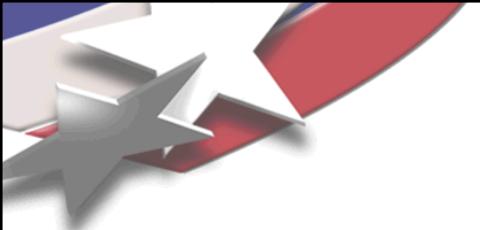# Network Algorithms for Information Analysis using the Titan Toolkit

## October 6, 2010

**William C. McLendon III**
Sandia National Laboratories

# Outline

- **Introduction to the toolkit**
  - Toolkit provisions & data pipelines

- **Applications**
  - Text Corpora Analysis
  - Fun with Web Crawls
  - Network Packet Streaming

- **Coding Example**
  - Using Python to load a graph from disk

# What is Titan?

**Data Structures**
- Table
- Tree
- DAG
- Directed Graph
- Undirected Graph
- Sparse N-way Array
- Dense N-way Array
- Unicode Text

**Database Drivers**
- MySQL
- Postgres
- Oracle
- SQLite
- ODBC
- Netezza

**Readers**
- Dimacs
- DOT
- GXL
- Chaco
- XML
- Tulip
- DelimitedText
- Unicode Delimited Text
- FixedWidth
- ISI, RIS
- Palantir XML

**Multidimensional Analysis**
- TPP / PARAFAC
- Trilinos Integration

**Graph Algorithms**
- Network Communities
- ST Search
- CSG Search
- Temporal Search
- Breadth First Search
- Connected Components
- Biconnected Components
- Brandes Centrality
- Subgraph Isomorphism

**MATLAB Integration**

**R Integration**

**Parallel Statistics Algorithms**
- Descriptive
- Order
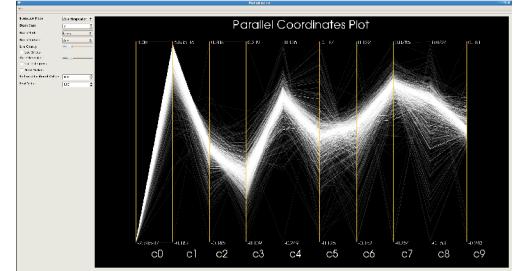- Correlative
- Contingency

**Linear-Time Graph Layouts**
- GSpace
- Hierarchical
- Clustered
- Three tree-based variants

**Text Analysis**
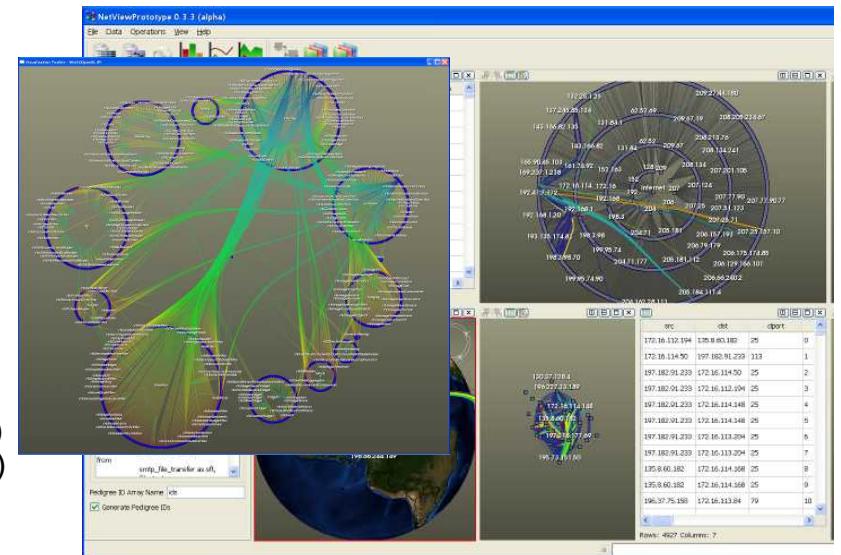- Latent Semantic Analysis (LSA)
- Latent Dirichlet Allocation (LDA)

**Multiple View Types**
- Render (3D)
- Graph
- Hierarchical Graph
- Tree
- Treemap
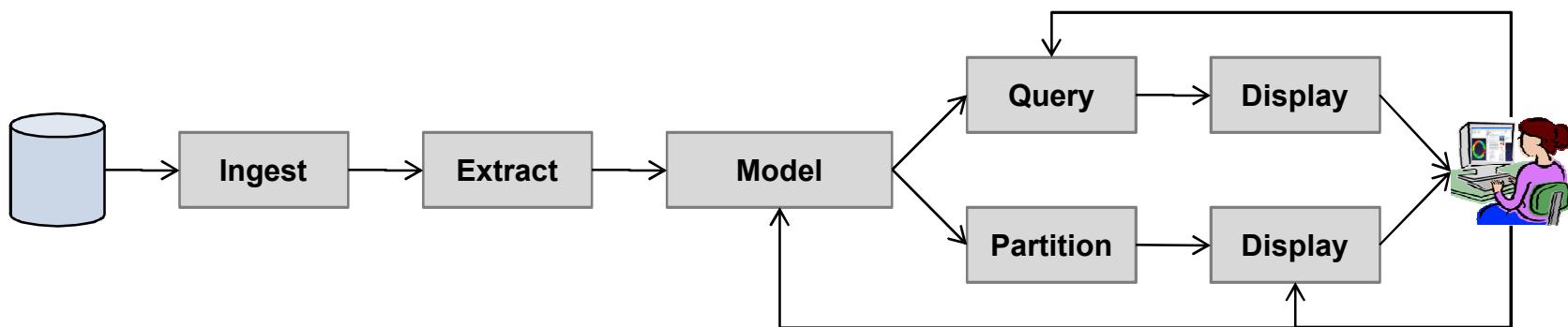- Georeferenced
- Parallel Coordinates
- Radial Tree Ring

**Multiple Platforms / Languages**
- Windows, Linux, OSX, HPC
- Write components in C++
- Use with C++, TCL, Python, Java, .NET, COM
- Use as OverView "plugins"

