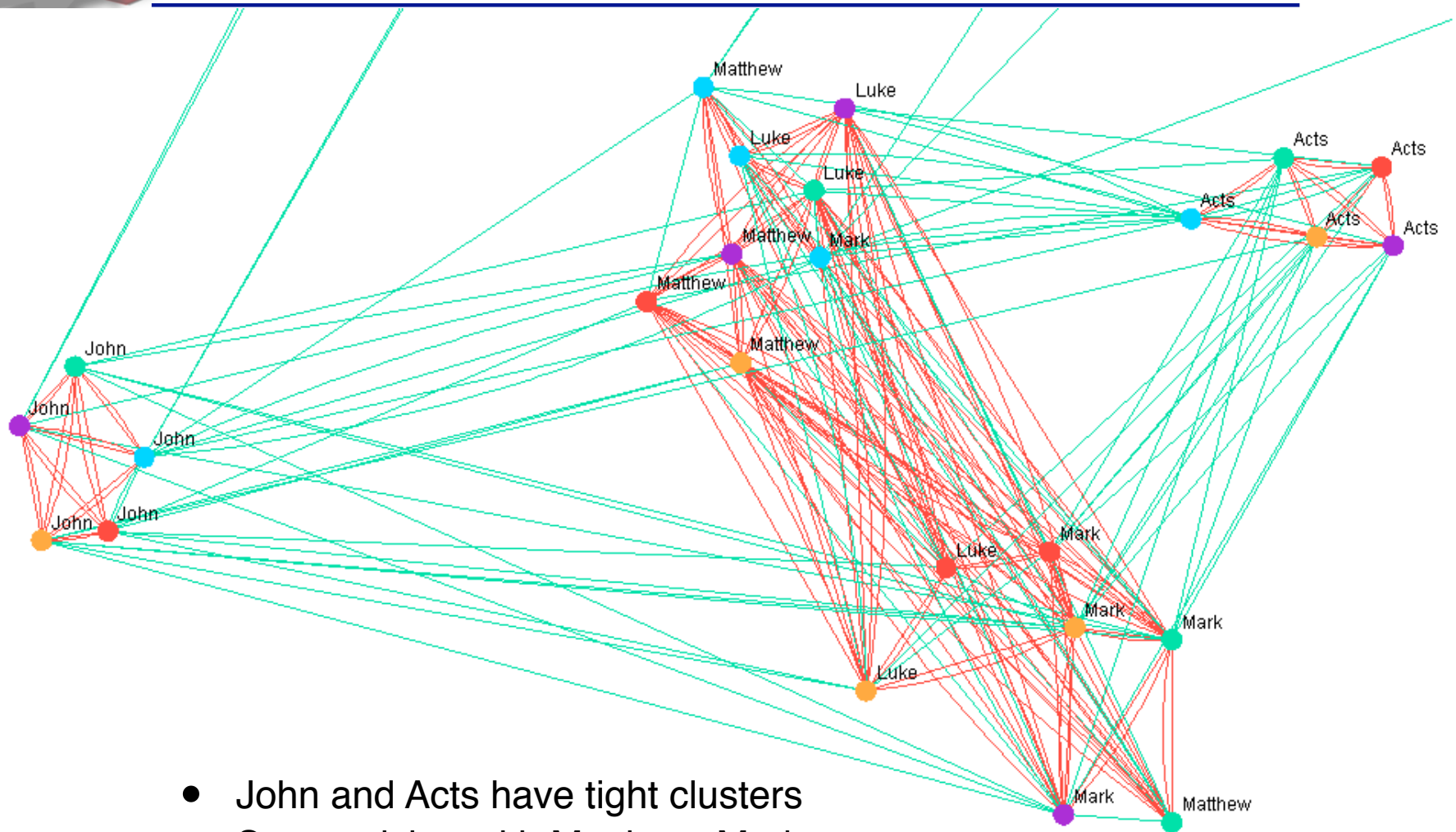
# Bible Clustering with LMSATA

- Books of the Bible color coded by language
- MP5 about 90%
- Books cluster first with their counterparts in other languages, then in larger clusters by topic



- Visualization software: Tamale 1.2
- Graph layout: VxOrd (Boyack et al., 2005)

# Clustering Close-up



- John and Acts have tight clusters
- Some mixing with Matthew, Mark, Luke (synoptic gospels - share a similar perspective)



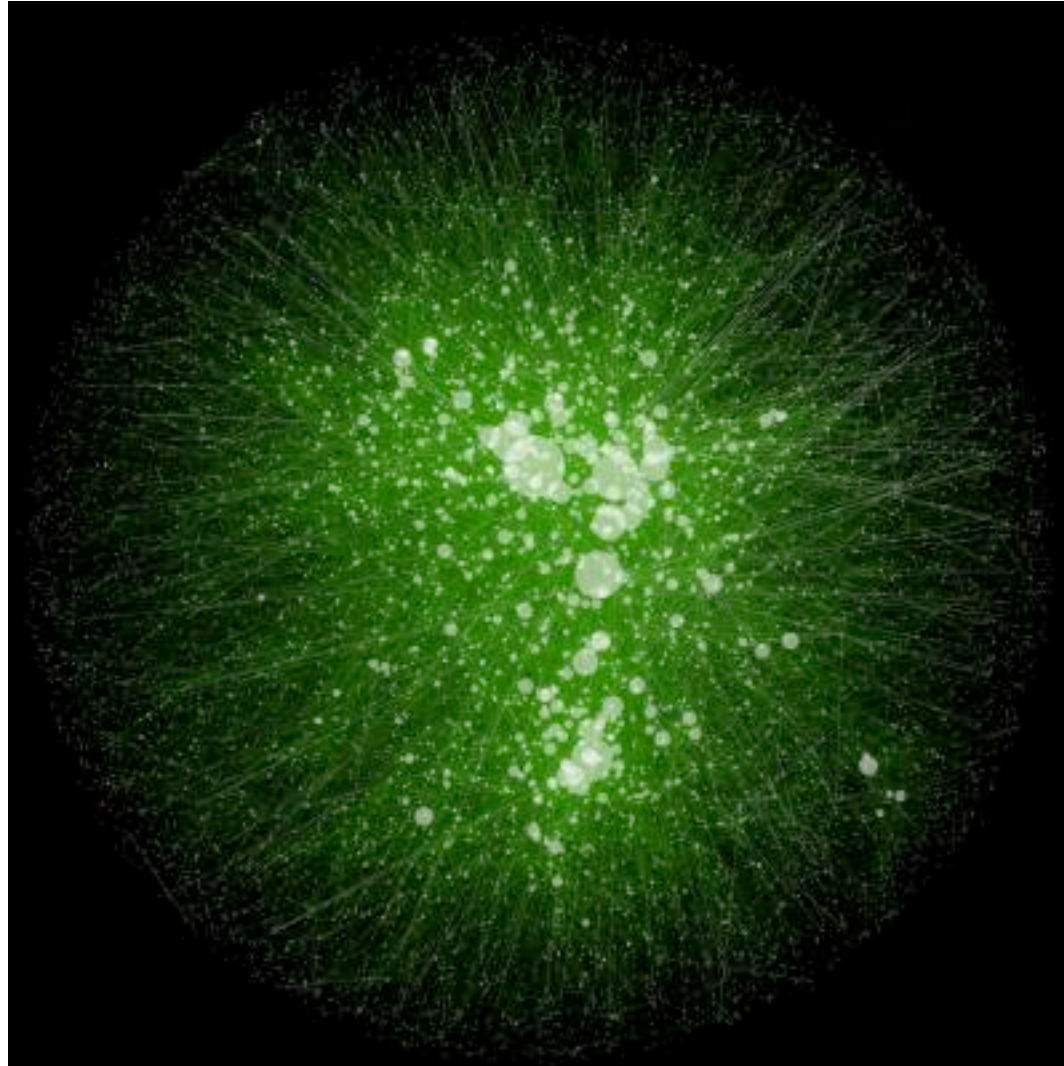


# All (Large) Data Analysis is Highly, Immediately Constrained

- Raw data must be processed
  - Example: database schema
- Analysis Algorithms op on specific data types
  - Graphs, tensors, images
- Advanced architectures support certain ops
  - XMT: multithreaded (graphs)
  - Netezza: hardware execution of SQL
  - Distributed memory: large tensor operations
- Analysis results viewed, queried many ways
  - Layout/vis introduces bias
- Human in the loop is a known unknown
  - *Especially in non-scientific arenas*

