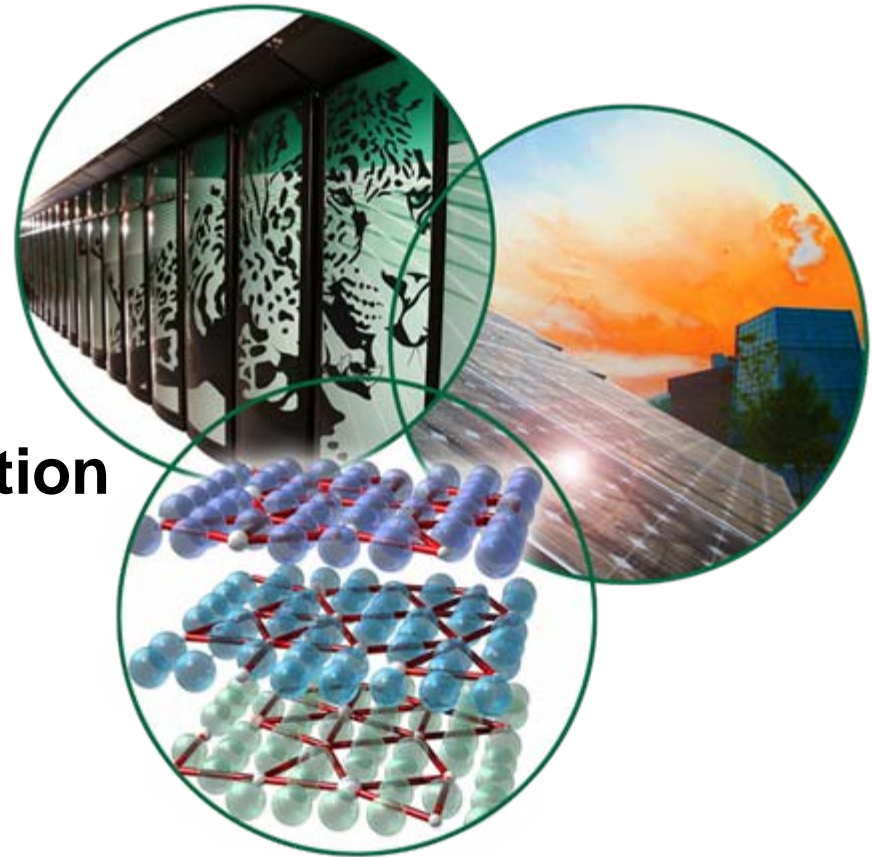# Oak Ridge National Laboratory
## —The "not so foggy" future!

**MG(R) Dennis K. Jackson**
**Director, Logistics Transformation**
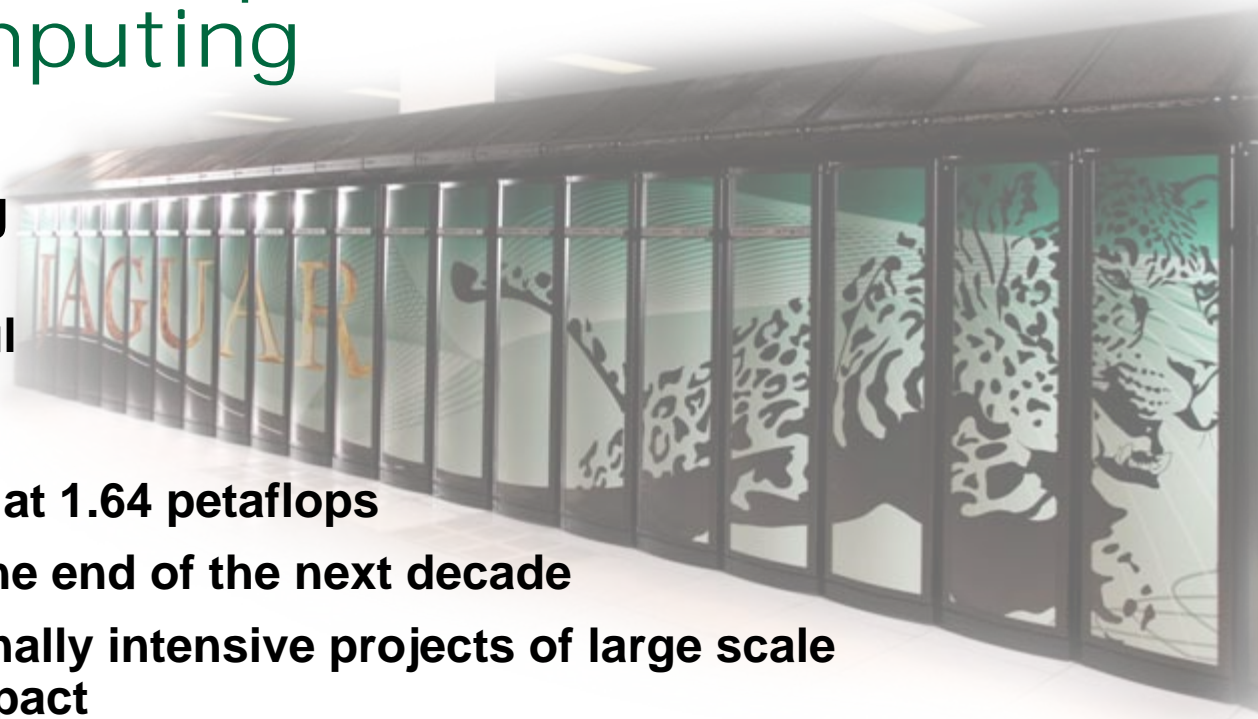**National Security Directorate**