# Building Applications Using Pipelines

**Applications are built by composing pipelines of filters.**



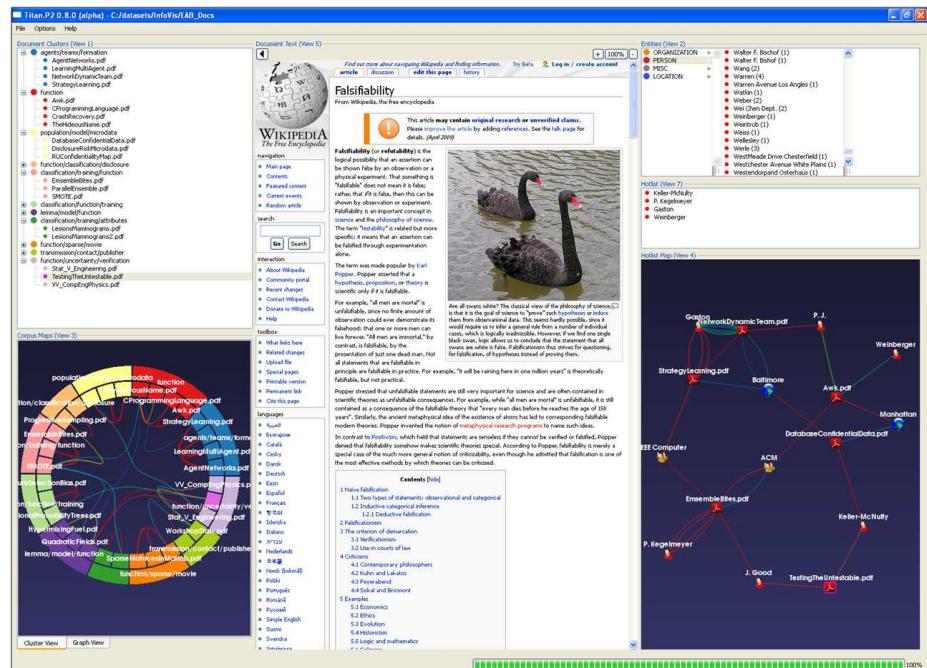**Provides a direct mapping between whiteboard and code.**

**Provides flexibility / exploration (for developers!)**

**Manages execution:**

• **Manages the "flow" of data through components.**

• **In what order are components executed?**

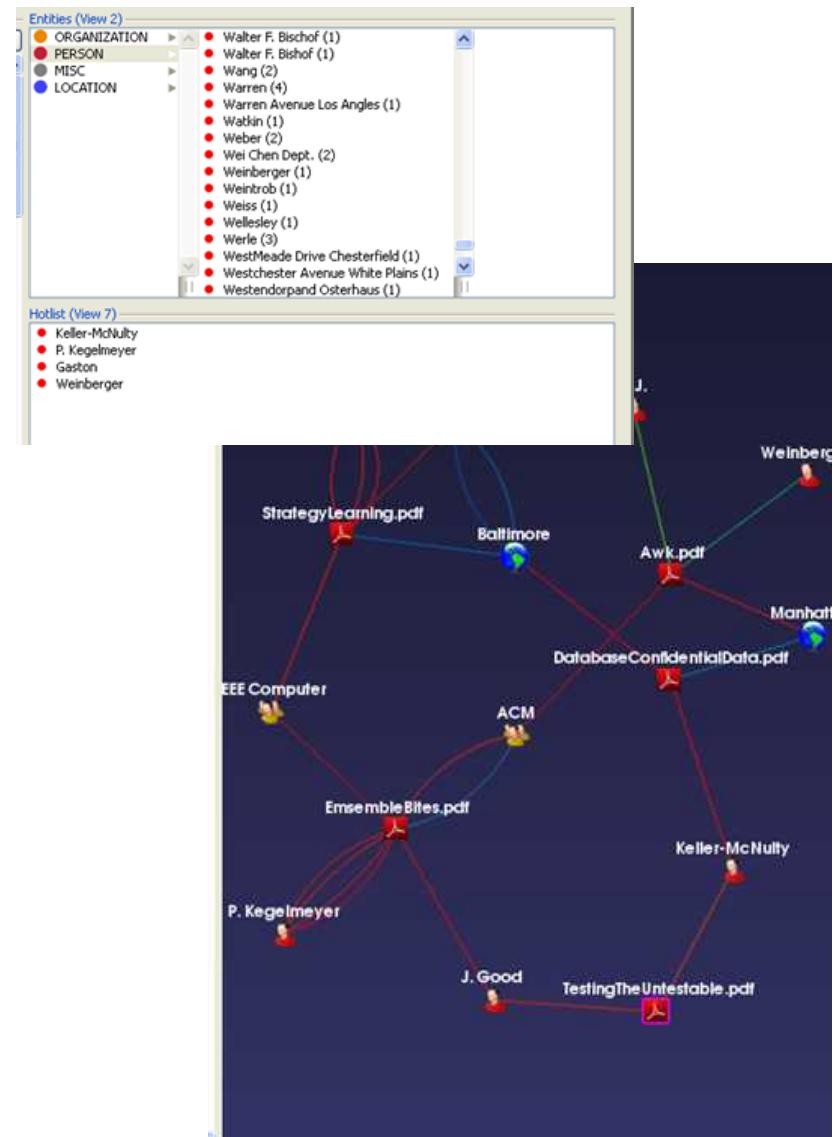• **What needs to be recomputed when parameters change?**

# Application: Text Corpora Analysis

- **Sense-making of a corpus of documents.**

- **Document-Document Similarity**
  - Computed using clustering algorithms (LDA, LSA, etc.)

