

# Humane Data Mining: The New Frontier

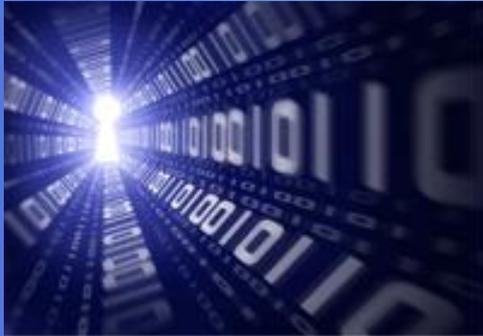
*Rakesh Agrawal*

*Microsoft Search Labs  
Mountain View, California*

*Updated version of the SIGKDD-06 keynote*



# Thesis



- Data Mining has made tremendous strides in the last decade
- It's time to take data mining to the next level of contributions
- We will need to develop new abstractions, algorithms and systems, inspired by new applications

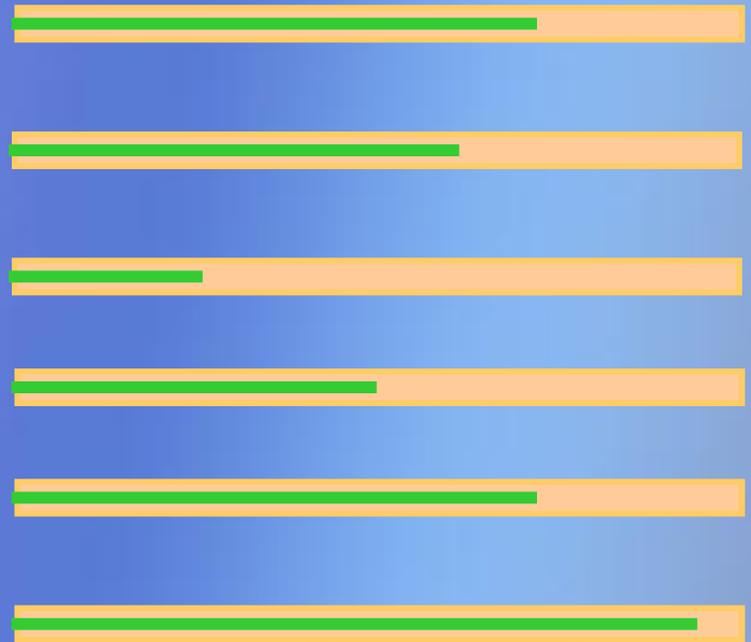
# Outline



- Progress report
- New Frontier

# A Snapshot of Progress

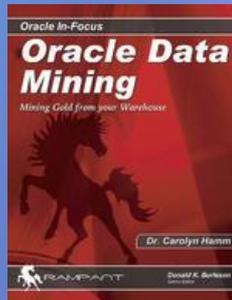
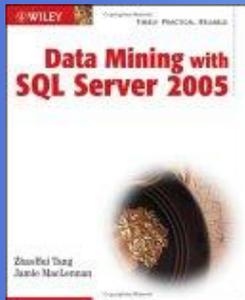
- System support
- Algorithmic innovations
- Foundations
- Usability
- Enterprise applications
- Unanticipated applications



# Database Integration



- Tight coupling through user-defined functions and stored procedures
- Use of SQL to express data mining operations
  - Composability: Combine selections and projections
  - Object-relational extensions enhance performance
  - Benefit of database query optimization and parallelism carry over
- SQL extensions

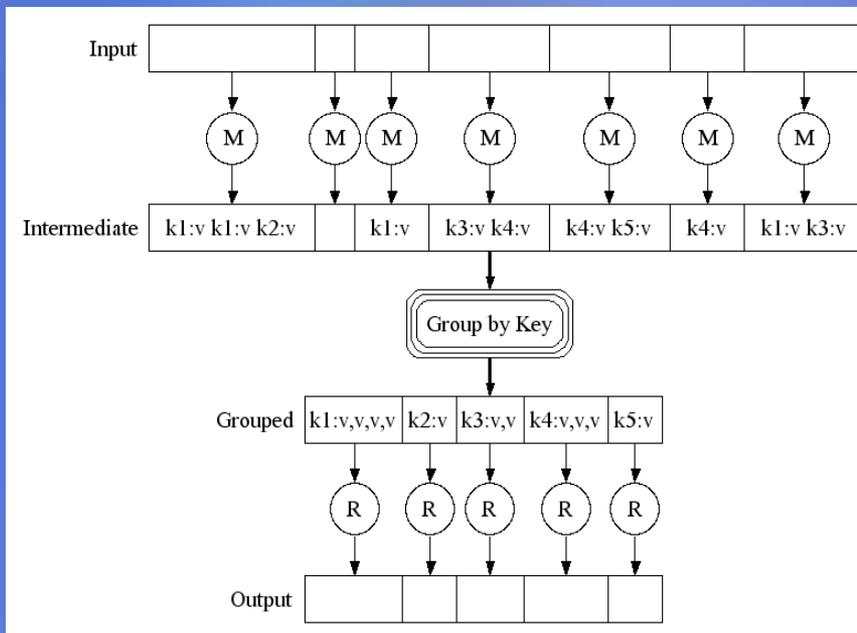


# Google's Data Mining Platform

## MapReduce<sup>1</sup>: Programming Model

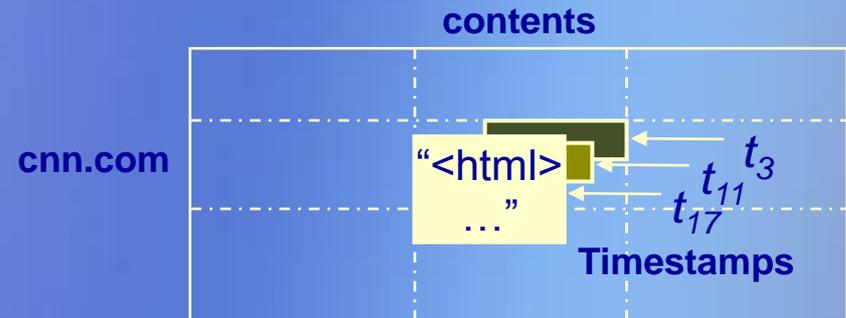
map(ikey, ival) -> list(okey, tval)

reduce(okey, list(tval)) -> list(oval)



- Automatic parallelization & distribution over 1000s of CPUs
- Log mining, index construction, etc

## BigTable<sup>2</sup>: Distributed, persistent, multi-level sparse sorted map



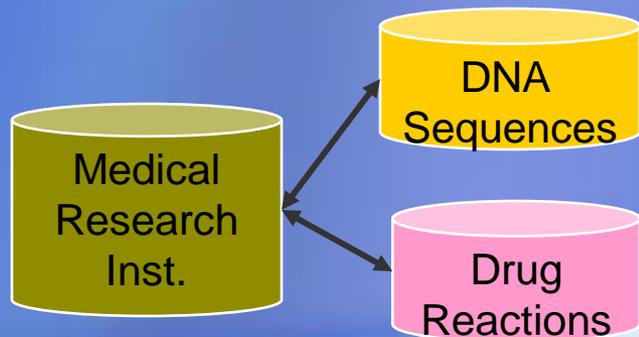
- Tablets, Column family
- >400 Bigtable instances
- Largest manages >300TB, >10B rows, several thousand machines, millions of ops/sec
- Built on top of GFS

<sup>1</sup>Dean et. al. "MapReduce: Simplified data processing on large clusters", OSDI 04.