**jacksondk@ornl.gov**
**(865) 574-7382**

U.S. DEPARTMENT OF **ENERGY**

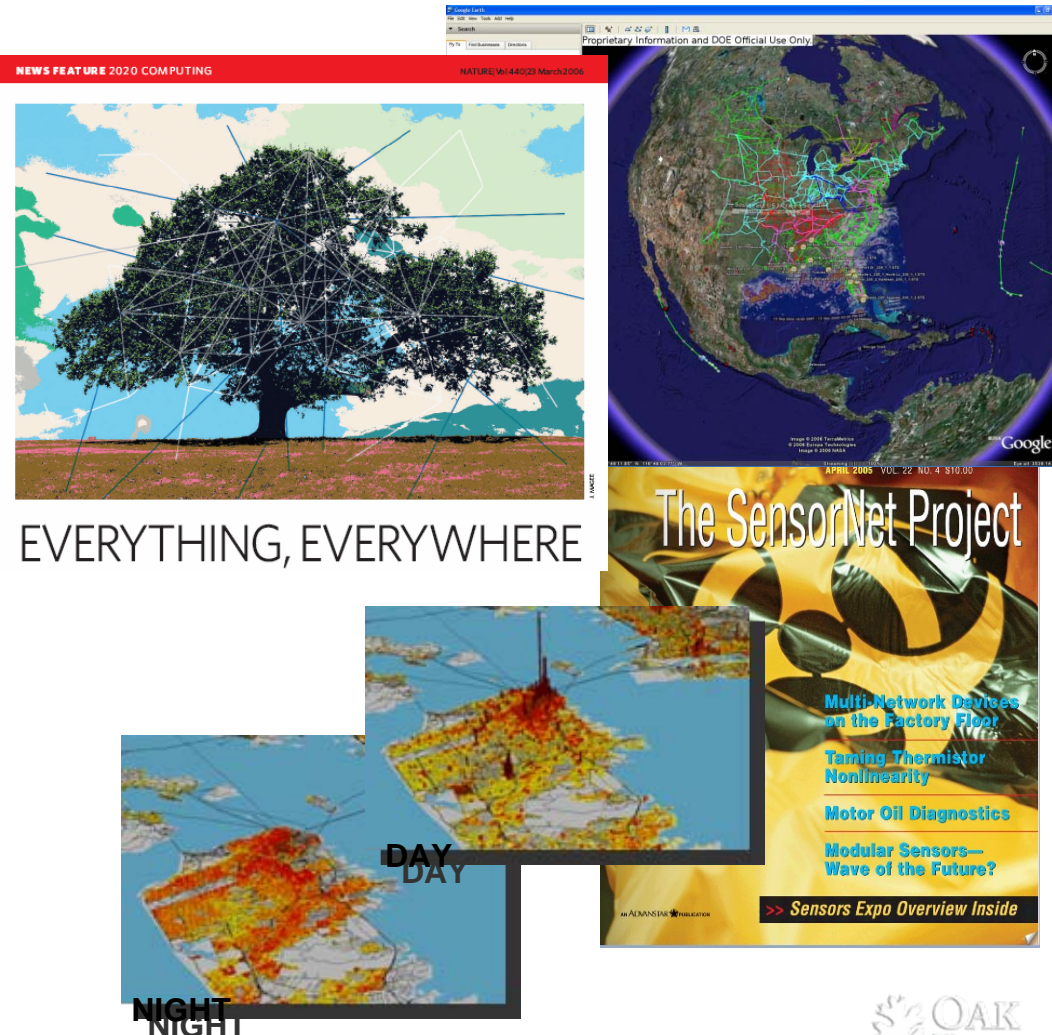OAK RIDGE National Laboratory

# Leading the development of ultrascale scientific computing

- **Leadership Computing Facility:**
  - **World's most powerful open scientific computing facility**
  - **Jaguar XT5 operating at 1.64 petaflops**
  - **Exascale system by the end of the next decade**
  - **Focus on computationally intensive projects of large scale and high scientific impact**

- **ORNL team won the Gordon Bell Prize at SC'08**

- **With the University of Tennessee, developing a second petascale computer for the National Science Foundation**

OAK RIDGE
National Laboratory

# ORNL Is Committed to the Knowledge Discovery Agenda

- **Entire Research Division Focused on Knowledge Discovery**
  - 130 full-time staff
  - 50 subcontractors
  - 50 students