- **Entity-Document Similarity**
  - Named Entity Recognition
  - Merged DD and DE graphs
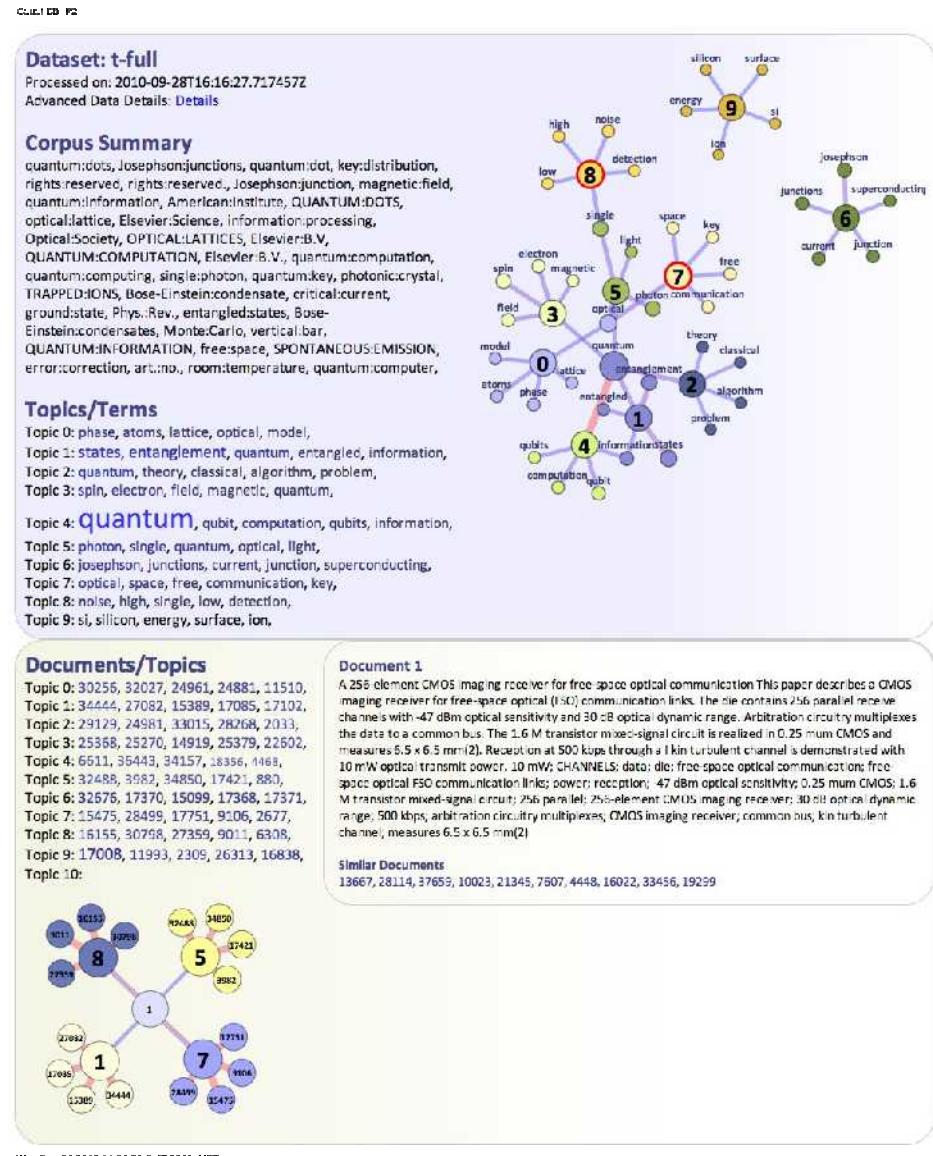  - Connection Subgraphs Searches

# Following a Soap Opera

- The "Hot List" is a subset of named-entities that are "interesting"…

- Generate a community graph that connects selected entities

- Why is this interesting?
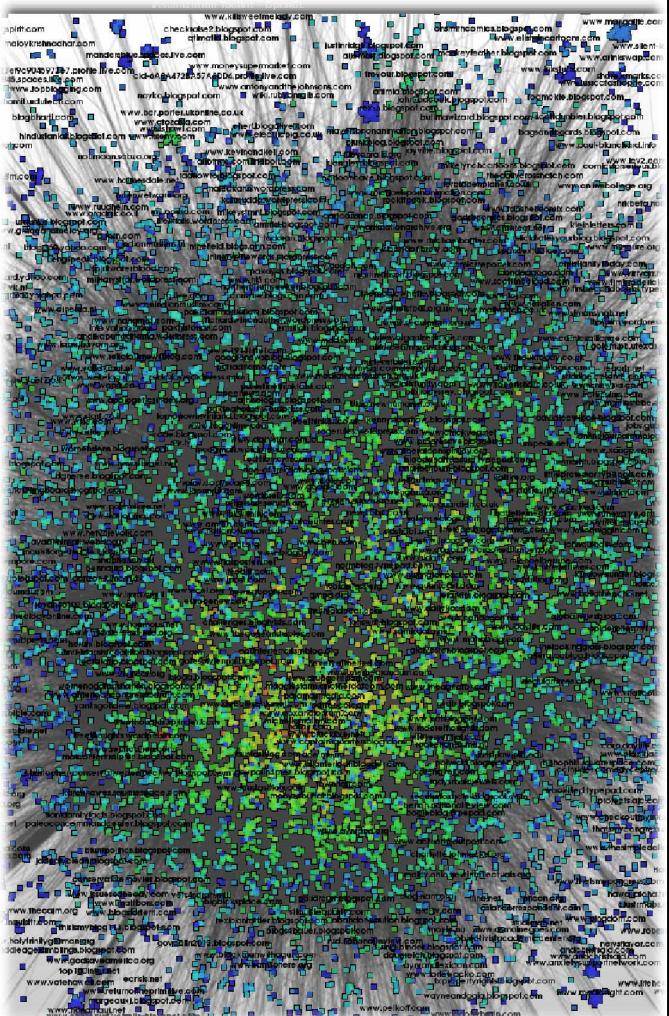  - Shows how entities are related and which documents the info came from

# More Corpora Analysis

- **Problem: Determine the topicality of a set of documents.**
- **Reuses many elements of the previous application data pipeline.**
- **Key differences:**
  - Implements a subset
  - Uses a web-delivery model (CouchDB)
  - 100% browser delivery
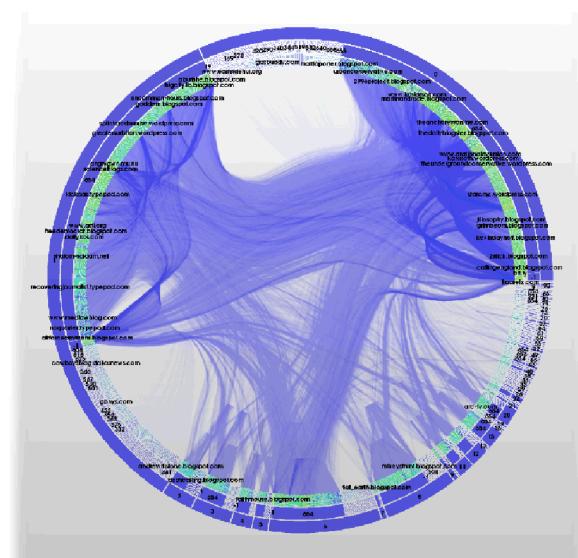  - Developed by two people using 10% time over a couple weeks.