<sup>2</sup>Hsieh. "BigTable: A distributed storage system for structured data", Sigmod 06.

# Sovereign Information Integration

- Separate databases due to statutory, competitive, or security reasons.
  - Selective, minimal sharing on a need-to-know basis.
- Example: Among those patients who took a particular drug, how many with a specified DNA sequence had an adverse reaction?
  - Researchers must not learn anything beyond counts.
- Algorithms for computing joins and join counts while revealing minimal additional information.



## Minimal Necessary Sharing

R	
a	
u	
v	
x	

S	
b	
u	
v	
y	

$R \Join S$

- R must not know that S has **b** and **y**
- S must not know that R has **a** and **x**

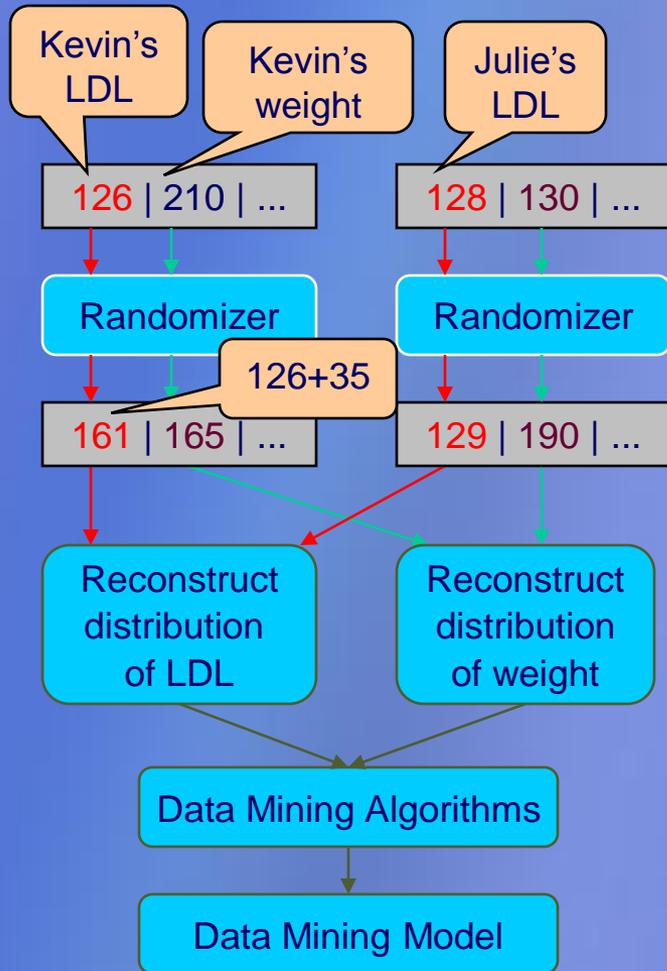
$R \Join S$

u
v

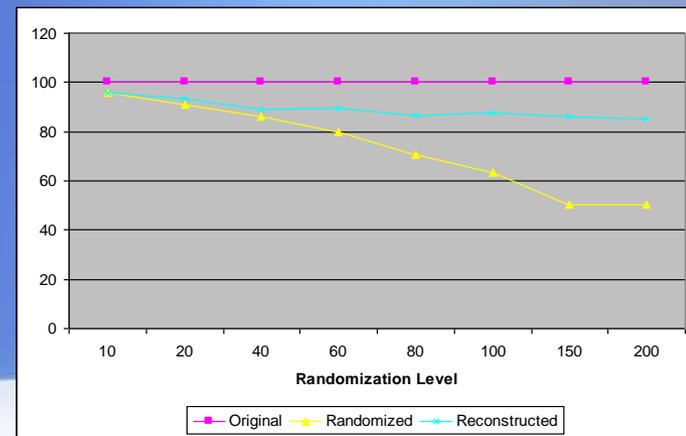
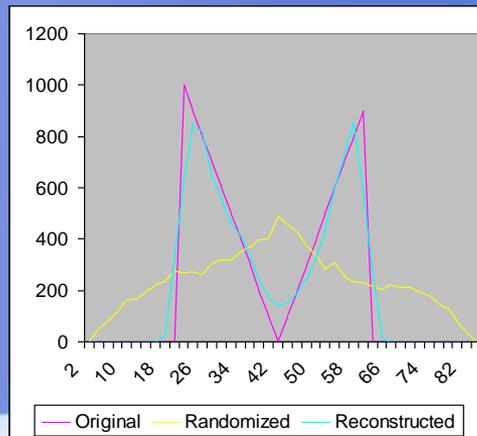
**Count ( $R \Join S$ )**

- R and S do not learn anything except that the result is 2.

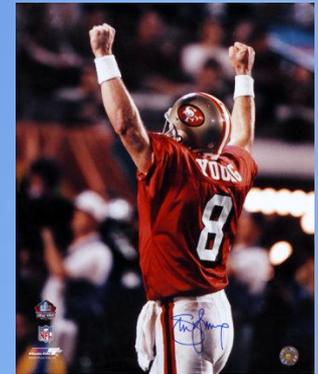
# Privacy Preserving Data Mining



- Preserves privacy at the individual patient level, but allows accurate data mining models to be constructed at the aggregate level.
- Adds random noise to individual values to protect patient privacy.
- EM algorithm estimates original distribution of values given randomized values + randomization function.
- Algorithms for building classification models and discovering association rules on top of privacy-preserved data with only small loss of accuracy.