- **Outstanding Resources:**

  **HPC, Networking, MRF, JICS**

- **LDRD Initiative in Knowledge Discovery**

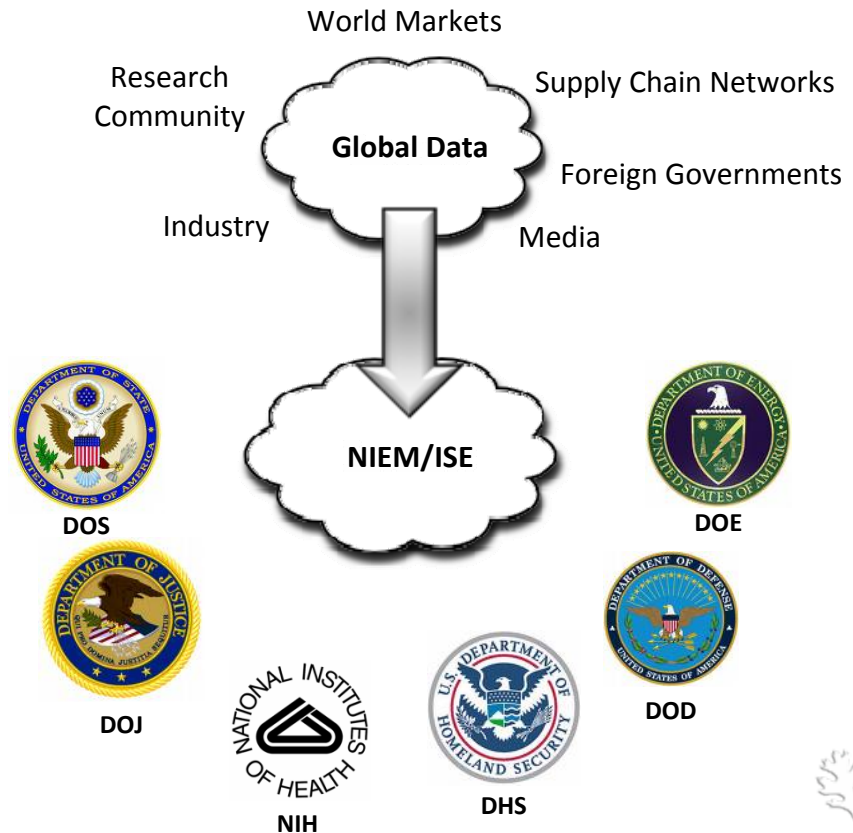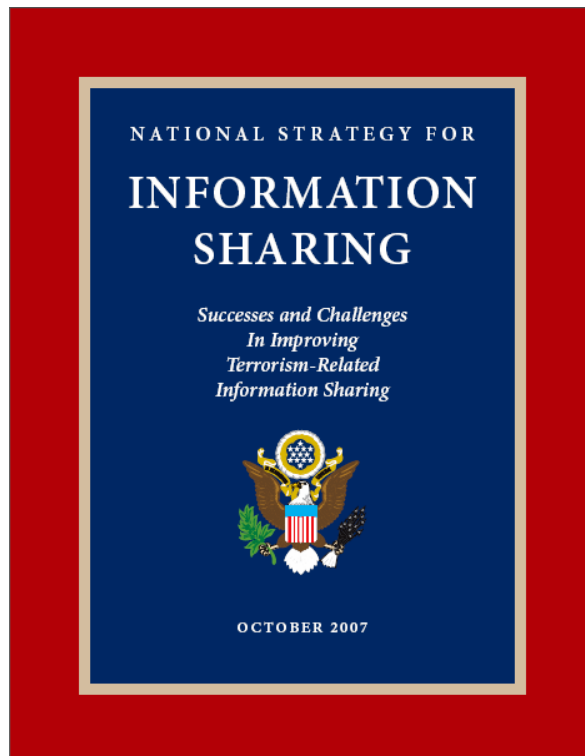- **Programmatic efforts well-aligned with this science agenda**
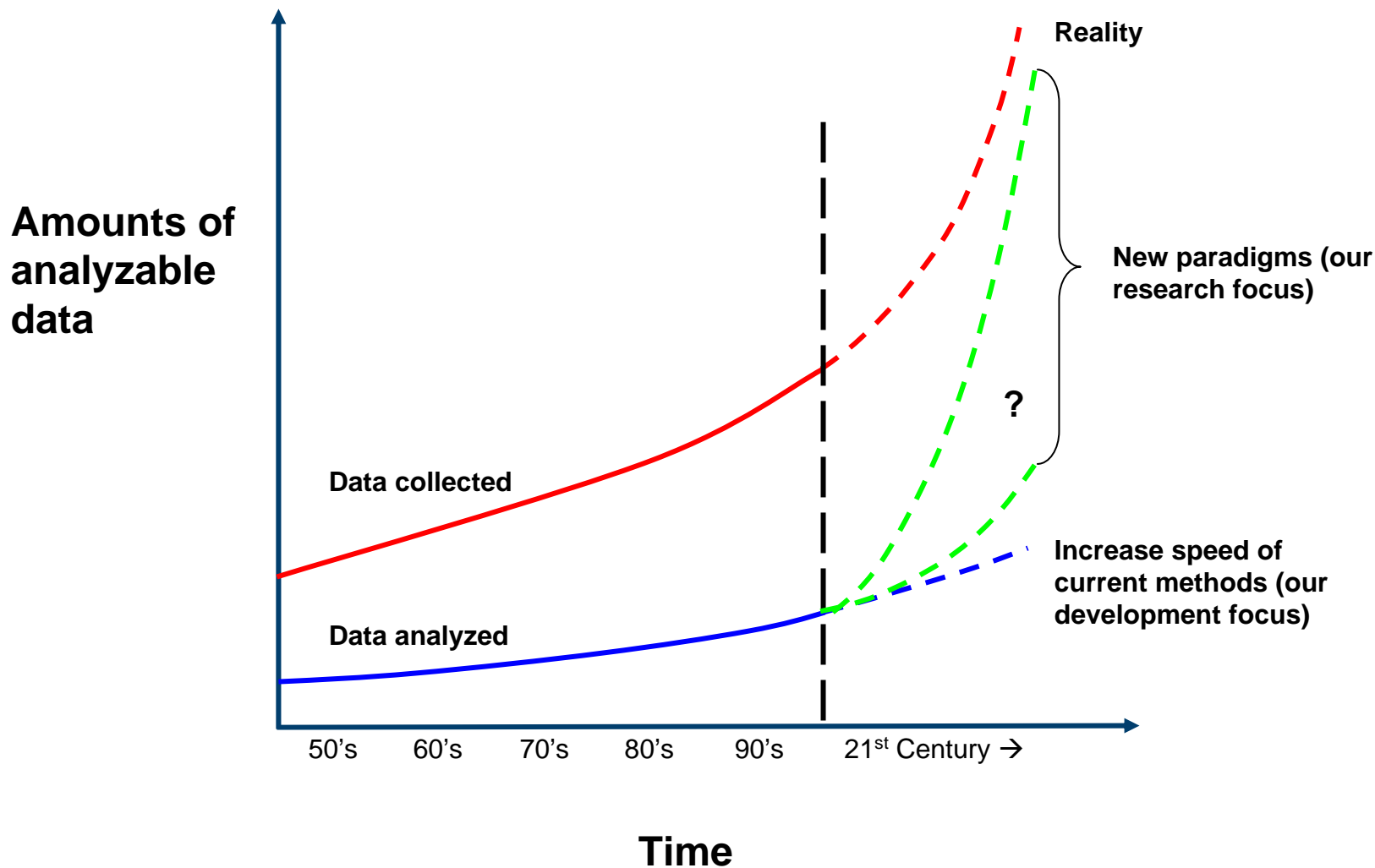
# ORNL's Focus in Knowledge Discovery…