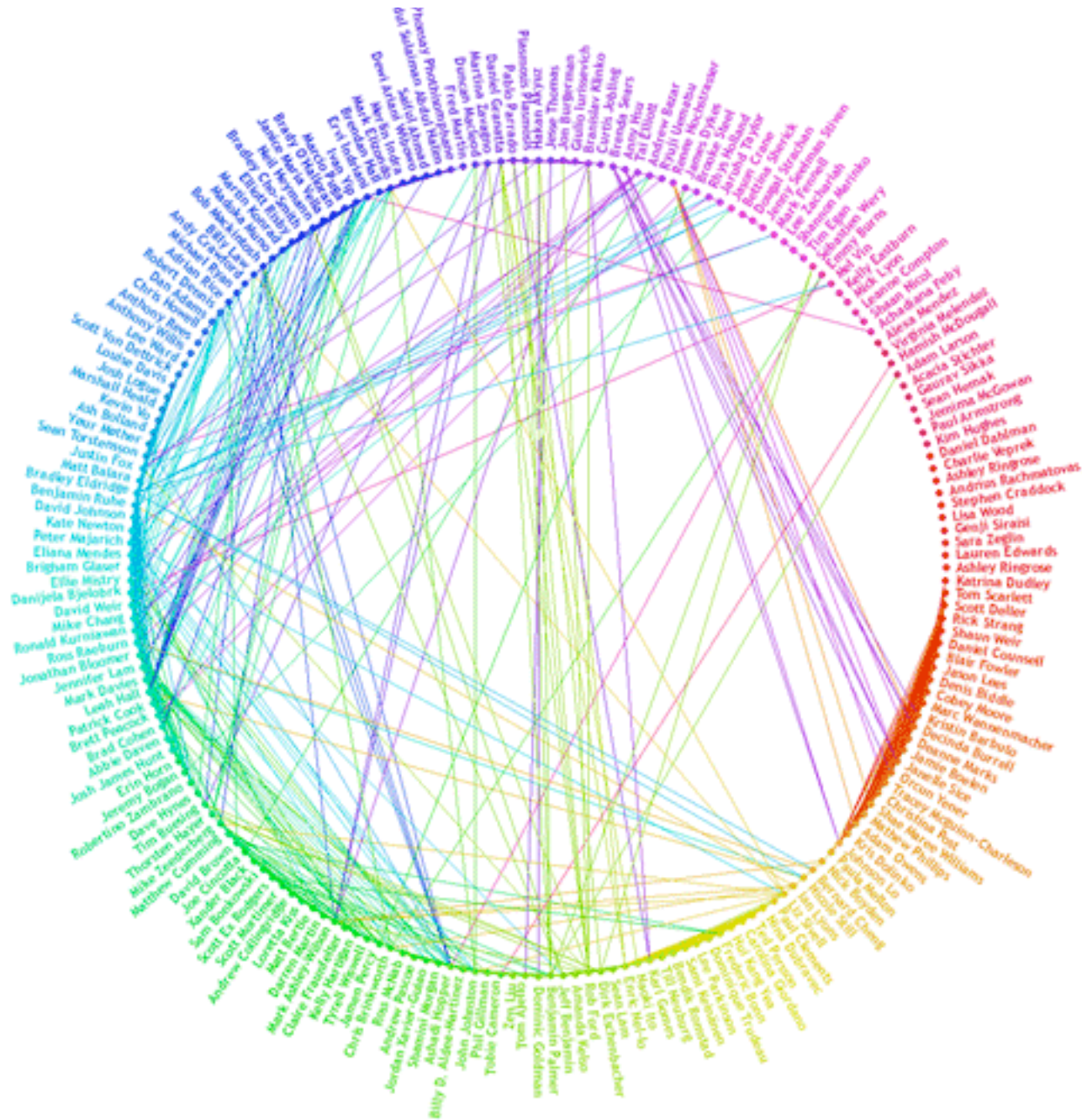


# The Hairball Principle



This image, created by Trey Ideker using Cytoscape, depicts thousands of known molecular and genetic interactions occurring inside the human body.

<http://www.physorg.com/news/2010-10-center-human-function.html>



[http://www.bannerblog.com.au/news/2008/08visualizing\\_your\\_life\\_the\\_non\\_new\\_age\\_way.php](http://www.bannerblog.com.au/news/2008/08visualizing_your_life_the_non_new_age_way.php)





Mapping of intensity of Facebook friends between pairs of cities, with the connections drawn as great arcs.

*After a few minutes of rendering, the new plot appeared, and I was a bit taken aback by what I saw. The blob had turned into a surprisingly detailed map of the world. Not only were continents visible, certain international borders were apparent as well. What really struck me, though, was knowing that the lines didn't represent coasts or rivers or political borders, but real human relationships. Each line might represent a friendship made while travelling, a family member abroad, or an old college friend pulled away by the various forces of life.*

**Paul Butler, in a blogpost:** <http://www.facebook.com/notes/facebook-engineering/visualizing-friendships/460716340019>



Sandia  
National  
Laboratories



## UX notes

- Real data is messy. *Very messy*
- Interesting data is hopeless to view in its entirety
  - This requires interactivity across CS
    - Algorithms, architectures, data movement, vis
- Understanding data is a process, not a product
- Customers must understand the algorithms, so they do not infer unsupported information from the presentation

***Elbow-to-elbow collaboration is essential to understanding large data – iterative development is essential.***

***The process and the data are messy – the software CANNOT be***



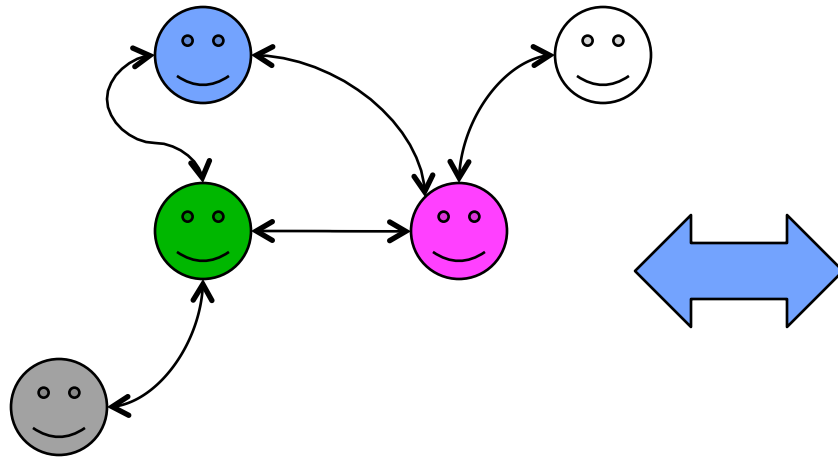


# Unifying Data Abstractions



# Semantic Graphs and Tensors are Sandia's Unifying Data Structures

*"The central hypothesis of this project is the idea that network structures extracted from large datasets can be subjected to mathematical analysis and testing to identify patterns of real-world behavior." NGC proposal*

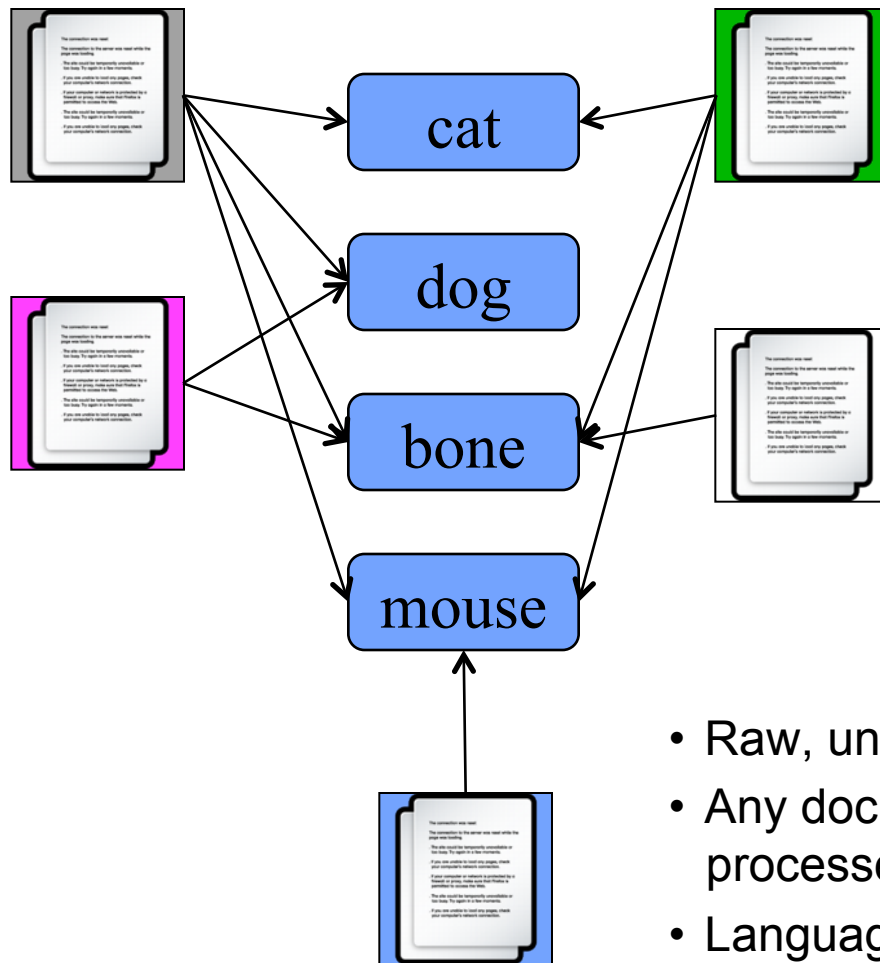


		1			
	1		1	1	
		1		1	
		1	1		1
				1	

- Networks can be of many types
  - Social network
  - Cyber traffic
  - Communications
- Graph and matrix/tensor are equivalent representations
  - Extends to multiple dimensions
- Data can be easily transformed
- Both graph and algebraic algorithms can be run on the same data
  - appropriate architectures