# Enterprise Applications Galore!



- Example: SAS Customer Successes

## Customer Relationship Management

[Claims Prediction](#) | [Credit Scoring](#) | [Cross-Sell/Up-Sell](#) |  
[Customer Retention](#) | [Marketing Automation](#) | [Marketing Optimization](#) |  
[Segmentation Management](#) | [Strategic Enrollment Management](#)

## Drug Development

## Financial Management

[Activity-Based Management](#) | [Fraud Detection](#)

## Human Capital Management

## Information Technology Management

[Charge Management](#) | [Resource Management](#) |  
[Service Level Management](#) | [Value Management](#)

## Performance Management

[Balanced Score-carding](#)

## Quality Improvement

## Regulatory Compliance

[Fair Banking](#)

## Risk Management

## Supplier Relationship Management

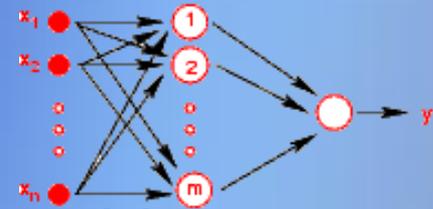
## Supply Chain Analysis

[Demand Planning](#) | [Warranty Analysis](#)

## Web Analytics

<http://www.sas.com/success/solution.html>

# Ordering Search Results



- Search results ranked using a learning algorithm (neural net).
- Training data: Query/document pairs labeled for relevance (perfect, excellent, good, etc.).

baseball	<a href="http://sports.espn.go.com/mlb/index">http://sports.espn.go.com/mlb/index</a>	Perfect
baseball	<a href="http://en.wikipedia.org/wiki/Baseball">http://en.wikipedia.org/wiki/Baseball</a>	Good

- Query independent (e.g. static page rank) as well as query dependent features (e.g. position of a query word in anchor text) for every document.

baseball <http://sports.espn.go.com/mlb/index> perfect f1 f2 ...

# Related Searches

The screenshot shows a search engine interface with a search bar containing the word "baseball". To the right of the search bar is a green "Search" button and links for "Advanced" and "Options". Below the search bar, the results are categorized under "Web results" with a count of "1-10 of 145,000,000". There are links for "Images", "Video", "News", "Maps", "MSN", and "More".

Below the main results, there are sponsored sites:
 

- MLB Tickets** - www.ticketsnow.com/mlb (Sponsored sites)
  - Lowest Prices Of The Year On MLB Tickets. Limited Time!
- Free Baseball Scores** - scoresandodds.com
  - Update Scores, Game Previews, Matchups, Injuries, Trends and More

At the bottom, there are local listings for "baseball near pleasanton, california" (1 more) with a link to "Swinging Santero's" (925) 846-1199 Pleasanton.

On the right side, there is a sidebar titled "Related searches" enclosed in a red rounded rectangle. It lists the following related queries:
 

- Major League Baseball
- Baseball Pictures
- Baseball Equipment
- Baseball Bats
- Baseball Hall Of Fame
- Baseball Cards
- Soccer
- Football

- Most popular queries containing the current query
- Analysis of how users reformulated their queries
  - Football —————> Soccer (whole query)
  - Wildflower cafe ———> Wildflower bakery (piecewise)
- Query click graph to find related queries



Query:

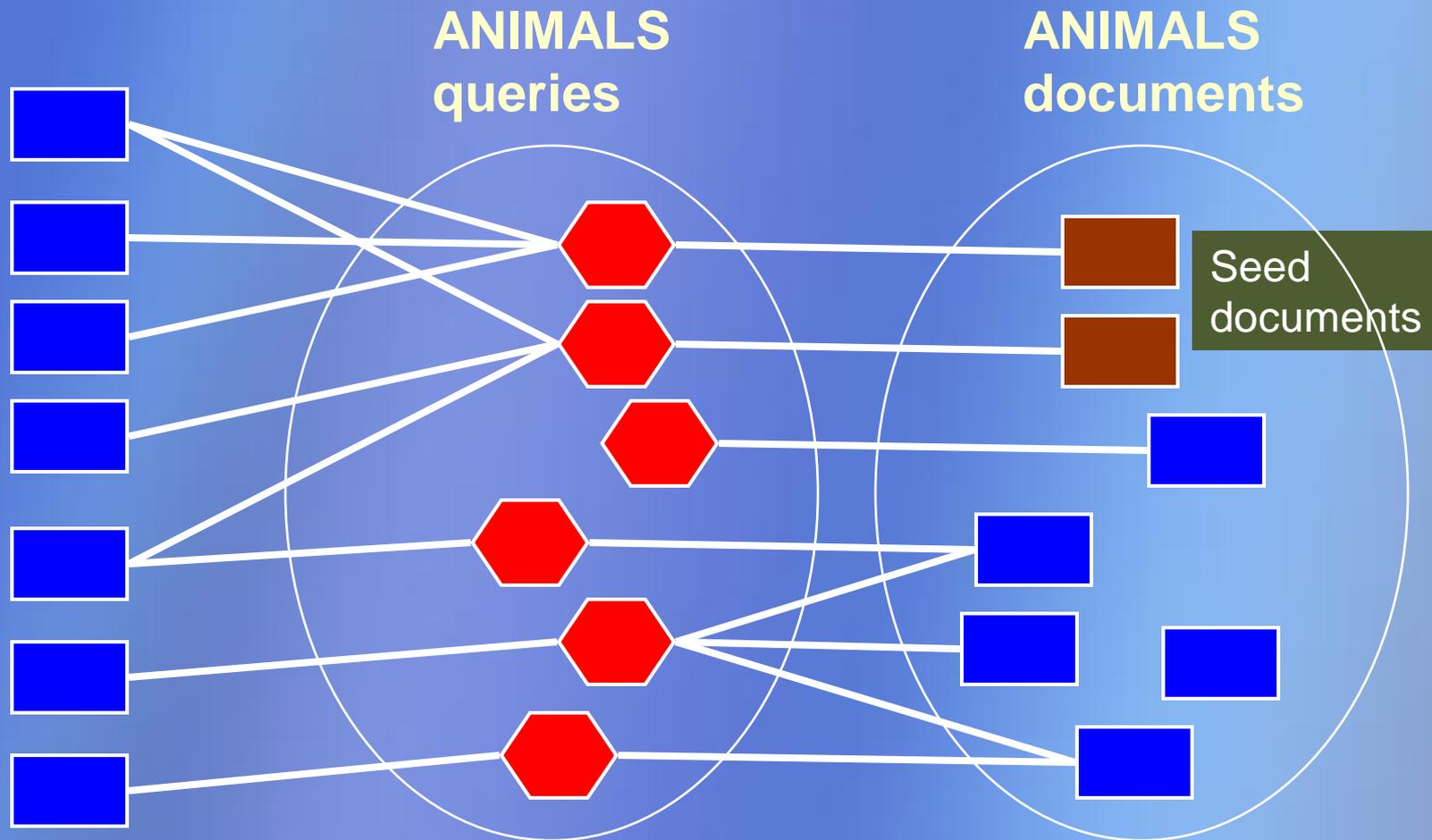
---

business/industries/transportation\_and\_logistics  
<http://www.panynj.gov/aviation/jfkframe.HTM>  
[http://queens.about.com/cs/transportation/a/jfk\\_airport.htm](http://queens.about.com/cs/transportation/a/jfk_airport.htm)

science\_and\_research/social\_sciences\_and\_humanities/history  
[http://en.wikipedia.org/wiki/JFK\\_%28film%29](http://en.wikipedia.org/wiki/JFK_%28film%29)  
[http://en.wikipedia.org/wiki/John\\_F.\\_Kennedy](http://en.wikipedia.org/wiki/John_F._Kennedy)  
[http://en.wikipedia.org/wiki/JFK\\_Reloaded](http://en.wikipedia.org/wiki/JFK_Reloaded)  
<http://www.coverups.com/jfk/index.htm>  
<http://mcadams.posc.mu.edu/home.htm>  
<http://www.csun.edu/~lms30894/index.html>  
<http://www.panynj.gov/aviation/jalframe.htm>