# Web-Crawl Data

- **Web-crawl data project**
- **Analyze a targeted web-crawl to:**
  - Identify topological 'communities' within the web network.
  - Perform text analysis on the page content (LSA, LDA) to pull topics of interest.
  - Cluster the pages around topics
  - Track the evolution of topics within communities over time.
  - Present results as an interactive set of charts and plots accessible through a common web browser.

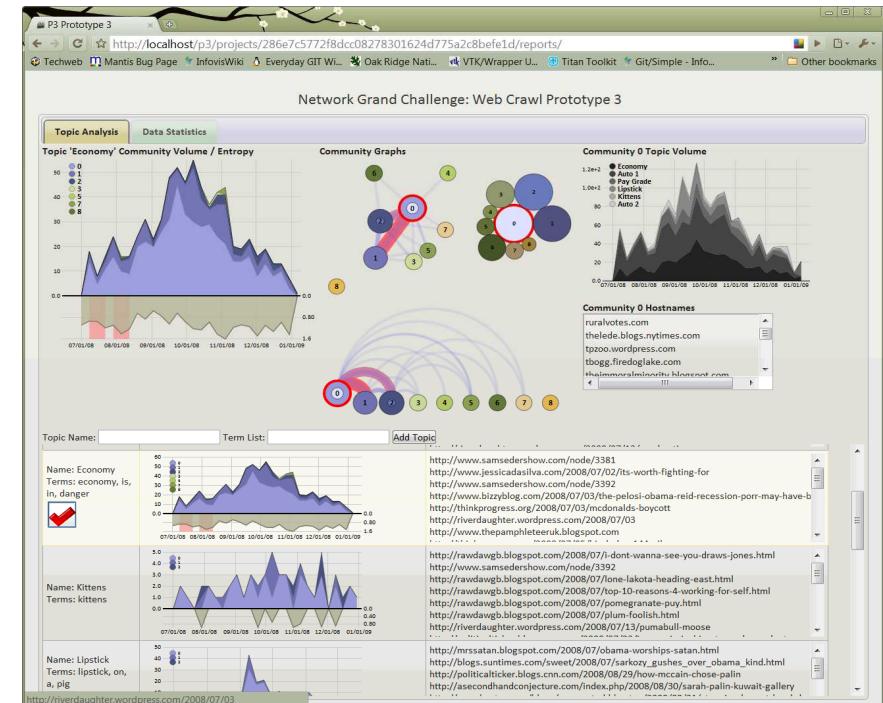# Validating a Graph

- **Often we receive data that we did not collect ourselves.**
- **We must validate the data**
  - Real world data can be messy and/or incomplete.
  - If the data does not conform to expectations, we must discover why.
- **In our web-crawl example, the initial data set was incorrect.**
  - Turned out to be a bug in the post processing scripts that pushed the crawl into a database.

# The Web Analysis Application

- **Toolkit breadth leveraged**
  - Database Drivers
  - Web Server / Client tiers in Titan
  - Latent Semantic Analysis (LSA)
  - Document Clustering
  - Statistics
  - Graph Algorithms
  - Protovis for charts
- **Results delivered using a browser.**

- **App developed in a short time period.**

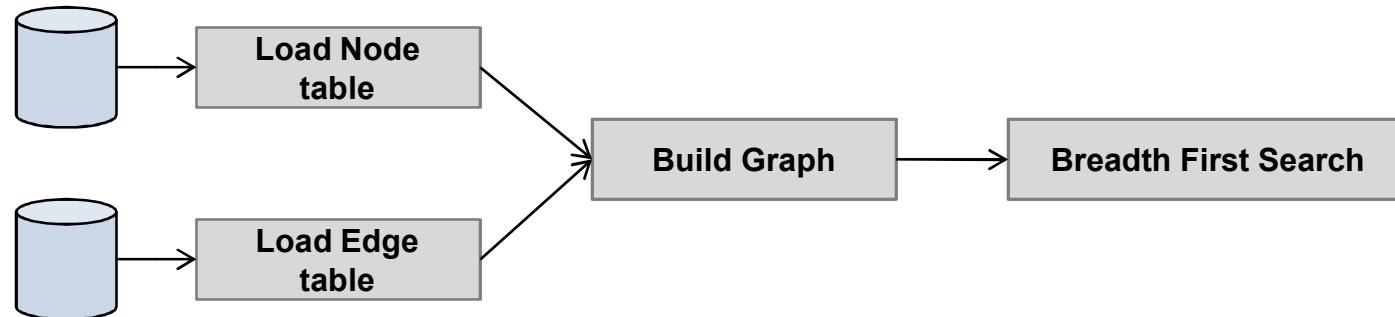# Network Packet Streaming

Titan Toolkit Demo

Realtime Network Packet
Processing and Visualization

15k packets/sec UDP
*(using Google Protocol Buffers)*

- **Real-Time Streaming App**

- **15k packets/second**

- **Combines monolithic app model with 'web' tech**

- **Real-Time modification of views by editing underlying javascript**

# Code Example: Loading a Network

- **The following few slides walks through a quick example that loads a graph using Titan via Python.**
  - Load a graph from two 'CSV' files
  - Run an algorithm on it, Breadth-First Search

# Load the Nodes Table

```
csv_to_graph.py

from vtk import *

nodes = vtkDelimitedTextReader()
nodes.SetFieldDelimiterCharacters(",")nodes.Se
tHaveHeaders(True)
nodes.SetDetectNumericColumns(True)
nodes.SetFileName("nodes.csv")
```
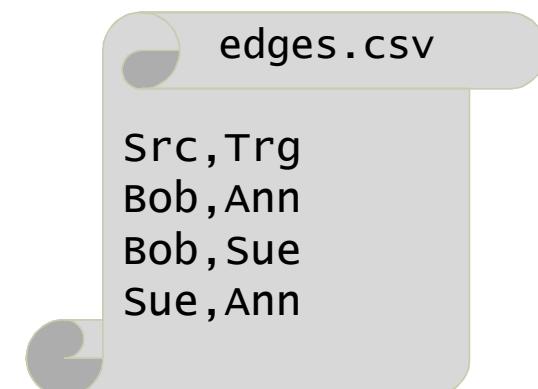
```
nodes.csv

ID,Name,Age
1,Bob,12
2,Ann,25
3,Sue,72
```

- **Explanation**
  - Loads the VTK libraries
  - Load the file *nodes.csv* into the *nodes* pipeline object using *vtkDelimitedTextReader*.