- **Actionable insights from massive, dynamic, disparate data sources**

  - **Knowledge representation of disparate data sources**

  - **High speed analysis and fusion of text, video, audio, and sensor data streams**

  - **Geospatial and temporal data science**

- **Ability to ask more complex questions and detect more complex processes using increasingly higher data resolution**

  - **Population models and population data development**

  - **Modeling and simulation of emerging behavior in complex systems (e.g., social systems)**

  - **Real-time data driven simulations (take advantage of data resolution and availability)**

Managed by UT-Battelle
for the Department of Energy

# Knowledge Discovery Challenge

**How to trigger and coordinate a discovery process across data held by industry, academia, and government agencies within and outside the United States**



NATIONAL STRATEGY FOR
**INFORMATION SHARING**

*Successes and Challenges
In Improving
Terrorism-Related
Information Sharing*

**OCTOBER 2007**

World Markets

Research Community

Supply Chain Networks

**Global Data**

Foreign Governments

Industry

Media

**NIEM/ISE**

DOS

DOE

DOJ

NIH

DHS

DOD

OAK RIDGE National Laboratory

# Knowledge Discovery Challenge

# Research and Development Focus Areas



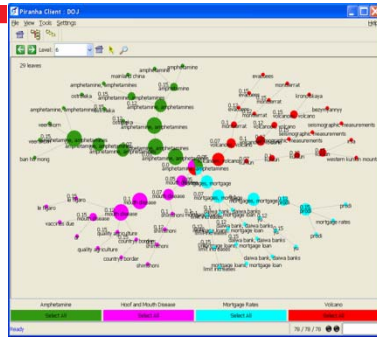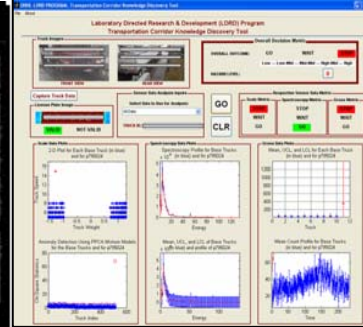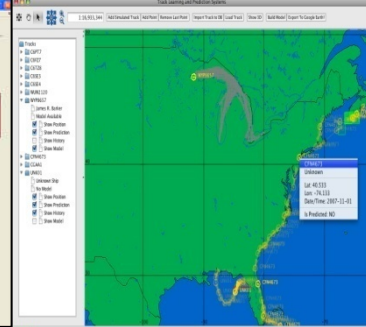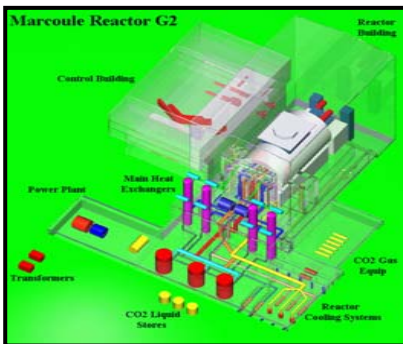**Sensor Networks**   **Analysis in Network**   **Persistent Surveillance**   **Data Fusion**   **Anomaly Detection**

**Predictive Analysis**   **Emergent Behavior**   **Population Dynamics**   **Social Data Analysis**   **Quantum Information**

OAK RIDGE National Laboratory

# Our largest set of projects relate to collection, analysis, and dissemination of sensor data.

- **Interdiction, detection, emergency response**
  - **Mobile, Transportation Corridors, Ports, Military Bases**