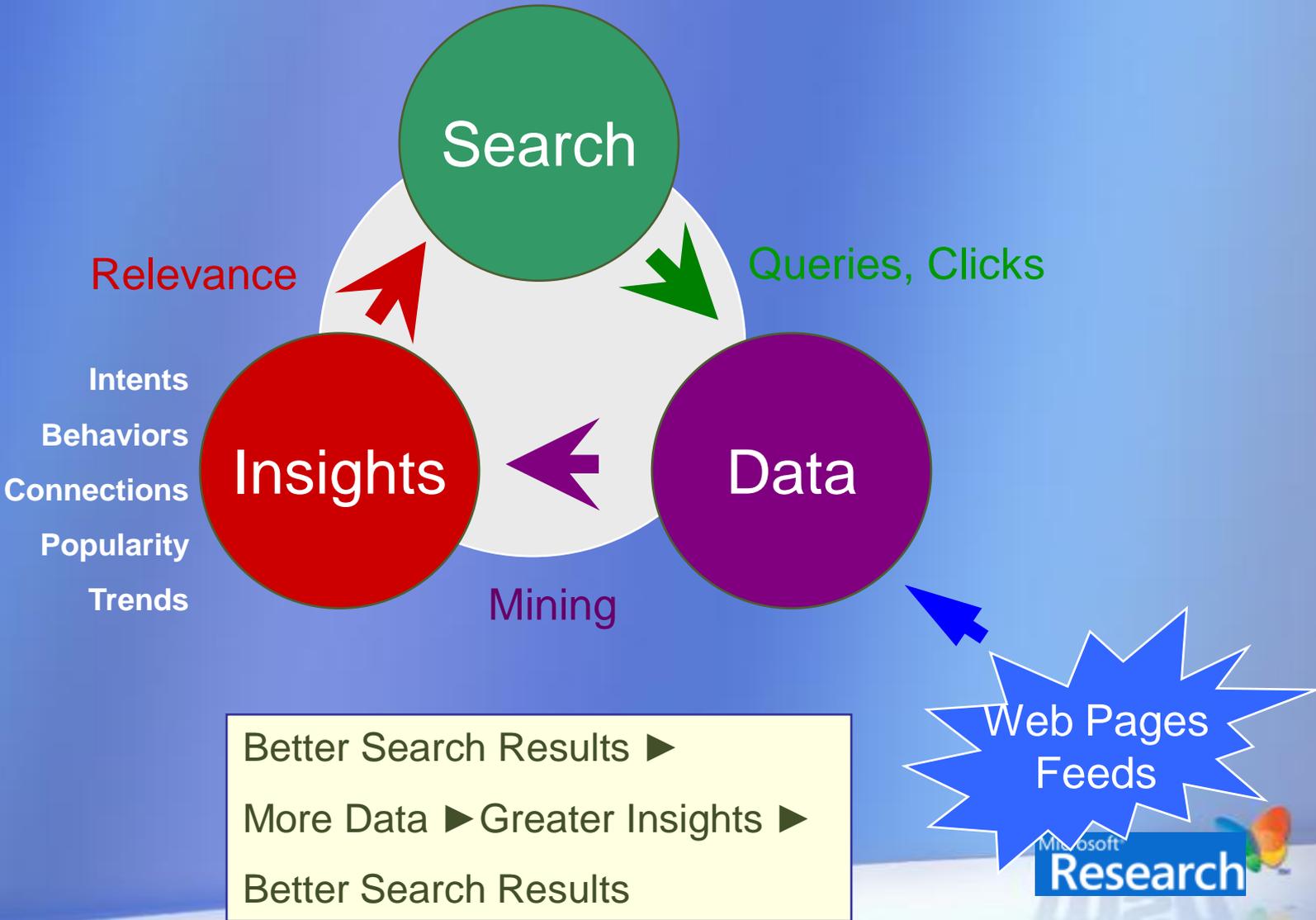
- Ideas from portfolio theory to allocate space to different result types
- Marginal utility of adding a document decreases if the result set already contains high quality documents of the same type
- Query and document classification using merged click logs