# Load the Edges Table

```
edges = vtkDelimitedTextReader()
edges.SetFieldDelimiterCharacters(",")edges.Se
tHaveHeaders(True)
edges.SetDetectNumericColumns(True)
edges.SetFileName("edges.csv")
```

edges.csv

```
Src,Trg
Bob,Ann
Bob,Sue
Sue,Ann
```

- **Explanation**
  - Essentially same as previous slide, but we're loading the edges file this time.

# Construct Graph from Nodes & Edges

```
Graph = vtkTableToGraph()
Graph.SetDirected(True)
Graph.AddInputConnection(edges,GetOutputPort())
Graph.SetVertexTableConnection(nodes.GetOutputPort())
Graph.AddLinkVertex("src","Name",False)
Graph.AddLinkVertx("trg","Name",False)
Graph.AddLinkEdge("src", "trg")
```

- **Explanation**
    - Uses vtkTableToGraph, which takes two tables (vertices, edges) as inputs.
    - Vertices are created from the nodes table
    - Edges are created by linking the "src" and "trg" entries in rows against nodes using the "Name" field.

# Run the BFS Algorithm

```
bfs = vtkBoostBreadthFirstSearch()
bfs.AddInputConnection(Graph.GetOutputPort())
bfs.SetOriginVertex("Bob")

bfs.Update()
```
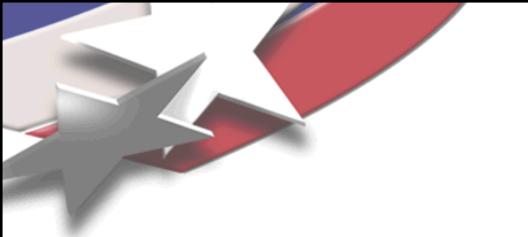
- **Explanation**
  - Attaches vtkBoostBreadthFirstSearch to the pipeline.
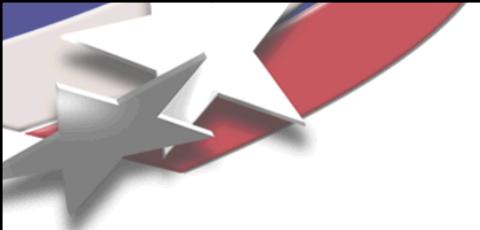  - Computes a BFS ordering on the graph.
- **Pipeline Construction**
  - The 'output port' from one filter is connected to the 'input port' of the next.
  - Once we've got the pipeline started, we can add additional filters to operate on the data.
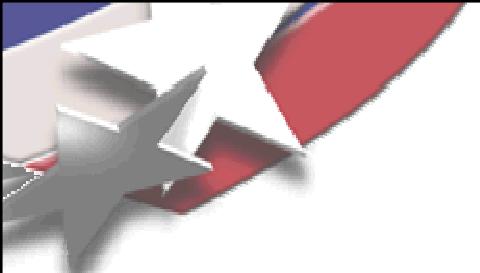
# **Conclusion**

- Titan is a toolkit which unifies many powerful toolkits around a flexible pipeline model.

- This flexibility allows development of applications to solve complex informatics problems with a small developer team.

- Multiple language support (Python, Java, etc.)

- The toolkit is fully open-source and available via Git.

http://titan.sandia.gov

# Questions / Discussion

# *Related Work*

**Prefuse Visualization Toolkit**

A Java-based toolkit for building interactive information visualization applications.

**Tulip Toolkit**

A C++ toolkit for building interactive information visualization applications.

**GraphViz**

A set of libraries specifically for the visualization of many different types of graphs.

**InfoVis Toolkit**

A Java-based toolkit for development of Information Visualization applications and components.

**InfoVis Cyberinfrastructure**

It is a set of libraries that provide a simple and uniform programming-interface to algorithms using the Eclipse Rich Client Platform (RCP).

**Piccolo Toolkit**

Piccolo is a layer built on top of a lower level graphics API. Currently supports Java and C# (Piccolo.Java Piccolo.NET)
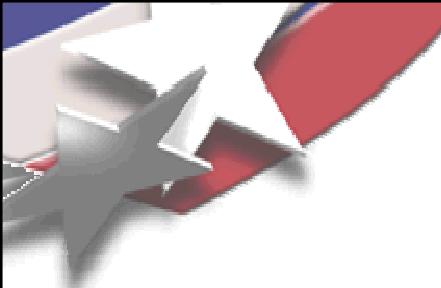
**GeoVista Studio**