- **Real-Time Data Management**
  - **Collection, Dissemination, Archiving**

- **Pre-deployment analysis**
  - **Cost, Performance Prediction, Risk vs Benefit**

- **Wide-area ubiquitous sensing, actuation, and deployment**
  - **Orchestrating the functionality across a large system of distributed sensors/processors (eg Electric Grid, Autonomous robotic systems)**

- **Cross-agency and cross-administrative boundary data-sharing and interoperability**
  - **Standards and policies**

- **Net-Centric Services**

- **Security, Access Controls**

# Social Network for Sharing Sensor Data



## SensorPedia

**Addresses the ability to access and fuse data from disparate sensor networks**

**Use of Web 2.0 "social networking" technologies (e.g., RSS, wikis, podcasts, mashups, blogs, and streaming video)**

**Key identity management and credentialing standards**

**Data owner controls publishing and subscribing**

**Explores how volunteered sensor data is being used and shared**

Managed by UT-Battelle
for the Department of Energy

# Knowledge Representation for Situation Awareness of the Electric Grid

**Wide-area Grid View**

**Outages**

**Weather**

**Impacts**

**Streaming Analysis**

- **Organize, stream, and fuse data from various sources through an analysis pipeline**

- **Present an intuitive visualization of the status to end-users**

OAK RIDGE
National Laboratory

# Where are all my local, state, and federal assets?

- **What assets can I track at all times?**

- **How well can I estimate the location of non-tracked assets?**

- **What computational resources will be required?**

- **What are the uncertainties?**

## Parallel/Distributed Discrete Event Simulation Engines

| Model Execution | Synchronization | Data Integration | Interoperability | ... | Multi-Scale |
|---|---|---|---|---|---|

| Super computers | Clusters | Multi-Cores | GPGPUs | ... | PDAs |
|---|---|---|---|---|---|

**IBM Blue Gene Award for scalable algorithms**

**Best Paper Award for agent-based methods**

**Tackling DTRA 10**5 persistent surveillance grand challenge**

OAK RIDGE National Laboratory

# Event Spectroscopy: Natural Disasters

ABC

NBC

CBS

PBS

CNN

FOX News

MS NBC

REUTERS

AP Associated Press

Google News

moreover

RSS
BLOG

Text
Time-series

## Textual Prism

| Technical | Social |
|-----------|--------|
| Power | Water (health) |
| Communications | Movement |
| Roads | Medical |
| Lights | Health Hazards |
| Water (Systems) | Crime |
| Outage | Shelters |

**Social Spectral Components**



Tornadoes
Jan 7, 2007

Ice Storm
Jan 14, 2007

**Technical Spectral Components**

OAK RIDGE
National Laboratory

# Piranha

## Knowledge Discovery Engine

*2007 R&D 100 Award Winner*

**Managed by UT-Battelle**
**for the Department of Energy**

# Textual Analysis

- **We understand this problem**
  - **8 years of research**
  - **40+ Papers, 3 patents**
  - **Hands on experience with DHS, Military, IC, and Industry**

- **We are very good at it**
  - **$15M in research investment**
  - **19 group members**
  - **R&D 100 Award (Oscars of invention) in 2007**

Presentation_name

OAK RIDGE
National Laboratory

# We can read a newspaper, but not a library … without help

# How can computers help?

- **The smartest computer can not read a simple first grade level book**

- **But simple computers can help us find what we need when we need it**

Managed by UT-Battelle
for the Department of Energy

Presentation_name

# Overview of Text Analysis

- **Keyword Methods – Very fast, good for millions**
  - **Search**
    - **"Seafood in DC"**
    - **Good if you know what you are looking for and can find it on the top of the result list**
  - **Unsupervised Classification**
    - **"What were the main topics in message traffic last month?"**
    - **Good to get a general overview a set of messages, though topics may not be valuable**
  - **Supervised Classification**
    - **"What explosive and trigger messages were in last months traffic?"**
    - **Good for finding topics of interest, provided you can describe the topics**

Managed by UT-Battelle
for the Department of Energy

Presentation_name

OAK RIDGE
National Laboratory

# Overview of Text Analysis

- **Full text methods – Slower, good for thousands**
  - **Clustering**
    - **"How are these set of documents related"**
    - **Good for organizing sets of documents done statistically, which may differ from human organization.**
  - **Term frequency Analysis**
    - **"What other words or concepts am I missing"**
    - **Good for linking terms and names, best suited for well written documents**
  - **Semantic Extraction – Slow but parallelizable**
    - **"I am out of ideas, what else can you find"**
    - **Good for the needle in a haystack analysis, but can be very slow.**

Presentation_name

OAK
RIDGE
National Laboratory

# How computers can help

**Document 1**

> The Army needs sensor technology to help find improvised explosive devices

**Terms**

Army
Sensor
Technology
Help
Find
Improvise
Explosive
device

**Term List**

Army
Sensor
Technology
Help
Find
Improvise
Explosive
Device
ORNL
develop
homeland
Defense
Mitre
won
contract