# Classification Using Click Graph



Algorithm: Random walk with absorbing states

# Search & Data: Virtuous Cycle

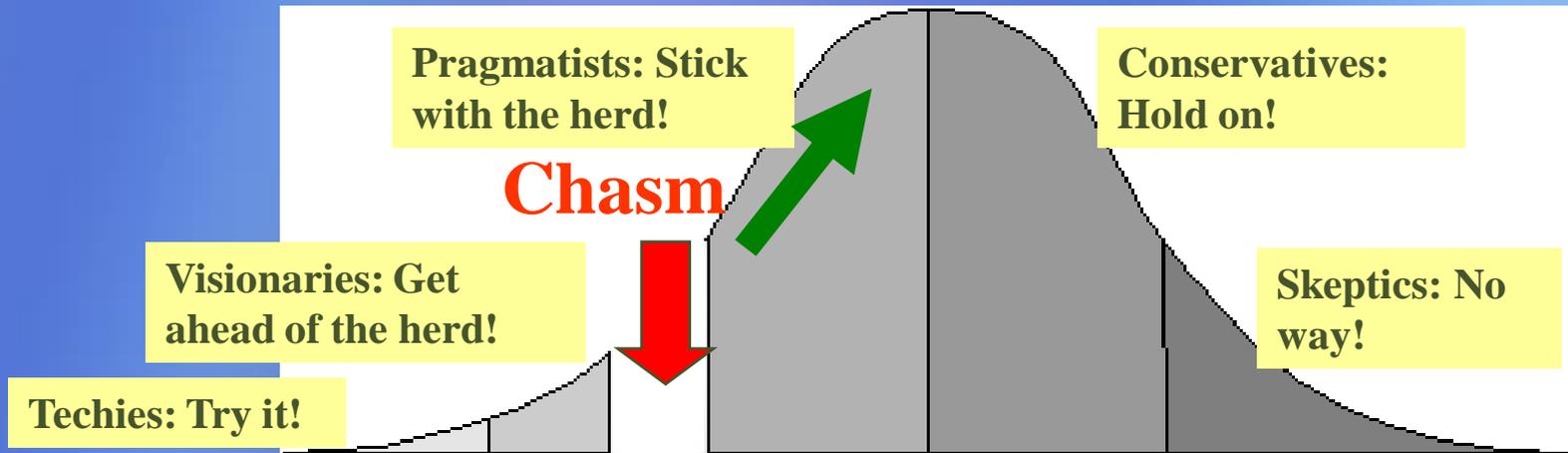


# Outline



- Progress Report
- New frontier

# Imperative Circa 2008



## Maintain upward trajectory:

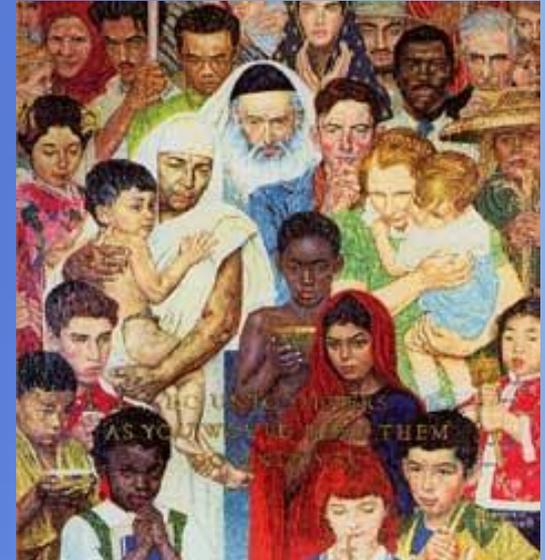
- Focus on a new class of applications, bringing into fold techies and visionaries, leading to new inventions and markets
- While continuing to innovate for the current mainstream market

# Humane Data Mining

“Is it right? Is it just?”

Is it in the interest of mankind?”

*Woodrow Wilson. May 30, 1919.*



## *Applications to Benefit Individuals*

Rooting our future work in this class of new applications, will lead to new abstractions, algorithms, and systems

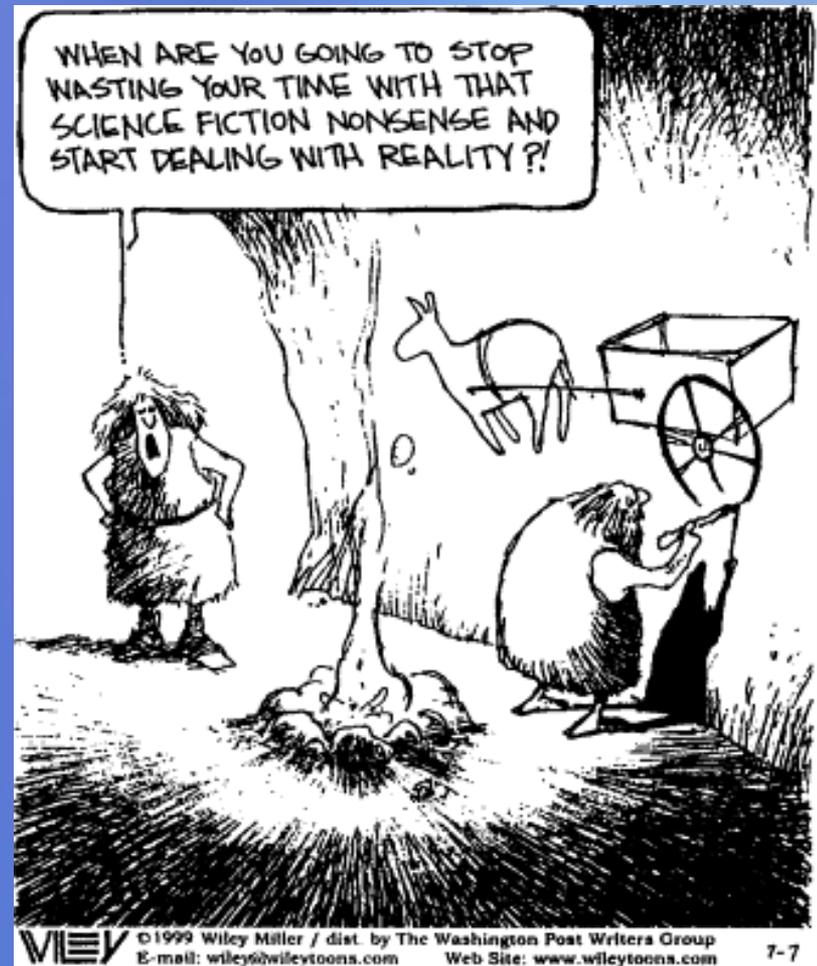
# An Expansive Definition of Data Mining

- Deriving value from a data collection by studying and understanding the structure of the constituent data



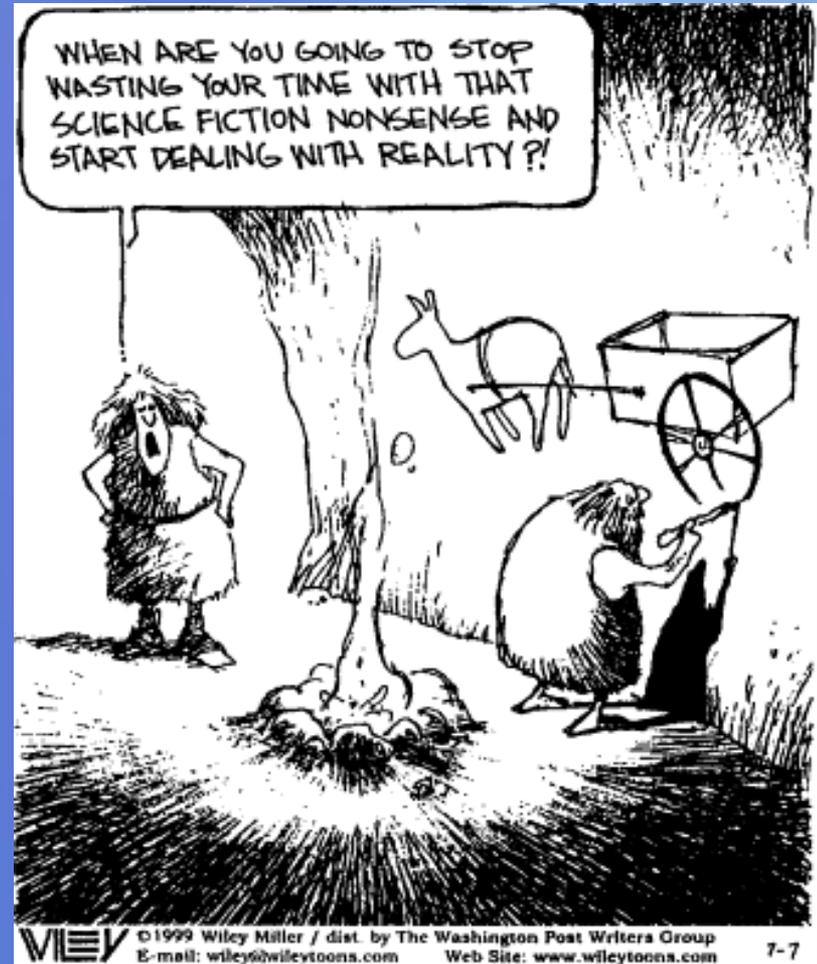
# Some Ideas

- Personal data mining
- Enable people to get a grip on their world
- Enable people to become creative
- Enable people to make contributions to society
- Data-driven science