# Text is Easily Transformed into Graph and Matrix Data Structures

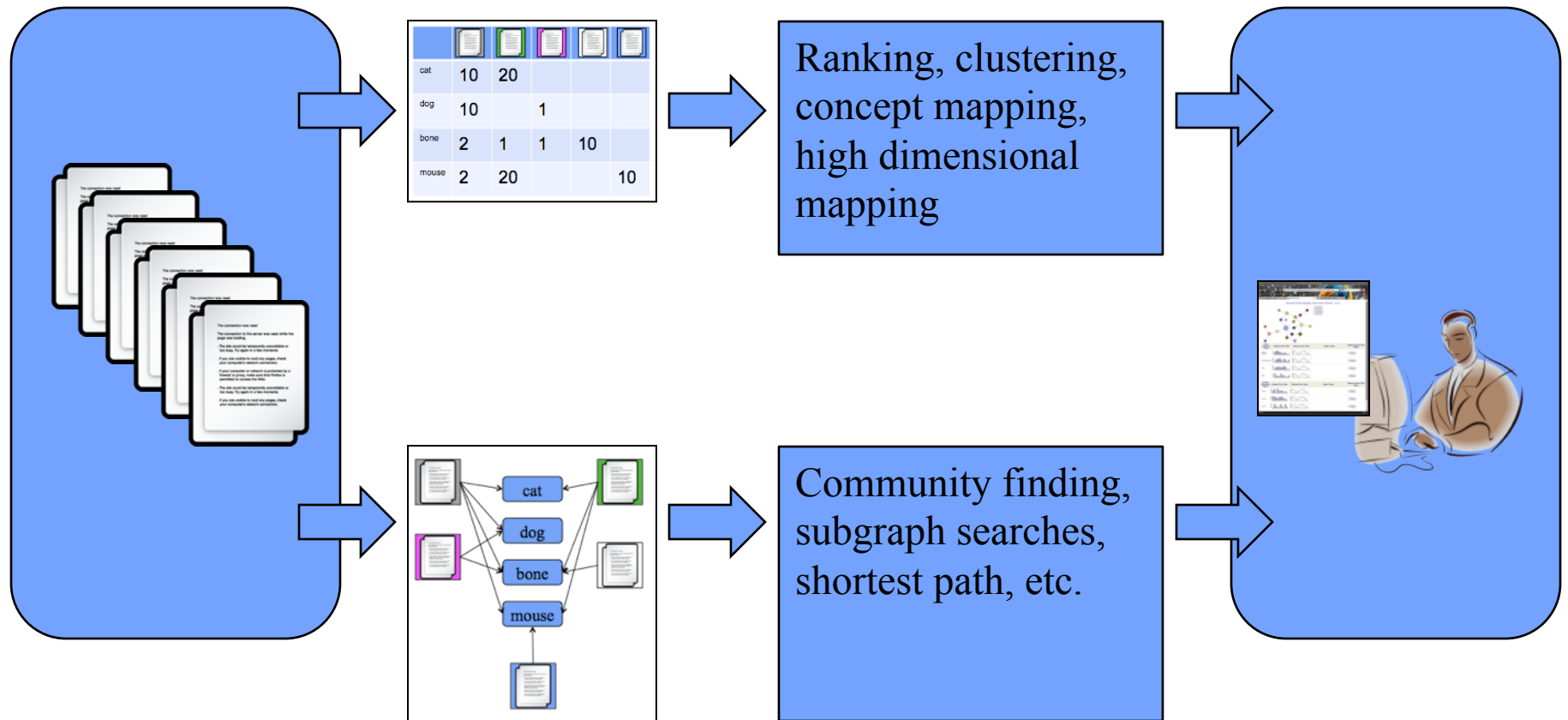


cat	10	20			
dog	10		1		
bone	2	1	1	10	
mouse	2	20			10

- Raw, unstructured text input
- Any document, in any language, can be processed into these data structures
- Language expertise is embedded in techniques
  - No language expertise required to use them



# ‘Connecting the Dots’ Benefits from a Range of Approaches



- Supports rich relationship-centered analysis
- Combines large, heterogeneous data corpora
- Different abstractions support different analytics



# PII: Integrating Algebraic and Graph Methods

Titan.P2 0.8.0 (alpha) - C:/datasets/InfoVis/EAB\_Docs

File Options Help

Document Clusters (View 1)

- agents/teams/formation
  - AgentNetworks.pdf
  - LearningMultiAgent.pdf
  - NetworkDynamicTeam.pdf
  - StrategyLearning.pdf
- function
  - Awk.pdf
  - CProgrammingLanguage.pdf
  - CrashRecovery.pdf
  - TheHidiousName.pdf
- population/model/microdata
  - DatabaseConfidentialData.pdf
  - DisclosureRiskMicrodata.pdf
  - RUConfidentialityMap.pdf
- function/classification/disclosure
  - EnsembleBites.pdf
  - ParallelEnsemble.pdf
  - SMOTE.pdf
- classification/function/training
  - lemma/model/function
  - classification/training/attributes
    - LesionsMammograms.pdf
    - LesionsMammograms2.pdf
- function/sparse/movie
  - transmission/contact/publisher
  - function/uncertainty/verification
    - Stat\_V\_Engineering.pdf
    - TestingTheUnintestable.pdf
    - VV\_CompEngPhysics.pdf

Document Text (View 5)

Find out more about navigating Wikipedia and finding information. Try Beta Log in / create account

article discussion edit this page history

## Falsifiability

From Wikipedia, the free encyclopedia

This article **may contain original research or unverified claims.** Please improve the article by adding references. See the talk page for details. (April 2009)

**Falsifiability** (or **refutability**) is the logical possibility that an assertion can be shown false by an observation or a physical experiment. That something is "falsifiable" does not mean it is false; rather, that *if* it is false, then this can be shown by observation or experiment. Falsifiability is an important concept in science and the philosophy of science. The term "testability" is related but more specific: it means that an assertion can be falsified through experimentation alone.

The term was made popular by Karl Popper. Popper asserted that a hypothesis, proposition, or theory is scientific only if it is falsifiable.

For example, "all men are mortal" is unfalsifiable, since no finite amount of observation could ever demonstrate its falsehood: that one or more men can live forever. "All men are immortal," by contrast, is falsifiable, by the presentation of just one dead man. Not all statements that are falsifiable in principle are falsifiable in practice. For example, "it will be raining here in one million years" is theoretically falsifiable, but not practical.

Popper stressed that unfalsifiable statements are still very important for science and are often contained in scientific theories as unfalsifiable consequences. For example, while "all men are mortal" is unfalsifiable, it is still contained as a consequence of the falsifiable theory that "every man dies before he reaches the age of 150 years". Similarly, the ancient metaphysical idea of the existence of atoms has led to corresponding falsifiable modern theories. Popper invented the notion of **metaphysical research programs** to name such ideas.

In contrast to **Positivism**, which held that statements are senseless if they cannot be verified or falsified, Popper denied that falsifiability somehow makes scientific theories special. According to Popper, falsifiability is merely a special case of the much more general notion of criticizability, even though he admitted that falsification is one of the most effective methods by which theories can be criticized.

**Contents** [hide]

- Naïve falsification
  - Two types of statements: observational and categorical
  - Inductive categorical inference
    - Deductive falsification
- Falsificationism
- The criterion of demarcation
  - Verificationism
  - Use in courts of law
- Criticisms
  - Contemporary philosophers
  - Kuhn and Lakatos
  - Feyerabend
  - Sokal and Bricmont
- Examples
  - Economics
  - Ethics
  - Evolution
  - Historicism
  - Logic and mathematics
  - Religion

Entities (View 2)

- ORGANIZATION
  - Walter F. Bischof (1)
  - Walter F. Bischof (1)
- PERSON
  - Wang (2)
  - Warren (4)
  - Warren Avenue Los Angeles (1)
  - Watkin (1)
  - Weber (2)
  - Wei Chen Dept. (2)
  - Weinberger (1)
  - Weintraub (1)
  - Weiss (1)
  - Wellesley (1)
  - Werle (3)
  - WestMeade Drive Chesterfield (1)
  - Westchester Avenue White Plains (1)
  - Westendorp and Osterhaus (1)
- MISC
- LOCATION

Hotlist (View 7)

- Keller-McNulty
- P. Kegelmeyer
- Gaston
- Weinberger

Hotlist Map (View 4)

Corpus Maps (View 3)

Navigation

- Main page
- Contents
- Featured content
- Current events
- Random article

Search

Go Search

Interaction

- About Wikipedia
- Community portal
- Recent changes
- Contact Wikipedia
- Donate to Wikipedia
- Help

Toolbox

- What links here
- Related changes
- Upload file
- Special pages
- Printable version
- Permanent link
- Cite this page

Languages

- العربية
- Ελληνικά
- Català
- Cesky
- Dansk
- Deutsch
- Eesti
- Español
- Français
- 한국어
- Íslenska
- Italiano
- עברית
- Nederlands
- 日本語
- Norsk (bokmål)
- Polski
- Português
- Română
- Pycckий
- Simple English
- Suomi
- Svenska
- Versailles

Cluster View Graph View

100%





- ## Document Display



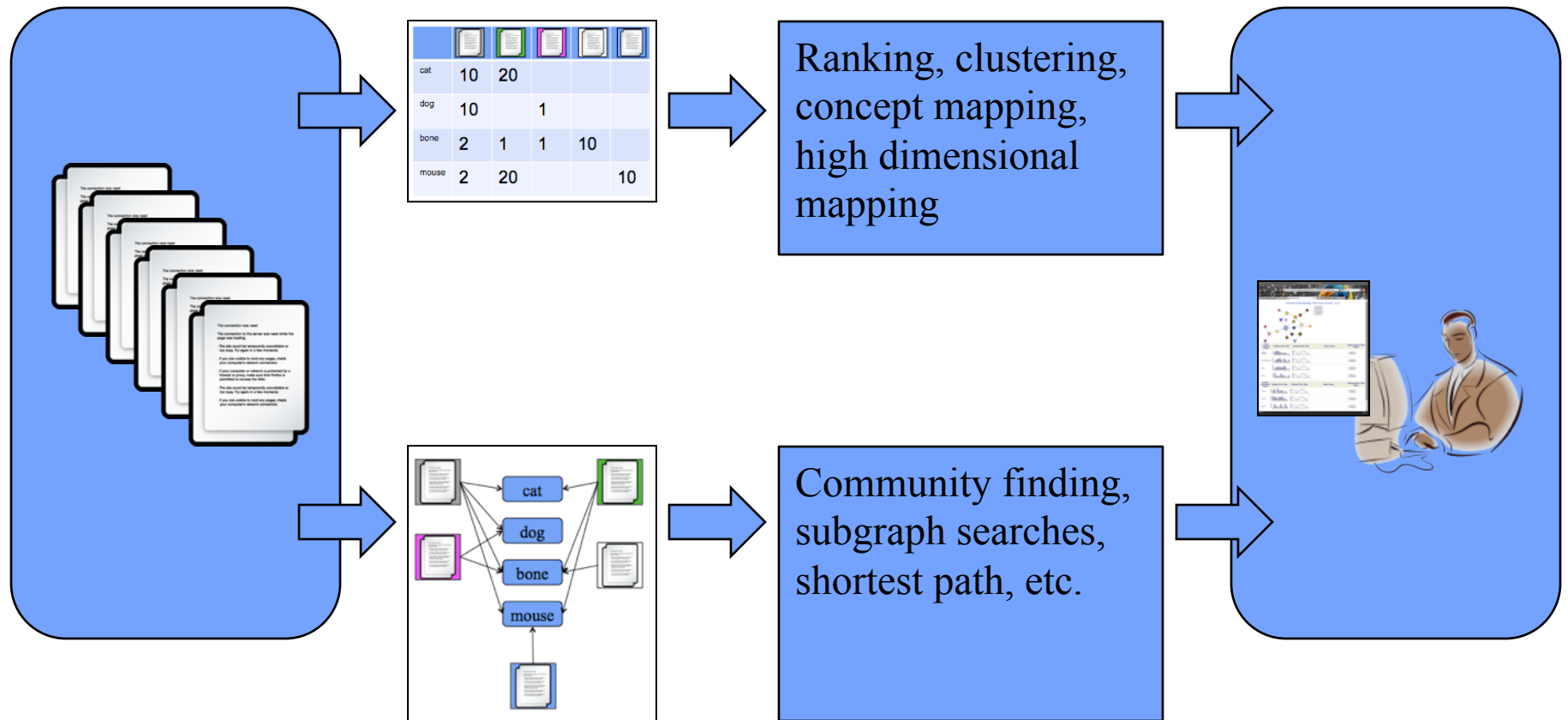
‘Soap opera’