GeoVISTA Studio is an open software development environment designed for geospatial data.
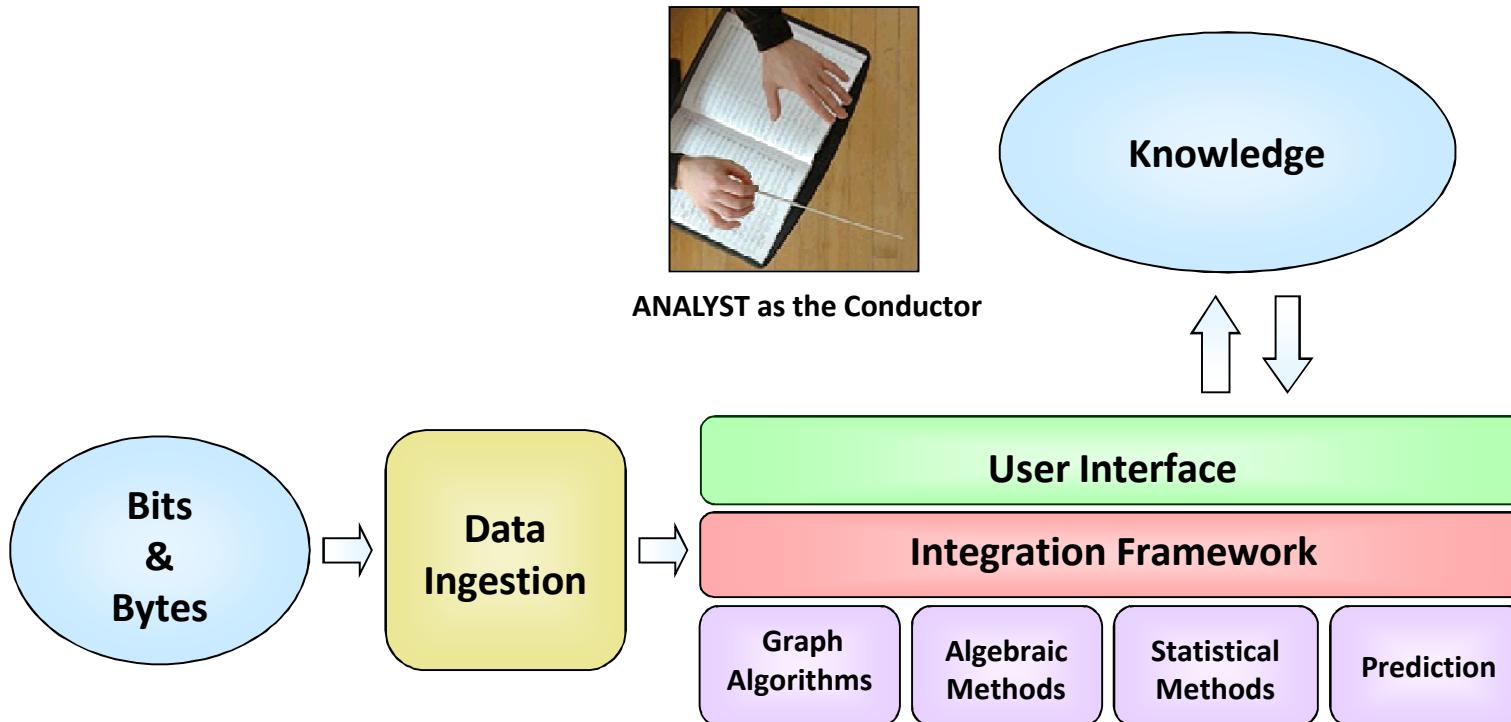
**Improvise**

A coordinated multi-view, Java-based, open source (GPL) information visualization application for MacOS X, Linux, and Windows.

**Protovis**

Open-source library using JavaScript and SVG for web-native visualizations; works with any modern web browser, no plugins required.

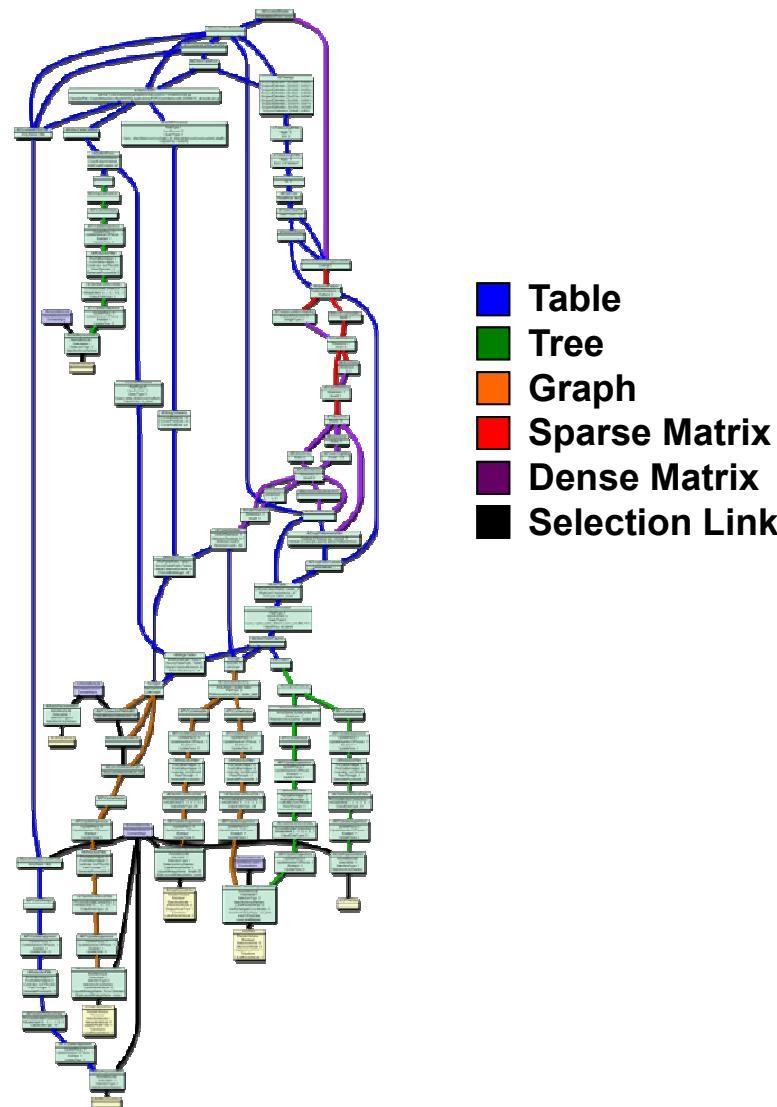# *Titan Informatics Toolkit*

**ANALYST as the Conductor**

**Knowledge**

**Bits & Bytes**

**Data Ingestion**

**User Interface**

**Integration Framework**

| Graph Algorithms | Algebraic Methods | Statistical Methods | Prediction |

# What is Titan?

## A flexible parallel pipeline architecture, written in C++



■ Table
■ Tree
■ Graph
■ Sparse Matrix
■ Dense Matrix
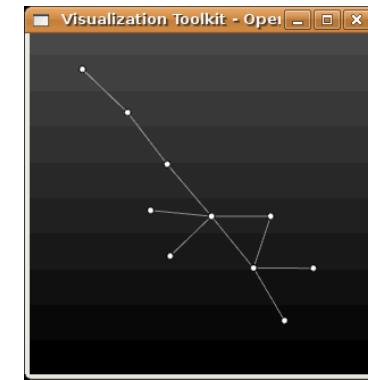■ Selection Link