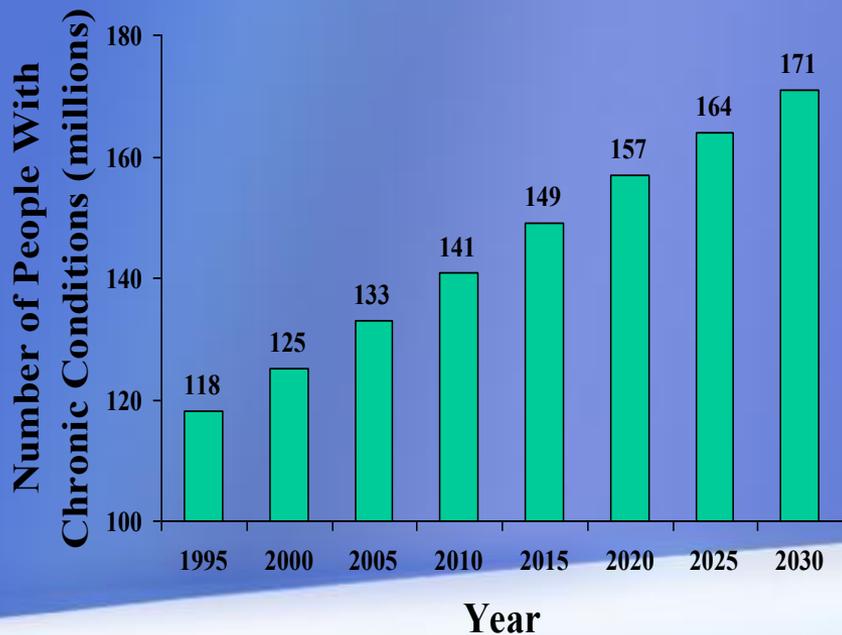
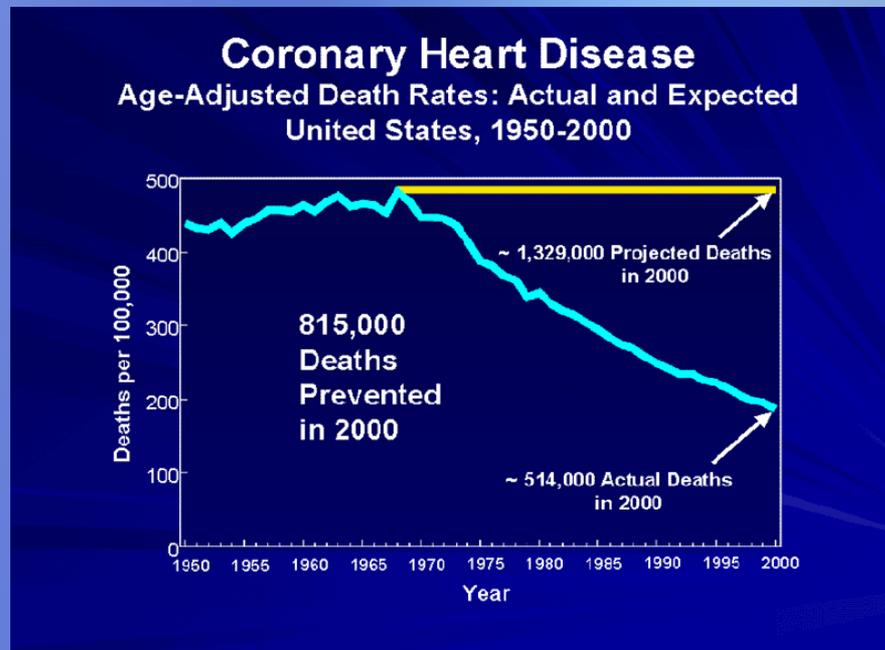
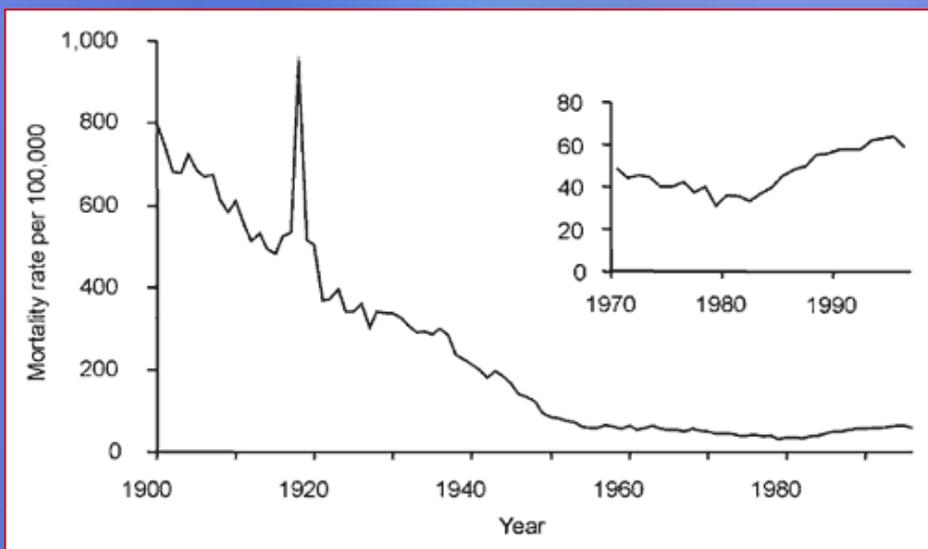


# Some Ideas

- Personal data mining
- Enable people to get a grip on their world
- Enable people to become creative
- Enable people to make contributions to society
- Data-driven science



# Changing Nature of Disease



- **New Challenge: chronic conditions:** illnesses and impairments expected to last a year or more, limit what one can do and may require ongoing care.
- In 2005, 133 million Americans lived with a chronic condition (up from 118 million in 1995).