Corpus query  
network

# Clustering



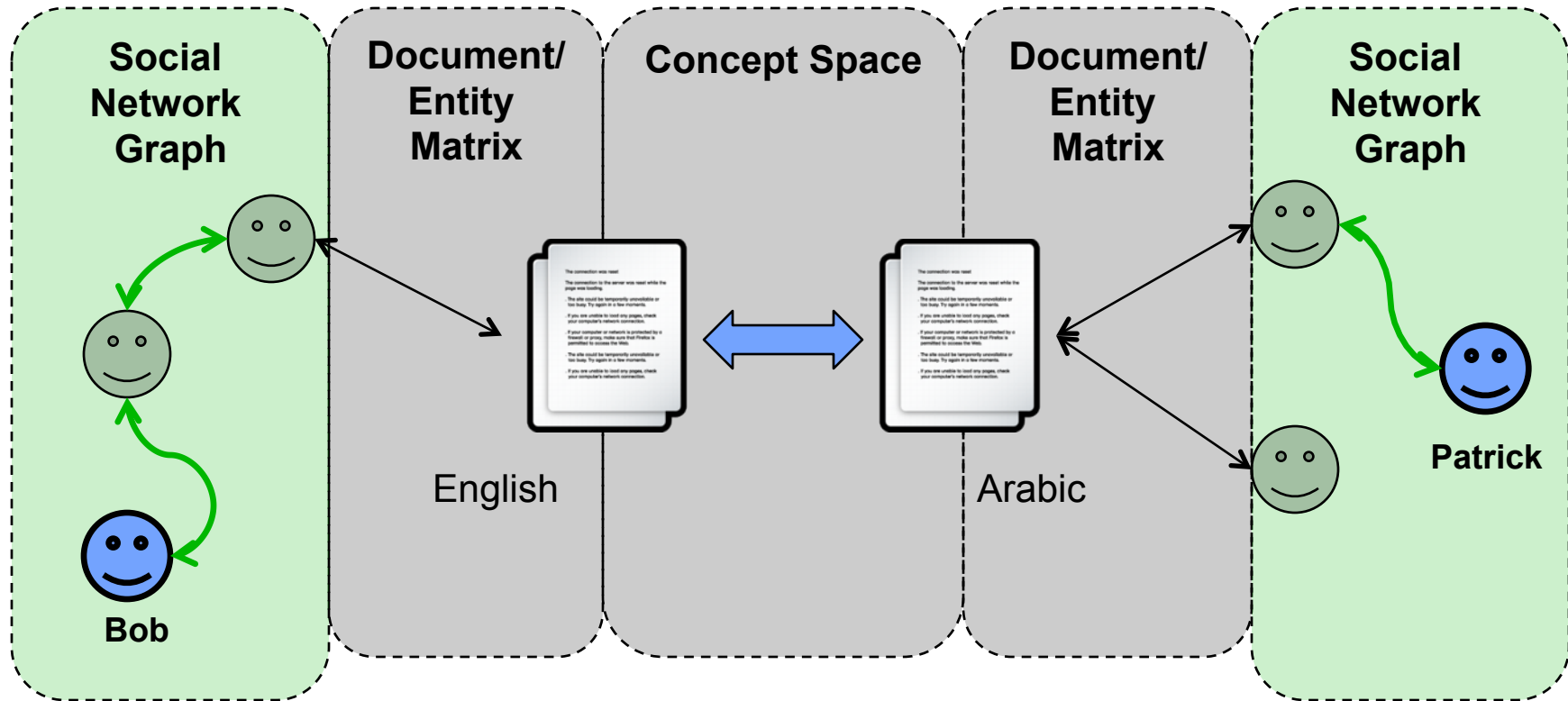
# ‘Connecting the Dots’ Benefits from a Range of Approaches



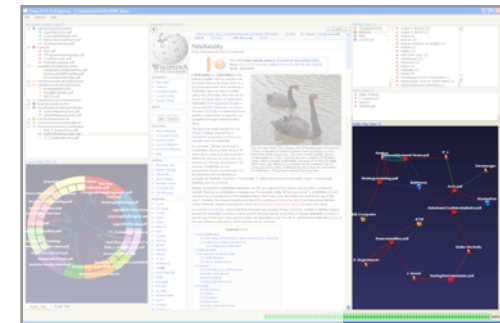
- Supports rich relationship-centered analysis
- Combines large, heterogeneous data corpora
- Different abstractions support different analytics



# Linked Graph and Algebraic Methods Provide Rich Analysis Capability



- Entities can be linked through social networks *and* concept space
- Provides rich connection data in a clear visual representation
- Promotes new paths of investigation



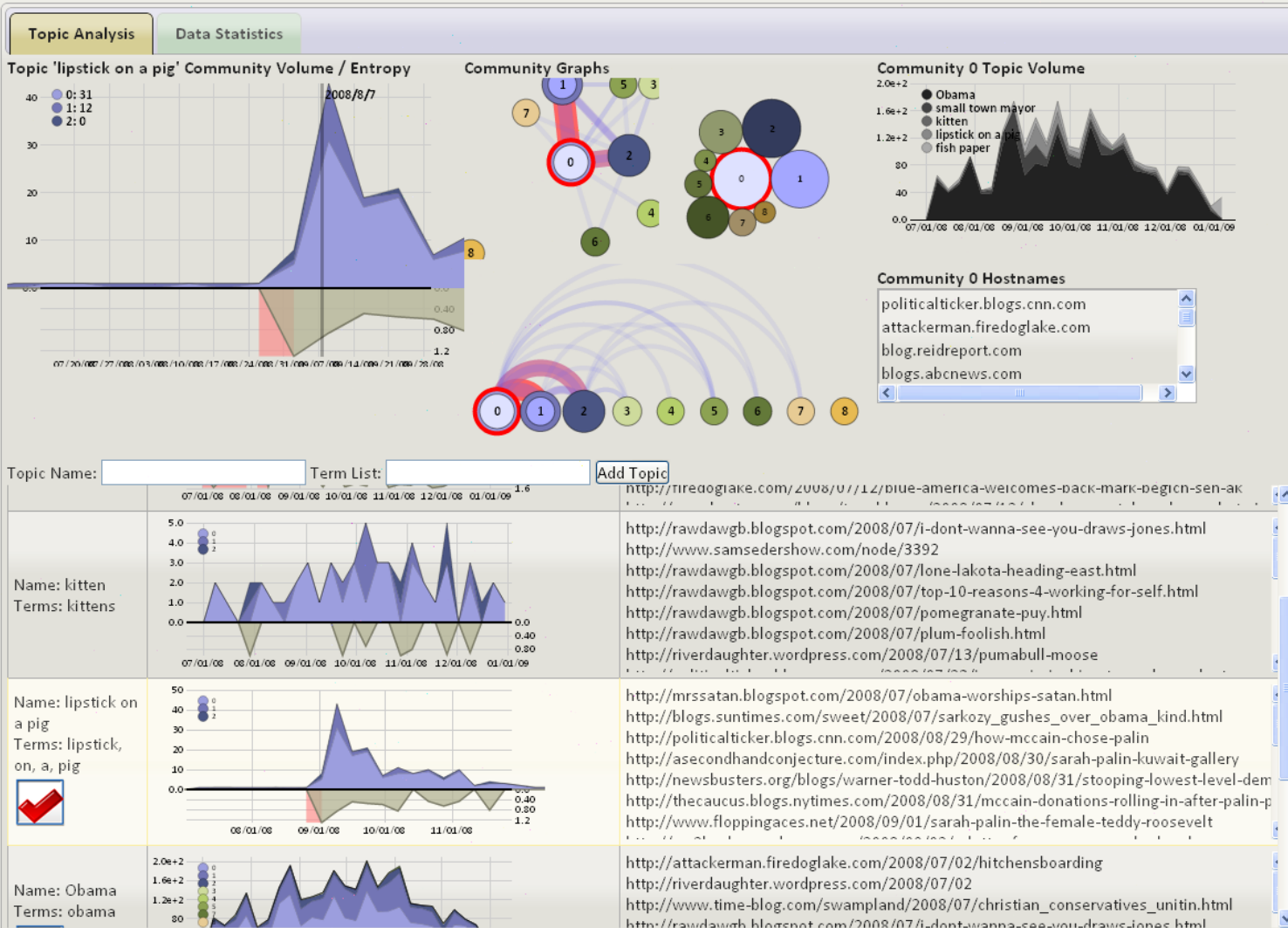
PII Prototype, emphasizing the combined graph and algebraic methods view



## UX notes

- Complexity arises from algorithms and data structures, as well as raw data
- Complexity comes from the human as well
  - Cognitive models, new intuitions, and *personal bias*

## Network Grand Challenge: Web Crawl <Known Topics>







## Themis Dashboard

*Themis was the ancient Greek Titan of good counsel ... to the ancient Greeks she was the organizer of the communal affairs of humans, particularly assemblies ...*

[Wikipedia](#)

### Users

#### Login

Status **Logged-in as tthead.**

Username

Password

Actions

#### Create User

Username

Password

Actions

#### Change Password

Username

Password

Actions

### Databases

Name

# Themis Ensemble "Automobile Performance"

## Import Variables

Delimited Text File

/Users/tshead/auto\_data.csv

Browse...