# Titan Language Bindings

**C++ on Windows XP**

**Python on OSX**

**Tcl/Tk on Linux**

**Java on Vista**

**COM and Excel, on Vista**

**.NET on Vista**

# Titan Programmable Filters

**Titan Data In**
(Tables, Trees, Graphs, etc)

**R Inside**

**R Programmable Filter**

**Titan Data Out**
(New / Updated Tables,
Trees,Graphs, etc)

**User-supplied
R expression**

**Titan Data In**
(Tables, Trees, Graphs, etc)

**Java™ Inside**

**Java Programmable Filter**

**Titan Data Out**
(New / Updated Tables,
Trees,Graphs, etc)

**User-supplied
Java Code**

**Titan Data In**
(Tables, Trees, Graphs, etc)

**XQuery / XSLT Inside**

**XML Programmable Filter**

**Titan Data Out**
(New / Updated Tables,
Trees,Graphs, etc)

**User-supplied
XQuery / XSLT expression**

**Not shown: Python & Matlab™ programmable filters …**

# R Programmable Filter Example

```
vtkRcalculatorFilter* rcf = vtkRcalculatorFilter::New();

rcf->SetInput(source->GetOutput());
rcf->PutTable("x");

rcf->SetRscript("m = do.call(cbind,x)\n \
                 cl <- kmeans(m,3)\n \
                 m = cbind(m,cl$cluster)\n \
                 colnames(m)[4] = \"cluster\"\n");

rcf->GetTable("m");

sink->SetInput(rcf->GetOutput());
```
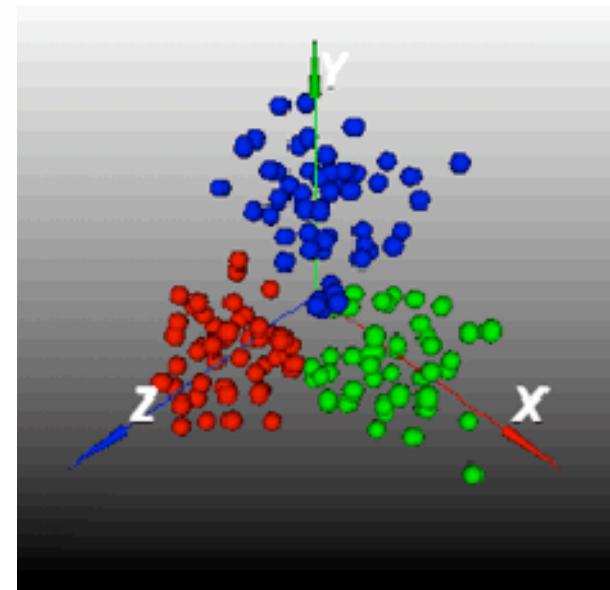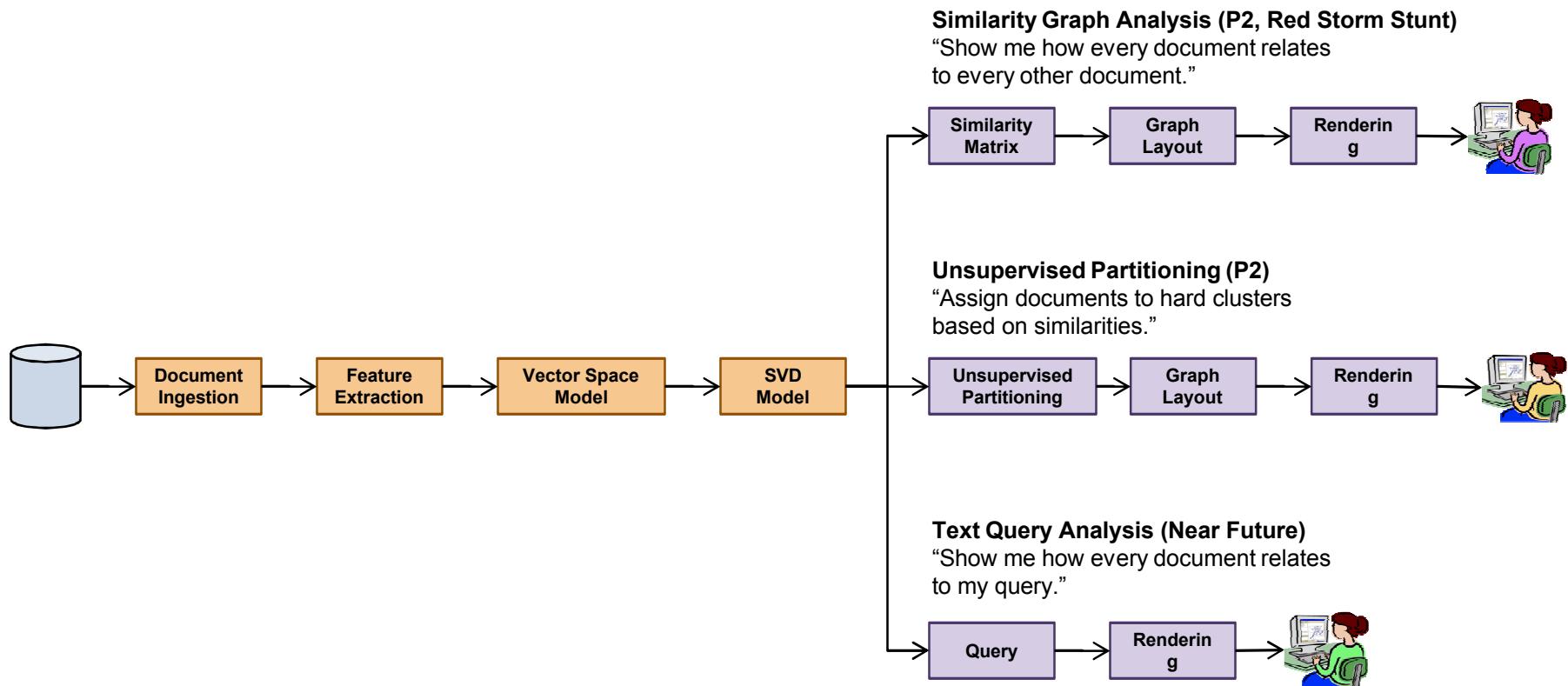
# Representative Titan Modeling + Analysis Pipelines