# Technology Trends

- Tremendous simplification in the technologies for capturing useful personal information

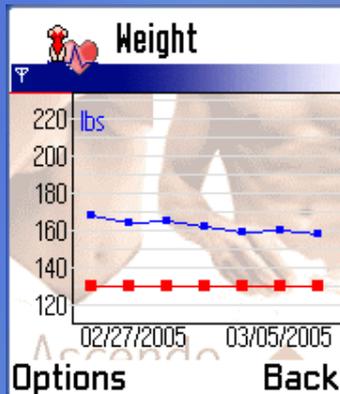
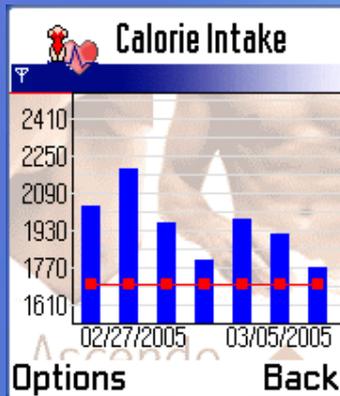


- Dramatic reduction in the cost and form factor for personal storage



- Cloud Computing

# Personal Health Analytics



**Fitness Assessment on 07/09/2000**

Client: **Paul Moore**

No of Prev Tests: **1**      Previous Tests: Last Test **14 Jun 2000**

**V02 - Bike Test Without Watts Meter**

	Actual	Predicted	Previous	L/min	Actual	Predicted	Previous
Load Kg	<input type="text" value="300"/>		N/A	Aerobic	<input type="text" value="3.58"/>		<input type="text" value="3.58"/>
Rpm	<input type="text" value="50"/>		N/A	ml/Kg/min	<input type="text" value="44"/>	<input type="text" value="33"/>	<input type="text" value="44"/>
Watts	<input type="text" value="150"/>	<input type="text" value="166"/>	N/A	V02			
		70%    90%		V02 Comment	<input type="text" value="Above Average V02"/>		
End	<input type="text" value="140"/>	<input type="text" value="128"/>	<input type="text" value="165"/>	<input type="text" value="140"/>			

Graph Options:  
 Aerobic Capacity     V02     Skip this Test

Navigation:

**Enter the Load in Kgs.**

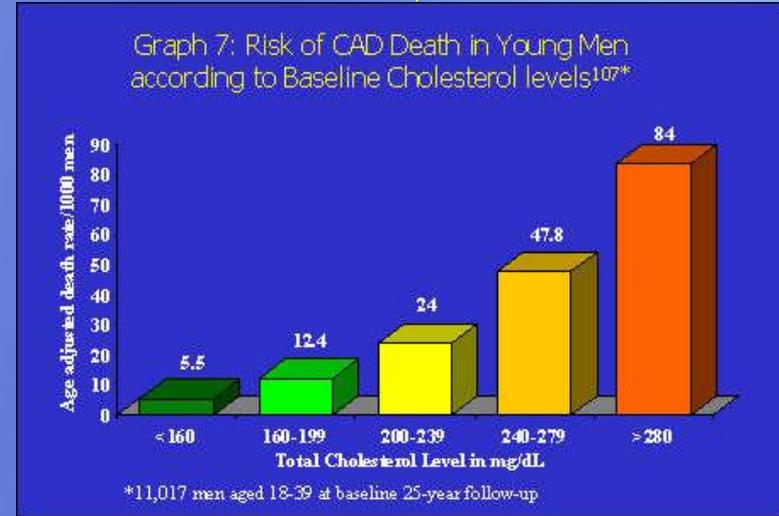
# Personal Data Mining



Charts for appropriate demographics?

Optimum level for Asian Indians: 150 mg/dL (much lower than 200 mg/dL for Westerners) Due to elevated levels of lipoprotein(a)\*

Distributed computation and selection across millions of nodes  
Privacy and security

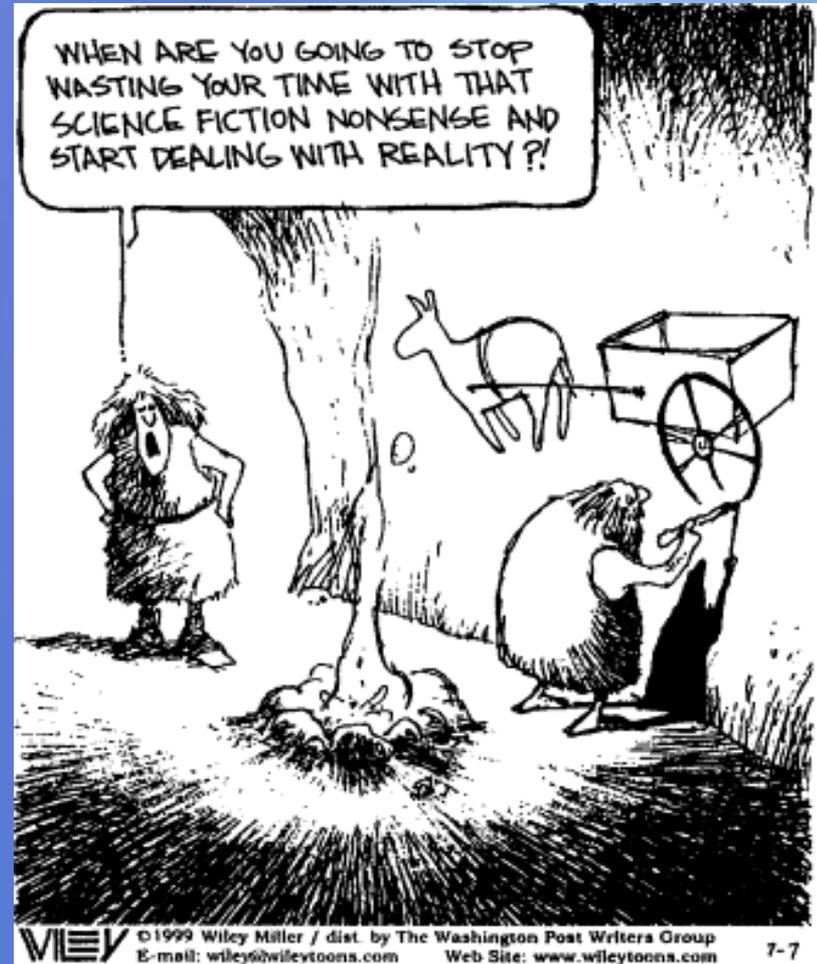


Microsoft  
Research

\*Enas et al. Coronary Artery Disease In Asian Indians. *Internet J. Cardiology*. 2001.

# Some Ideas

- Personal data mining
- Enable people to get a grip on their world
- Enable people to become creative
- Enable people to make contributions to society
- Data-driven science



# The Tyranny of Choice

## ANATOMY OF THE LONG TAIL

Online services carry far more inventory than traditional retailers. Rhapsody, for example, offers 19 times as many songs as Wal-Mart's stock of 39,000 tunes. The appetite for Rhapsody's more obscure tunes (charted below in yellow) makes up the so-called Long Tail. Meanwhile, even as consumers flock to mainstream books, music, and films (right), there is real demand for niche fare found only online.