**Document 2**

> ORNL has developed sensor technology for homeland defense

ORNL
develop
sensor
technology
homeland
defense

**Document 3**

> Mitre has won a contract to develop ho defense se for explos devices

Mitre
won
contract
develop
...
devices

## Vector Space Model

|            | Doc 1 | Doc 2 | Doc 3 |
|------------|-------|-------|-------|
| **Army**       | 1 | 0 | 0 |
| **Sensor**     | 1 | 1 | 1 |
| **Technology** | 1 | 1 | 0 |
| **Help**       | 1 | 0 | 0 |
| **Find**       | 1 | 0 | 0 |
| **Improvise**  | 1 | 0 | 0 |
| **Explosive**  | 1 | 0 | 1 |
| **Device**     | 1 | 0 | 1 |
| **ORNL**       | 0 | 1 | 0 |
| **develop**    | 0 | 1 | 1 |
| **homeland**   | 0 | 1 | 1 |
| **Defense**    | 0 | 1 | 1 |
| **Mitre**      | 0 | 0 | 1 |
| **won**        | 0 | 0 | 1 |
| **contract**   | 0 | 0 | 1 |

## *Documents to vectors*

OAK RIDGE National Laboratory

# Textual Clustering

**Vector Space Model**

|            | Doc 1 | Doc 2 | Doc 3 |
|------------|-------|-------|-------|
| Army       | 1     | 0     | 0     |
| Sensor     | 1     | 1     | 1     |
| Technology | 1     | 1     | 0     |
| Help       | 1     | 0     | 0     |
| Find       | 1     | 0     | 0     |
| Improvise  | 1     | 0     | 0     |
| Explosive  | 1     | 0     | 1     |
| Device     | 1     | 0     | 1     |
| ORNL       | 0     | 1     | 0     |
| develop    | 0     | 1     | 1     |
| homeland   | 0     | 1     | 1     |
| Defense    | 0     | 1     | 1     |
| Mitre      | 0     | 0     | 1     |
| won        | 0     | 0     | 1     |
| contract   | 0     | 0     | 1     |

**Similarity Matrix**

|       | Doc 1 | Doc 2 | Doc 3 |
|-------|-------|-------|-------|
| Doc 1 | 100%  | 17%   | 21%   |
| Doc 2 |       | 100%  | 36%   |
| Doc 3 |       |       | 100%  |

*Documents to Documents*

**Cluster Analysis**

D1    D2    D3

*Most similar documents*

**TFIDF**

$$W_{ij} = \log_2\left(f_{ij} + 1\right) * \log_2\left(\frac{}{n}\right)$$

$$d_2(\mathbf{x}_i, \mathbf{x}_j) = \left(\sum_{k=1} (x_{i,k} - x_{j,k})\right)$$

*Vectors to trees*

**Time Complexity**

$O(n^2 \text{Log } n)$

# Challenge

- **Current computer algorithms that process text work well for small sets of documents**
  - **Average newspaper story .0001 seconds**

- **Not as well for medium size sets**
  - **Encyclopedia Britannica 2.3 days**

- **Infeasible for large sets.**
  - **British newspapers from 1800 – 1900 requires 317 years of processing**

OAK
RIDGE
National Laboratory

# ORNL Breakthrough...

$$W_{ij} = \log_2(f_{ij} + 1) * \log_2\left(\frac{C+1}{c+1}\right)$$
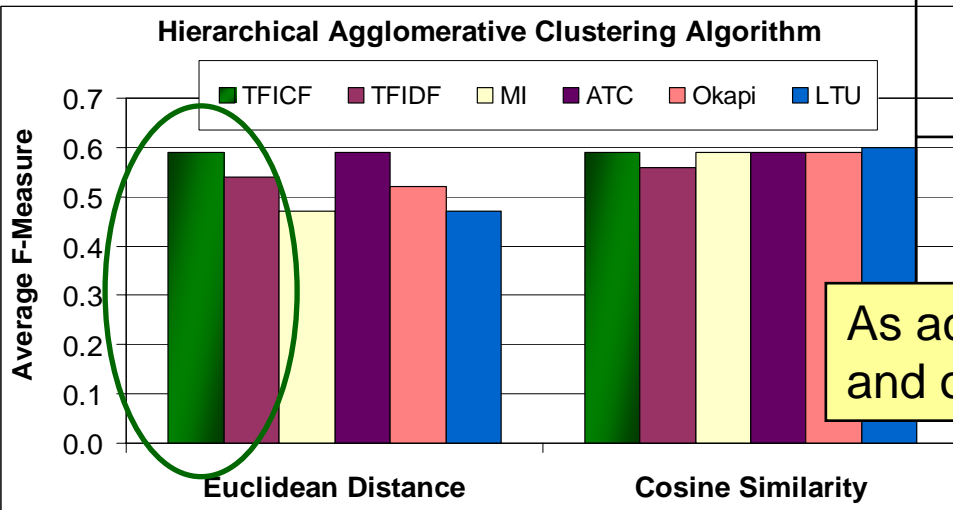
Inverse Corpus Frequency

## Term Weighting Schemes

| Name | Term Weighting Scheme |
|------|----------------------|
| TF-IDF | $wij = \log(fij) \times \log(N/nj)$ |
| MI | $wij = \log \dfrac{\frac{fij}{N}}{\frac{\sum_{i=1}^{N} fij}{N} \times \frac{\sum_{j=1}^{M} fij}{N}}$ |
| ATC | $wij = \dfrac{\left(0.5 + 0.5 \times \frac{fij}{max\_f}\right)\log\left(\frac{N}{nj}\right)}{\sqrt{\sum_{i=1}^{N}\left[\left(0.5 + 0.5 \times \frac{fij}{max\_f}\right)\log\left(\frac{N}{nj}\right)\right]^2}}$ |
| Okapi | $wij = \left(\dfrac{fij}{0.5 + 1.5 \times \frac{dl}{avg\_dl} + fij}\right)\log\left(\dfrac{N - nj + 0.5}{fij + 0.5}\right)$ |
| LTU | $wij = \dfrac{(\log(fij) + 1.0)\log\left(\frac{N}{nj}\right)}{0.8 + 0.2 \times \frac{dl}{avg\_dl}}$ |