**Similarity Graph Analysis (P2, Red Storm Stunt)**
"Show me how every document relates to every other document."

| Similarity Matrix | → | Graph Layout | → | Rendering |

**Unsupervised Partitioning (P2)**
"Assign documents to hard clusters based on similarities."

| Document Ingestion | → | Feature Extraction | → | Vector Space Model | → | SVD Model |

| Unsupervised Partitioning | → | Graph Layout | → | Rendering |

**Text Query Analysis (Near Future)**
"Show me how every document relates to my query."

| Query | → | Rendering |

# Representative Titan Modeling + Analysis Pipelines

**Similarity Graph Analysis (P2, Red Storm Stunt)**
"Show me how every document relates to every other document."

| Similarity Matrix | → | Graph Layout | → | Rendering |

**Unsupervised Partitioning (P2)**
"Assign documents to hard clusters based on similarities."

| Document Ingestion | → | Feature Extraction | → | Vector Space Model | → | SVD Model |

| Unsupervised Partitioning | → | Graph Layout | → | Rendering |

**Text Query Analysis (Near Future)**
"Show me how every document relates to my query."

| Query | → | Rendering |

Terabytes    Gigabytes    Megabytes

## "Data Reduction / Aggregation"

# Representative Titan Modeling + Analysis Pipelines

**Similarity Graph Analysis (P2, Red Storm Stunt)**
"Show me how every document relates to every other document."

```
Similarity   →   Graph      →   Renderin
Matrix           Layout         g
```

**Unsupervised Partitioning (P2)**
"Assign documents to hard clusters based on similarities."

```
Document   →  Feature     →  Vector Space  →  SVD    →
Ingestion     Extraction     Model            Model

              Unsupervised  →  Graph    →  Renderin
              Partitioning     Layout      g
```

**Text Query Analysis (Near Future)**
"Show me how every document relates to my query."

```
Query   →   Renderin
            g
```

## "Computational Complexity"

# Representative Titan Modeling + Analysis Pipelines

**Similarity Graph Analysis (P2, Red Storm Stunt)**
"Show me how every document relates to every other document."

| Similarity Matrix | → | Graph Layout | → | Rendering |

**Unsupervised Partitioning (P2)**
"Assign documents to hard clusters based on similarities."

Document Ingestion → Feature Extraction → Vector Space Model → SVD Model →

| Unsupervised Partitioning | → | Graph Layout | → | Rendering |

**Text Query Analysis (Near Future)**
"Show me how every document relates to my query."

| Query | → | Rendering |

Hours ————————————————————— Seconds   Milliseconds

## "Latency"

# Representative Titan Modeling + Analysis Pipelines

**Similarity Graph Analysis (P2, Red Storm Stunt)**
"Show me how every document relates
to every other document."

| Similarity Matrix | Graph Layout | Rendering |

**Unsupervised Partitioning (P2)**
"Assign documents to hard clusters
based on similarities."

| Document Ingestion | Feature Extraction | Vector Space Model | SVD Model |

| Unsupervised Partitioning | Graph Layout | Rendering |

**Text Query Analysis (Near Future)**
"Show me how every document relates
to my query."

| Query | Rendering |

**"Interaction & Feedback"**

# Where We Are Today

**Local Filesystem or Database**



## NGC P2

**Serial Option: put the entire pipeline in one process (P2)**

Zero administration, easy to use, online computation, small data.

# Modest Client / Server Capability for Organizations without HPC

**Server Filesystem or Database**

**Titan Web Service**

**(Based on ParaText)**

**HTTPS**

**NGC P3**

**Client / Server Option: model generation and some analysis on server, remaining analysis in client (P3)**

Modest administration, slightly more complex, some computation offline, modest data sizes.

# Modest Client / Server Capability for Organizations without HPC



**Server Filesystem or Database**

**Web Server**

**(Based on ParaText)**

**HTTPS**

**NGC P3**

**NGC P3**

**Client / Server Option: model generation and some analysis on server, remaining analysis in client (P3)**

Modest administration, slightly more complex, some computation offline, modest data sizes.

# How the NGC will deliver "HPC informatics capabilities that are both *usable* and *useful* to analysts."

Portals, shared file system, or shared database.

**HTTPS**

**HPC Filesystem or Database**

**HPC Compute Nodes**

**(Based on ParaText and the Red Storm Stunt)**

**HPC Service Node**

**NGC P3**

**NGC P3**

**Client X**

**HPC Option: model generation on HPC compute nodes, some analysis on a service node, remaining analysis in client (P3)**

More administration, more offline computation, largest data sizes.

# HPC Informatics for Intelligence Analysis

## HPC Computation Platform

| HPC Database | Query Database and Load Balance | Data Transformations | Algorithms | Layout and Rendering | Standard 100-T Network | Presentation And Interaction |
|---|---|---|---|---|---|---|

**Query**
- *Select * from network_packet_table where date > 7/1/2008 and date < 7/8/2008*

**Data Transformations**
- *vtkTableToTree*
- *vtkTableToGraph*
- *vtkTableToSparseArray*
- *vtkTableToDenseArray*

**Algorithms**
- *Statistics*
- *Linear Algebra*
- *Tensor Methods*
- *Graph Algorithms*
- *Matlab and "R"*
- *MapReduce*

**Layout and Rendering**
- *Tree Layout*
- *Tree Map*
- *Graph Layout*
- *Hierarchical*
- *Geodesic*

**Presentation and Interaction**
- *Client/Server*
- *Geometry and Image Delivery*
- *Windows/Unix/MAC cross platform UI*
- *Linked Selection*

# Web Interfaces