Field Delimiter

,

String Delimiter

"

Upload

## Create Model

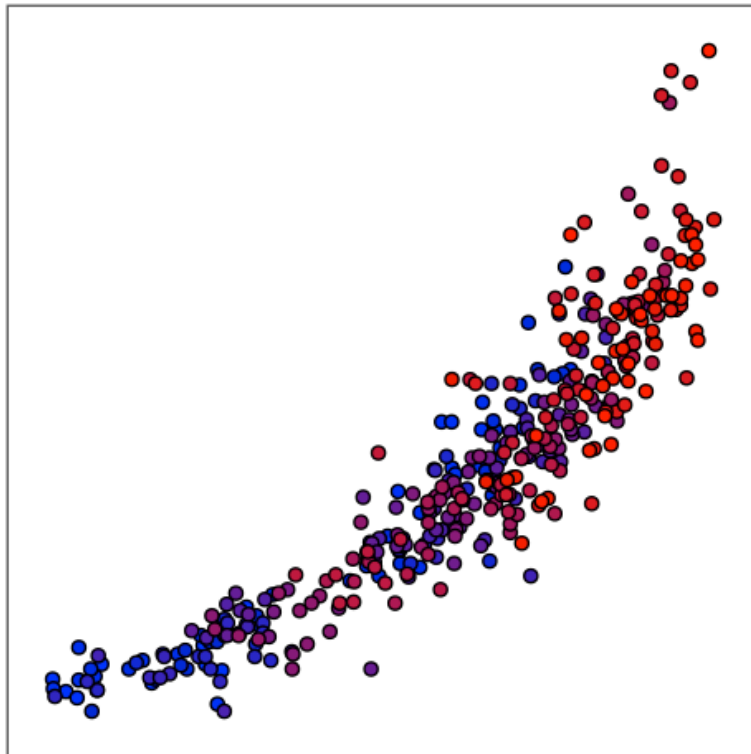
All - to - MPG

Variable	Input	Output
MPG	<input type="checkbox"/>	<input checked="" type="checkbox"/>
Cylinders	<input checked="" type="checkbox"/>	<input type="checkbox"/>
Weight	<input checked="" type="checkbox"/>	<input type="checkbox"/>
HP	<input checked="" type="checkbox"/>	<input type="checkbox"/>
Displacement	<input checked="" type="checkbox"/>	<input type="checkbox"/>
AccelTime	<input type="checkbox"/>	<input type="checkbox"/>
Model Year	<input checked="" type="checkbox"/>	<input type="checkbox"/>

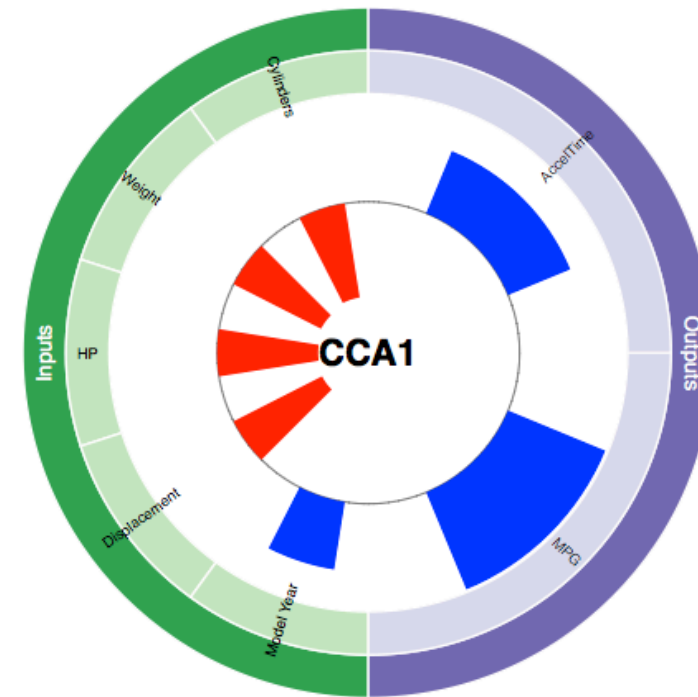
Create

[Design View](#)

## Themis Model "All - to - all"

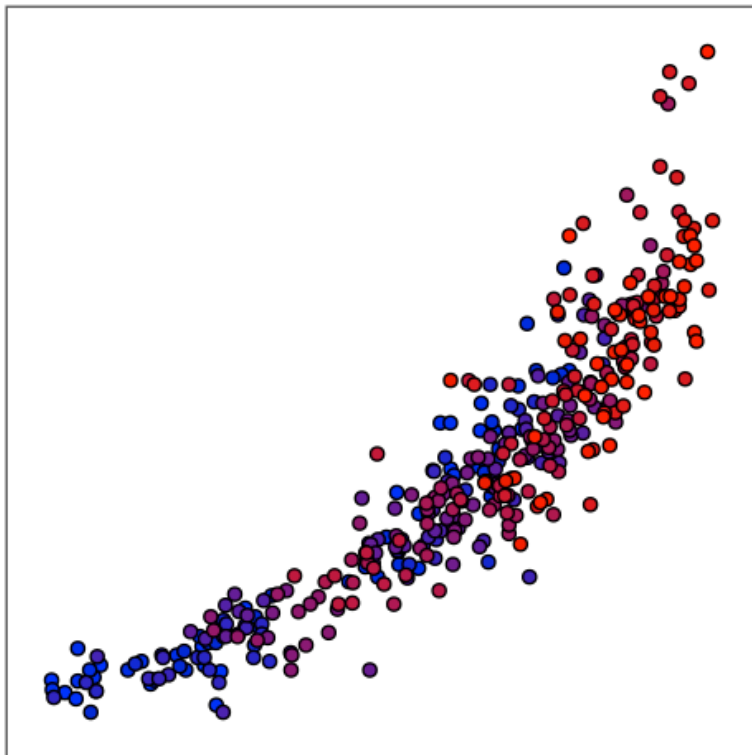


Component	R <sup>2</sup>	P
CCA1	0.910	0.089
CCA2	0.698	0.512

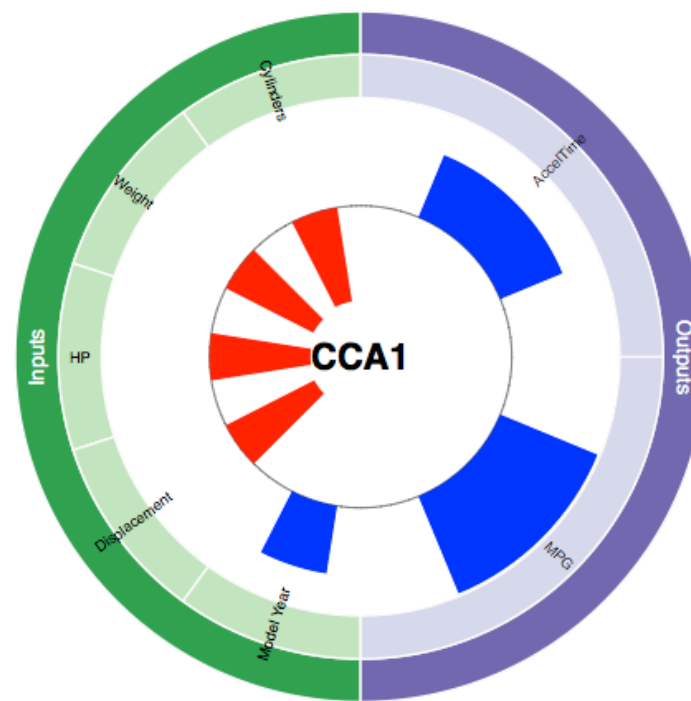


Simulation ID	Cylinders	Weight	HP	Displacement	Model Year	MPG	AccelTime
0	8.00	3.50e+3	130	307	70.0	18.0	12.0
1	8.00	3.69e+3	165	350	70.0	15.0	11.5
2	8.00	3.44e+3	150	318	70.0	18.0	11.0
3	8.00	3.43e+3	150	304	70.0	16.0	12.0

## Themis Model "All - to - all"



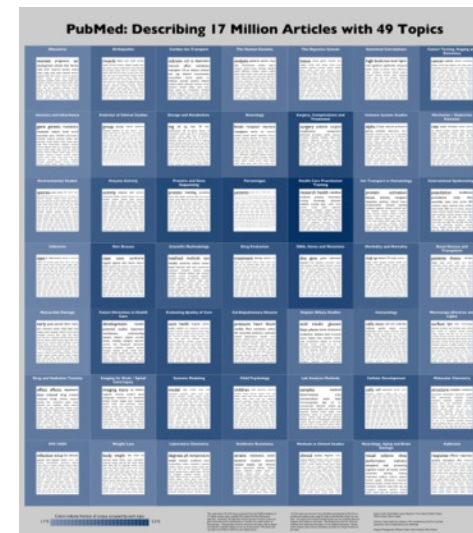
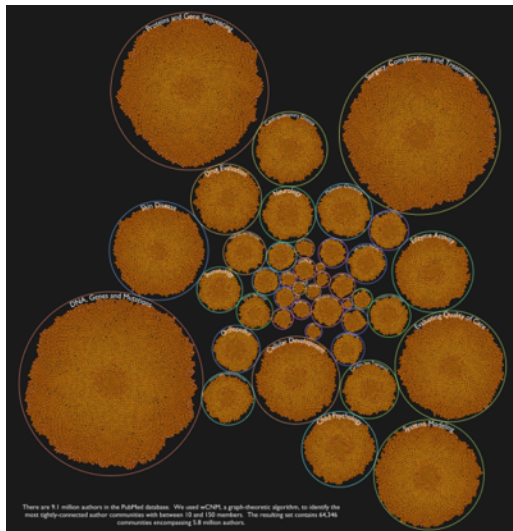
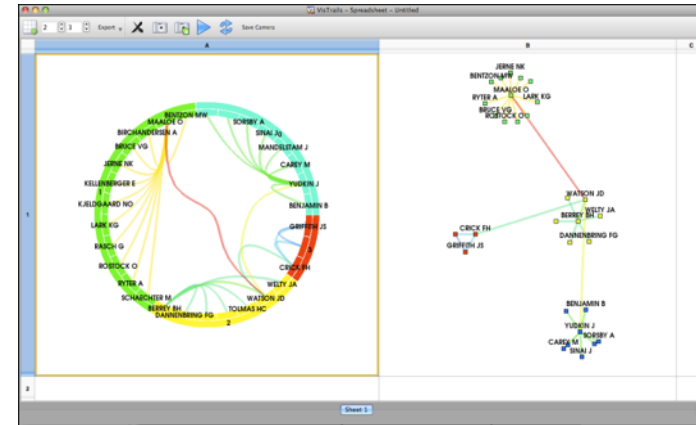
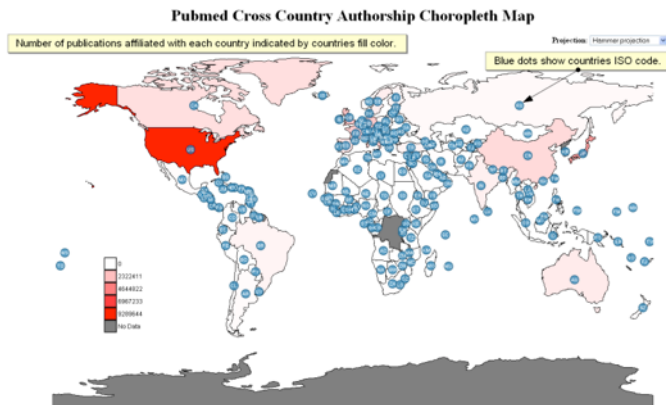
Component	$R^2$	P
CCA1	0.910	0.089
CCA2	0.698	0.512