## Test Data Sets

| Data Set | # of Docs | # of Classes | Largest Class | Smallest Class |
|----------|-----------|--------------|---------------|----------------|
| Reuters | 2349 | 58 | 1041 | 1 |
| SMART | 3891 | 3 | 1460 | 1033 |
| 20 News | 4650 | 12 | 399 | 385 |



**Hierarchical Agglomerative Clustering Algorithm**

Legend: TFICF, TFIDF, MI, ATC, Okapi, LTU

Average F-Measure (0.0–0.7)

Euclidean Distance    Cosine Similarity

As accurate as current methods and orders of magnitude faster

Presentation_name

OAK RIDGE National Laboratory

# Capability overview

| Capability | Capacity in documents | Piranha | Search Engines | Natural Language Processing Tools |
|---|---|---|---|---|
| Search | 100M+ | Yes | Yes | No |
| Unsupervised classification | 1M | Yes | Some | No |
| Supervised classification | 1M | Yes | No | No |
| Clustering | 100K | Yes | No | No |
| Term Frequency Analysis | 100K | Yes | Yes, but not available to user | Yes |
| Semantic Extraction | 1000 | Yes | No | Yes |

Managed by UT-Battelle
for the Department of Energy

Presentation_name

OAK RIDGE
National Laboratory

# Large scale data exploration constrained by wall-clock time to provide decision support.

- **Detect anomalies**

- **Data dip into structured and unstructured data**

- **Inductive hypothesis generation**

- **Human interaction enhanced by real-time data support**

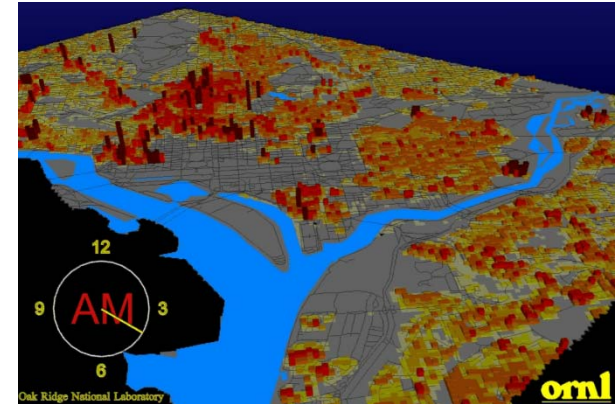- **Threat anticipation**

# Population Data and Models



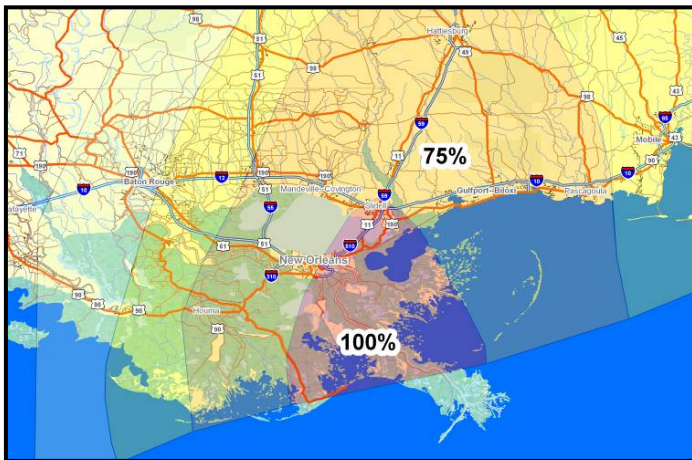Population
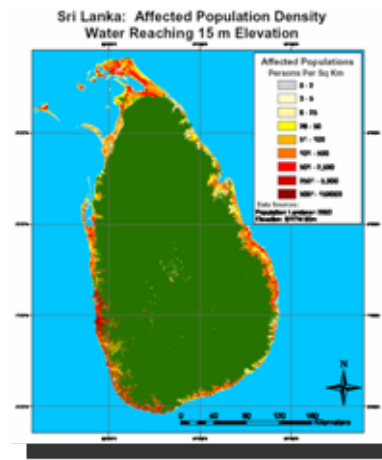ORNL LandScan Global Population Project