Simulation ID	Cylinders	Weight	HP	Displacement	Model Year	MPG	AccelTime
0	8.00	3.50e+3	130	307	70.0	18.0	12.0
1	8.00	3.69e+3	165	350	70.0	15.0	11.5
2	8.00	3.44e+3	150	318	70.0	18.0	11.0
3	8.00	3.43e+3	150	304	70.0	16.0	12.0



# We have successfully clustered and visualized 17 million records.



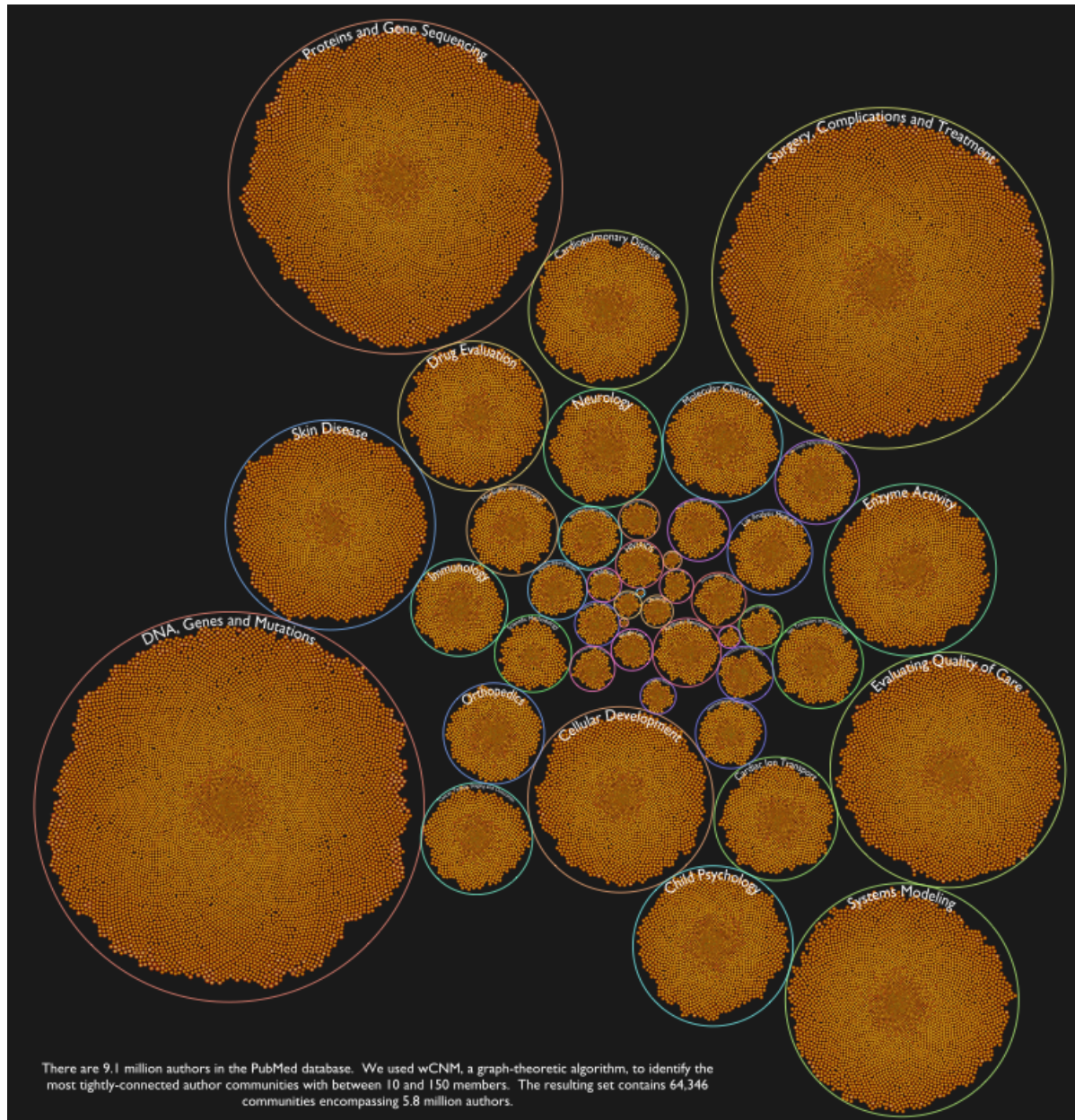
3/10/11

44





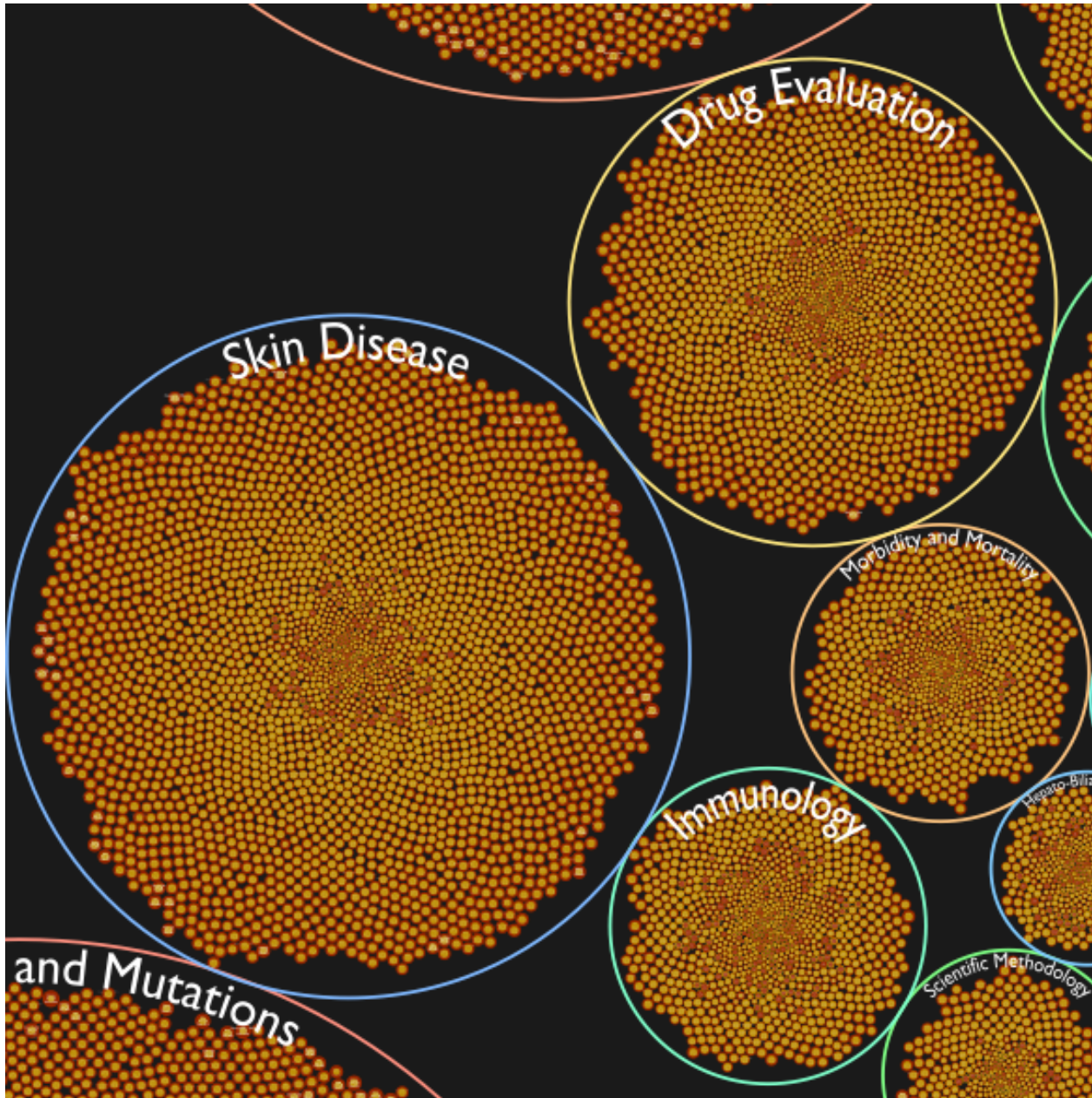
# Author Communities Visualized



- Each blob is a topic area
- Each tiny orange dot is a single community
- 44 topics, 64,314 communities
  - 10-150 authors in each community

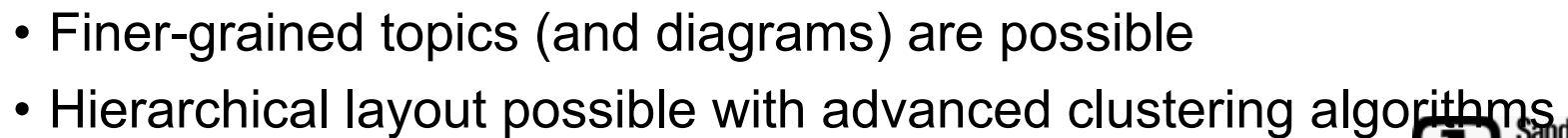


## Author Communities: Zoom 1



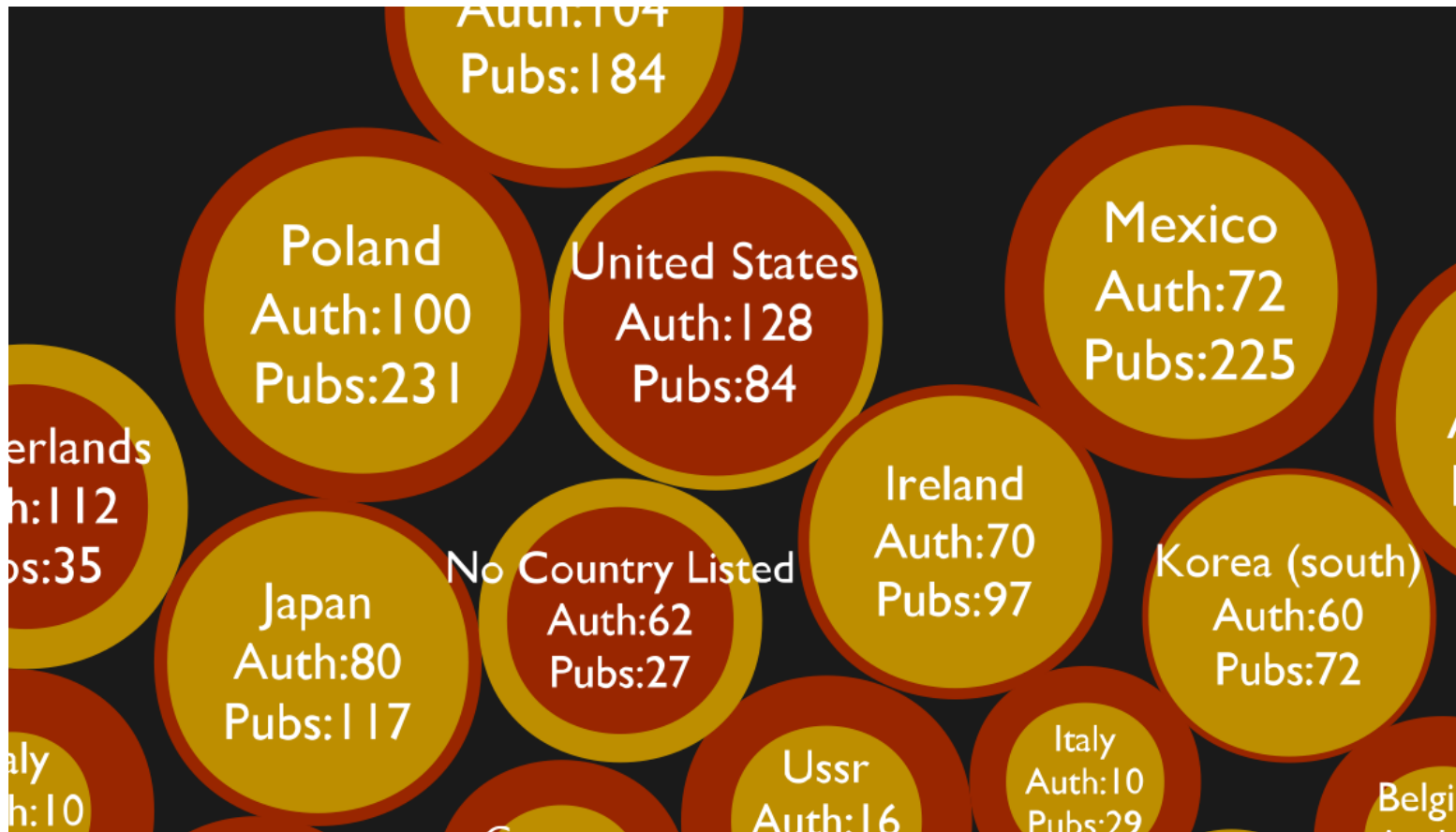
- Topic areas contain between 5 and 9800 communities
- Each community is sized according to its author/article count







## Author Communities: Zoom 3



- Community glyphs indicate relative proportion of authors to articles
- Red = articles, yellow = authors
- Red outside yellow means more articles than authors (and vice versa)
- Labels show country where most papers are published

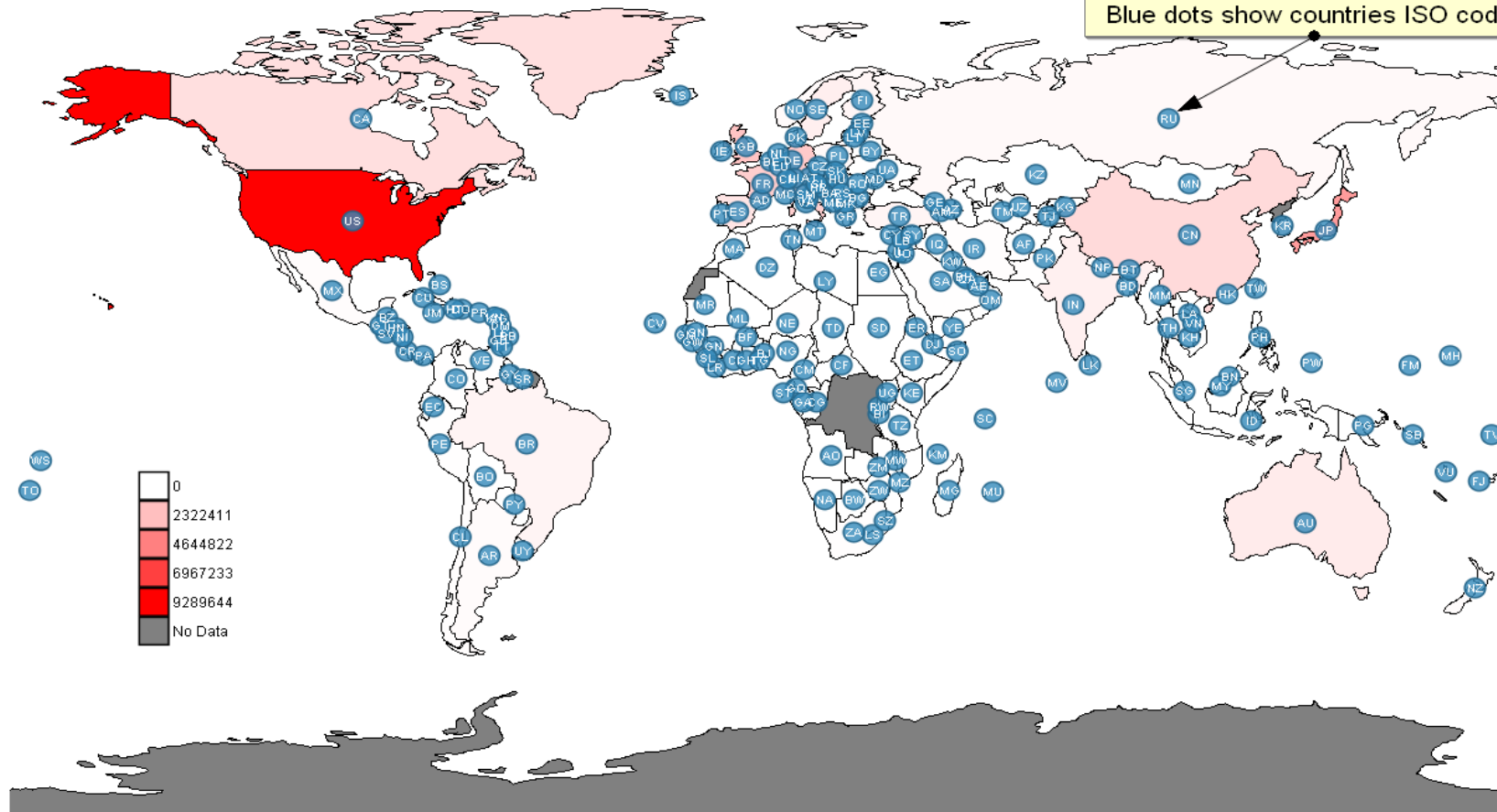


## Pubmed Cross Country Authorship Choropleth Map

Number of publications affiliated with each country indicated by countries fill color.

Projection: Hammer projection

Blue dots show countries ISO code.