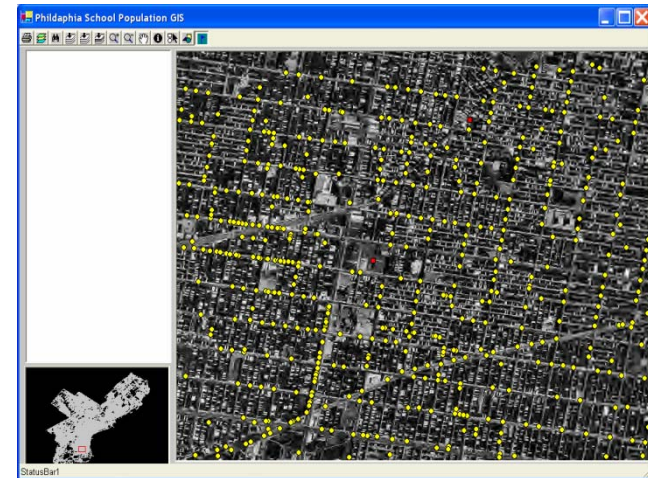**LandScan Global 30"x30"**

**LandScan USA Day/Night 3"x3"**

**Nominal 24-hour variation**

**Hurricane Impacts**
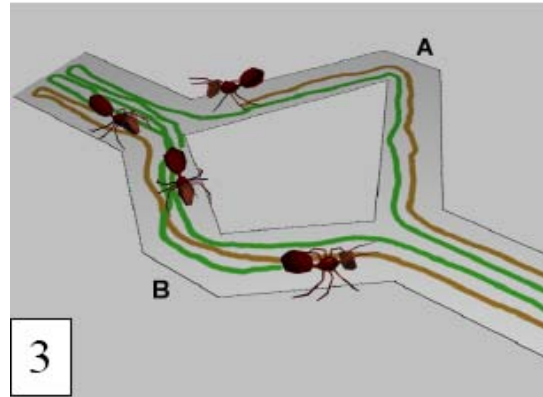
**Tsunami Impacts**

**Exposure Impacts**

# Emergent Behavior in Social Systems



**Birds flocking**



**Ant pathways**



**Human response**
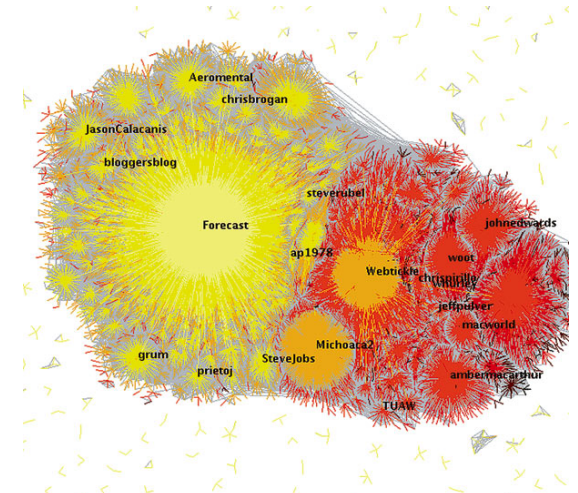
*Agent-based simulations*
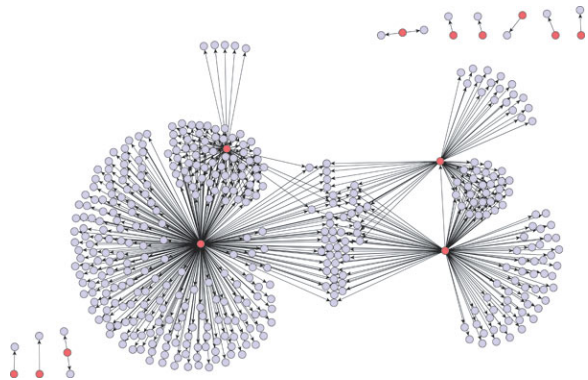
*Discrete-event simulations*
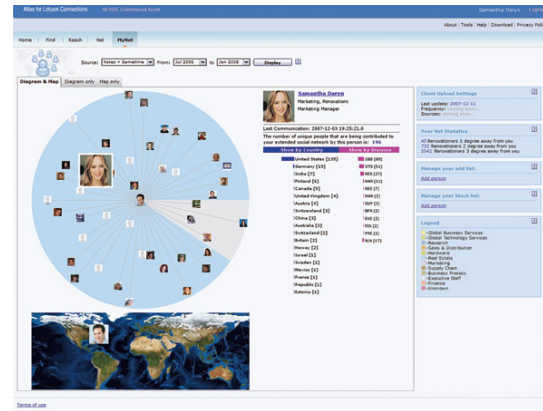
# Social Networks Analysis



**Blogosphere**



**Comment Flow**



**Twitter Social Network**



**Viral Marketing**



**Workplace Networks**

Images from Technology Review, Vol 111/No 2.  March/April 2008

OAK RIDGE National Laboratory

# Virtual Worlds to Explore Social Behaviors



**Second Life – Linden Lab**

**Education**

**Tourism**

**Collaboration**

**Shopping**

**Interviews**

OAK RIDGE National Laboratory

# Achieving Systematic Situation Awareness

**Collect** (1) → **Integrate** (2) → **Detect** (3) → **Network/Share** (4)



Syndromic Surveillance



Field Sensors



Monitoring



Knowledge Base



Anomaly Detection



Across Agencies

**Reconcile** (7) ← **Respond** (6) ← **Evaluate/Assert** (5)







Cross Check

Managed by UT-Battelle
for the Department of Energy

# Summary

- **Current technology cannot yet solve emerging national challenges in knowledge discovery**

- **Intelligent software agents and associated research areas comprise <u>significant</u> breakthrough technology**

- **Results indicate <u>high-potential</u> to help solve these national challenges**

- **We have a progression of significant and successfully deployed agent systems and research to our credit**

Managed by UT-Battelle
 for the Department of Energy

Presentation_name

OAK
RIDGE
National Laboratory