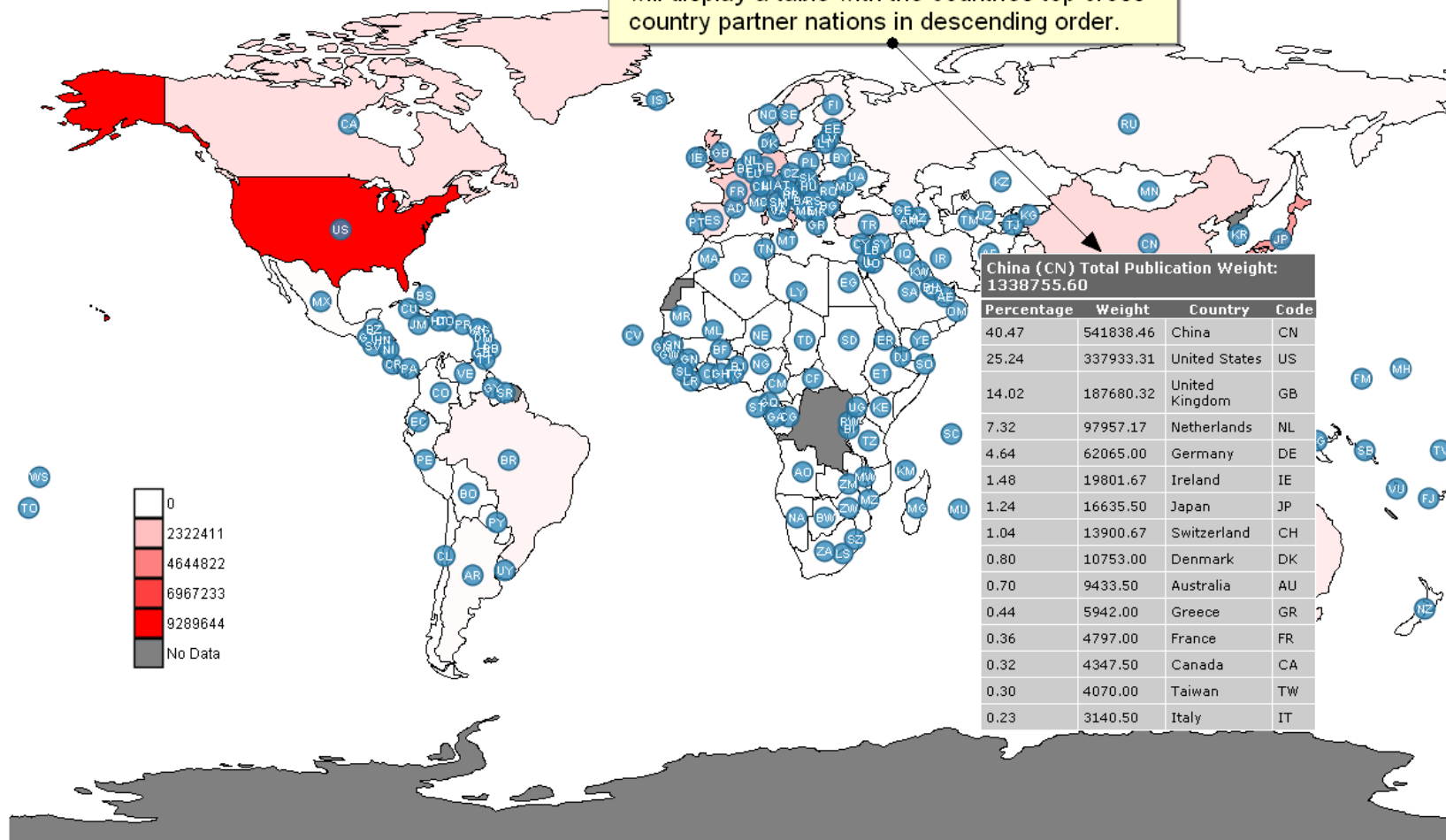




## Pubmed Cross Country Authorship Choropleth Map

Hovering the mouse over a countries blue dot will display a table with the countries top cross country partner nations in descending order.

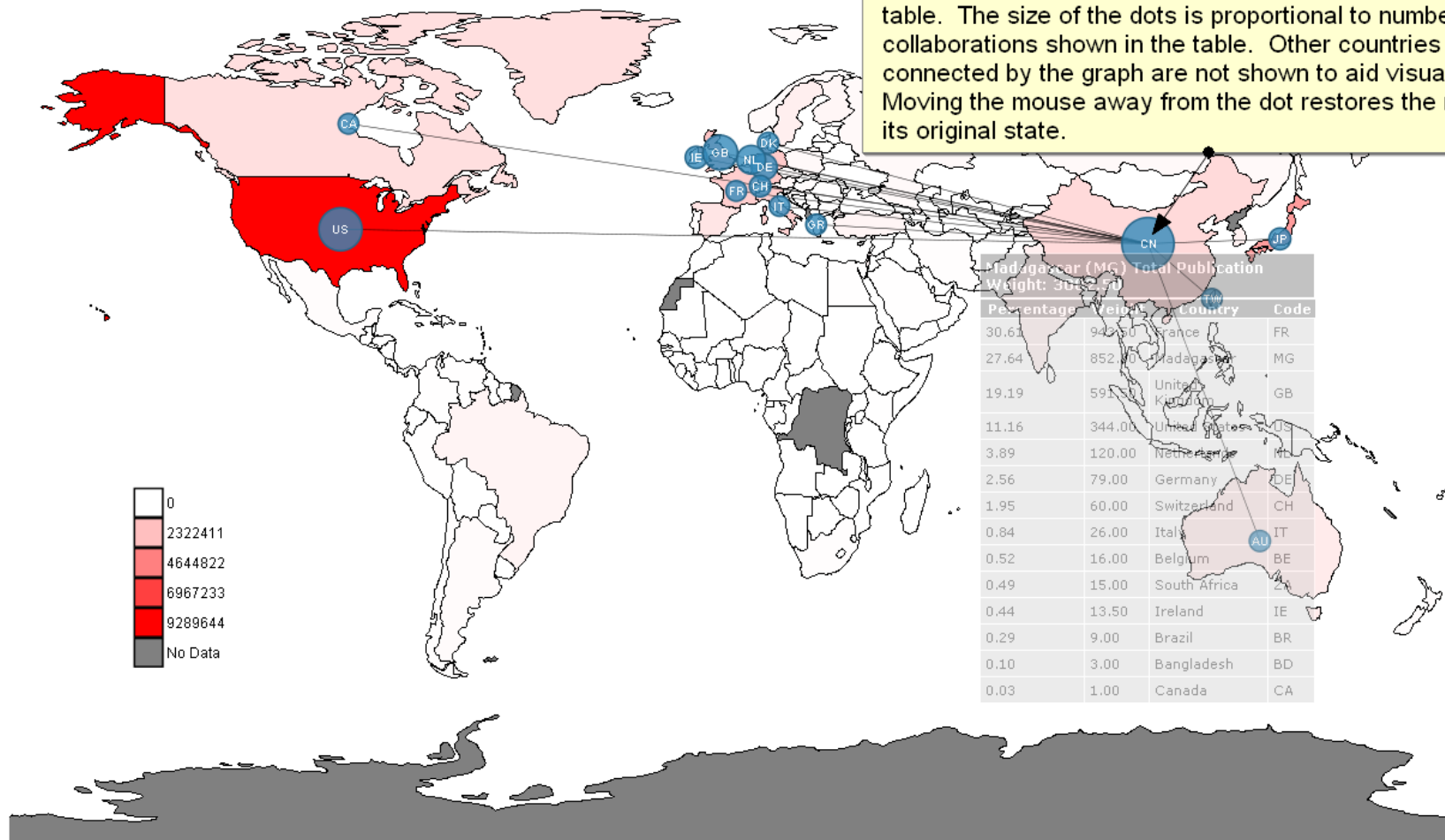
Projection: Hammer projection





## Pubmed Cross Country Authorship Choropleth Map

Clicking the mouse on a countries blue dot will produce a graph from the information contained in the displayed table. The size of the dots is proportional to number of collaborations shown in the table. Other countries not connected by the graph are not shown to aid visual clarity. Moving the mouse away from the dot restores the map to its original state.





# **All (Large) Data Analysis is Highly, Immediately Constrained**

- **Raw data must be processed**
  - Example: database schema
- **Analysis Algorithms op on specific data types**
  - Graphs, tensors, images
- **Advanced architectures support certain ops**
  - XMT: multithreaded (graphs)
  - Netezza: hardware execution of SQL
  - Distributed memory: large tensor operations
- **Analysis results viewed, queried many ways**
  - Layout/vis introduces bias
- **Human in the loop is a known unknown**
  - *Especially in non-scientific arenas*



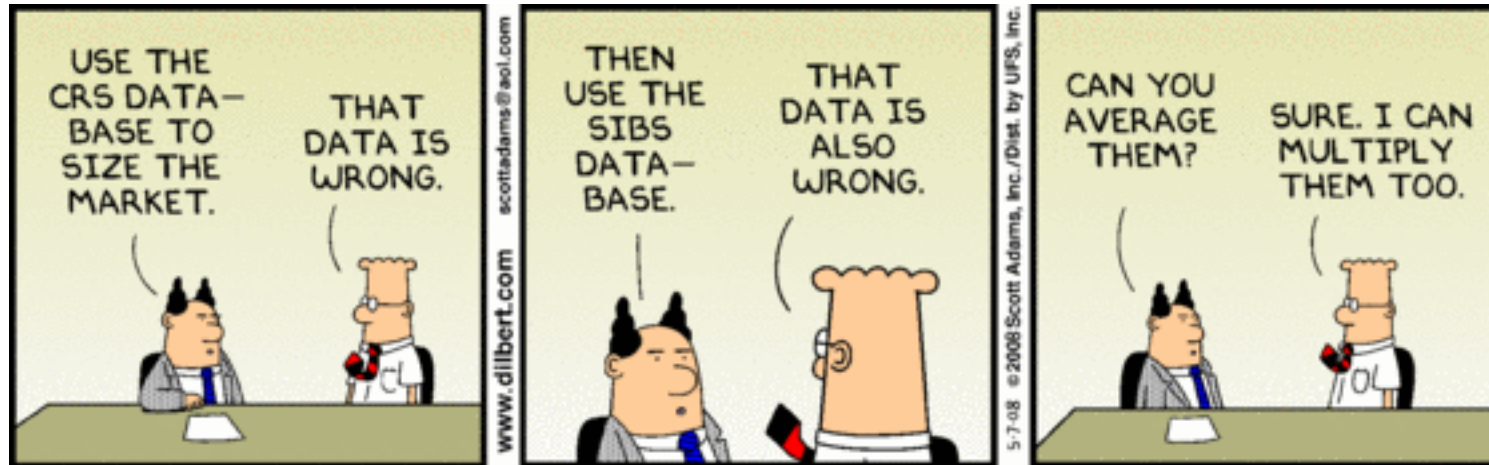
# Analysis Path Forward

- **Tight coupling among architectures, algorithms and applications**
  - Customers must understand algorithms
- ***Everything must be interactive***
  - But we must make task-specific tradeoffs
- **Good software foundations are crucial to collaboration in the field**
  - <http://titan.sandia.gov>
- **It takes interactive research to find things that are valuable enough to throw away**
  - Finding viable products is worse
- **Good news: this is so difficult and important that visualization, UI, and design are recognized as core to success**



# An important next step

Make sure your users understand the math ...



- We find that users need to understand the underlying algorithms – everything from analysis to layout (visualization)
  - Prevents incorrect conclusions from the data
- Analysis V&V needs study
  - How do we quantify the accuracy of the conclusions that are made?
- Close connection to analysts (Human Factors Team) is critical to ensure relevance
  - An untrusted algorithm is useless (literally won't get used)
  - A trusted algorithm should be understood
- This data-algorithm-visualization-analyst cycle is crucial
  - Subject of another talk





**Questions?**