### RHAPSODY

TOTAL INVENTORY:  
735,000 songs



### AMAZON.COM

TOTAL INVENTORY:  
2.3 million books



### NETFLIX

TOTAL INVENTORY:  
25,000 DVDs

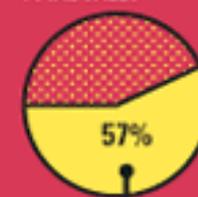


**THE NEW GROWTH MARKET:**  
**OBSCURE PRODUCTS YOU CAN'T GET ANYWHERE BUT ONLINE**

TOTAL SALES



TOTAL SALES



TOTAL SALES



product not available in offline retail stores



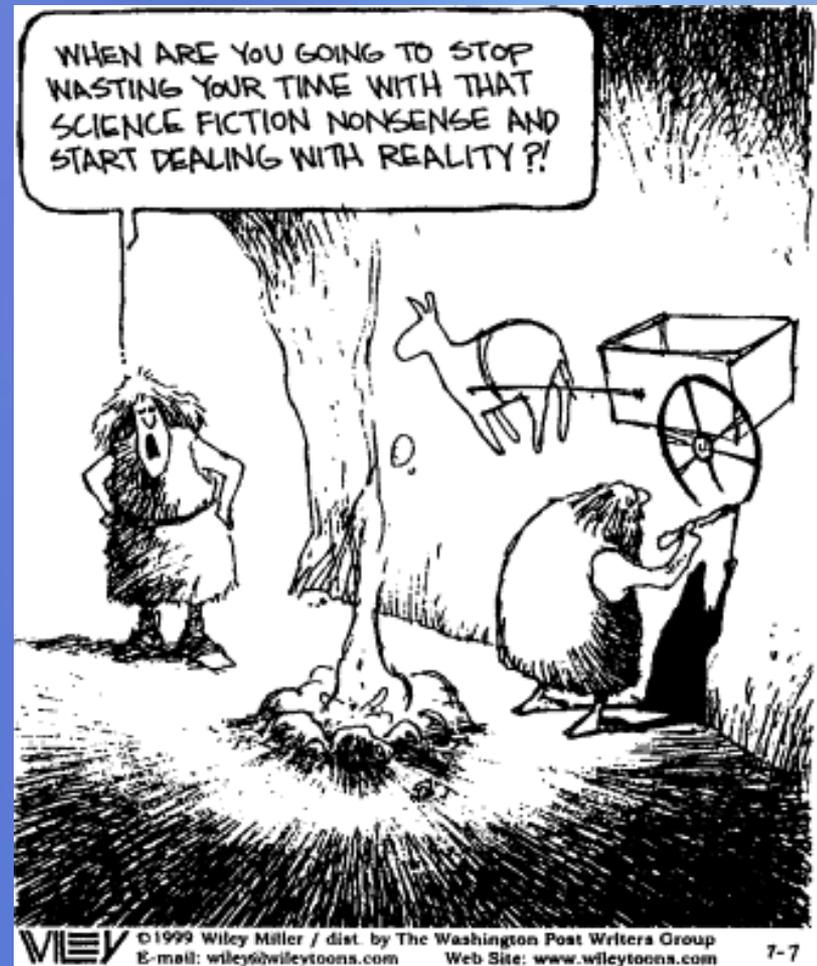
Sources: Erik Brynjolfsson and Jeffrey Hu, MIT, and Michael Smith, Carnegie Mellon; Barnes & Noble; Netflix; RealNetworks

Chris Anderson. The Long Tail. 2006.

See also: Anita Elberse. Should You Invest in the Long Tail? HBR, 2008.

# Some Ideas

- Personal data mining
- Enable people to get a grip on their world
- **Enable people to become creative**
- Enable people to make contributions to society
- Data-driven science



# Tools to Aid Creativity



- Bawden’s four kinds of information to aid creativity: Interdisciplinary, peripheral, speculative, exceptions and inconsistencies
- Intriguing work of Prof Swanson: Linking “non-interacting” literature
  - $L_1$ : Dietary fish oils lead to certain blood and vascular changes
  - $L_2$ : Similar changes benefit patients with Raynaud's syndrome,  $L_1 \cap L_2 = \phi$ .
  - Corroborated by a clinical test at Albany Medical College
  - Similarly, magnesium deficiency & Migraine (11 factors) ; corroborated by eight studies.
- Will we provide the tools?

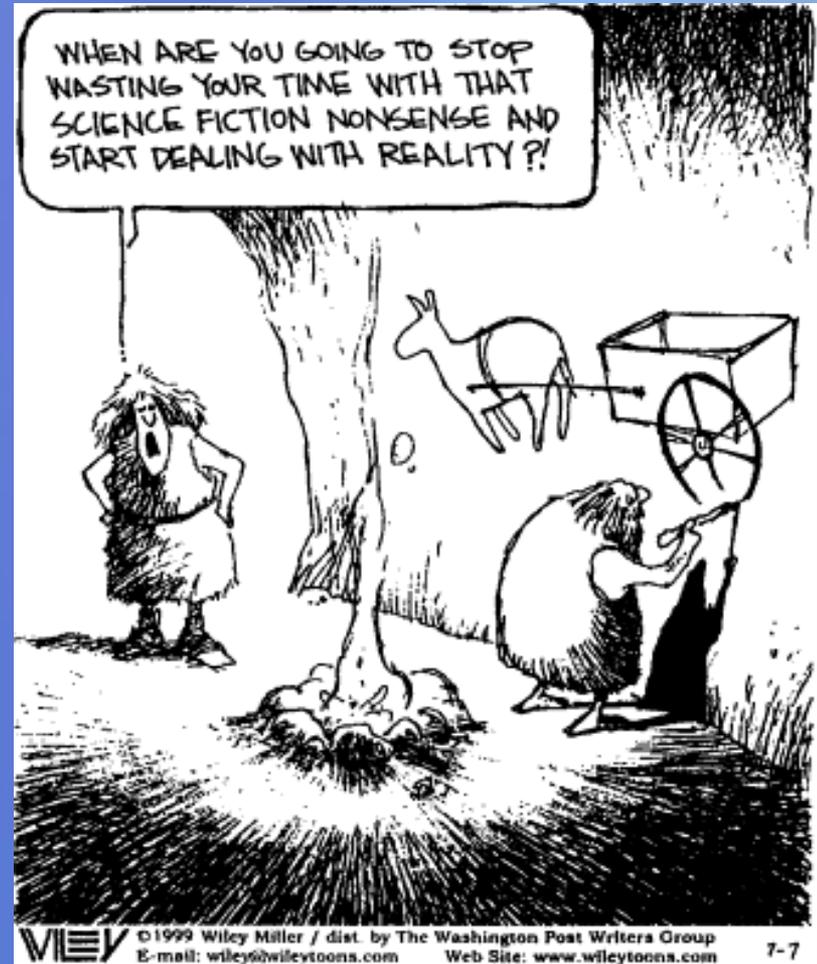


Bawden. “Information systems and the stimulation of the creativity”. Information Science 86.

Swanson. “Medical literature as a potential source of new knowledge”. Bull Med Libr Assoc. 90

# Some Ideas

- Personal data mining
- Enable people to get a grip on their world
- Enable people to become creative
- Enable people to make contributions to society
- Data-driven science

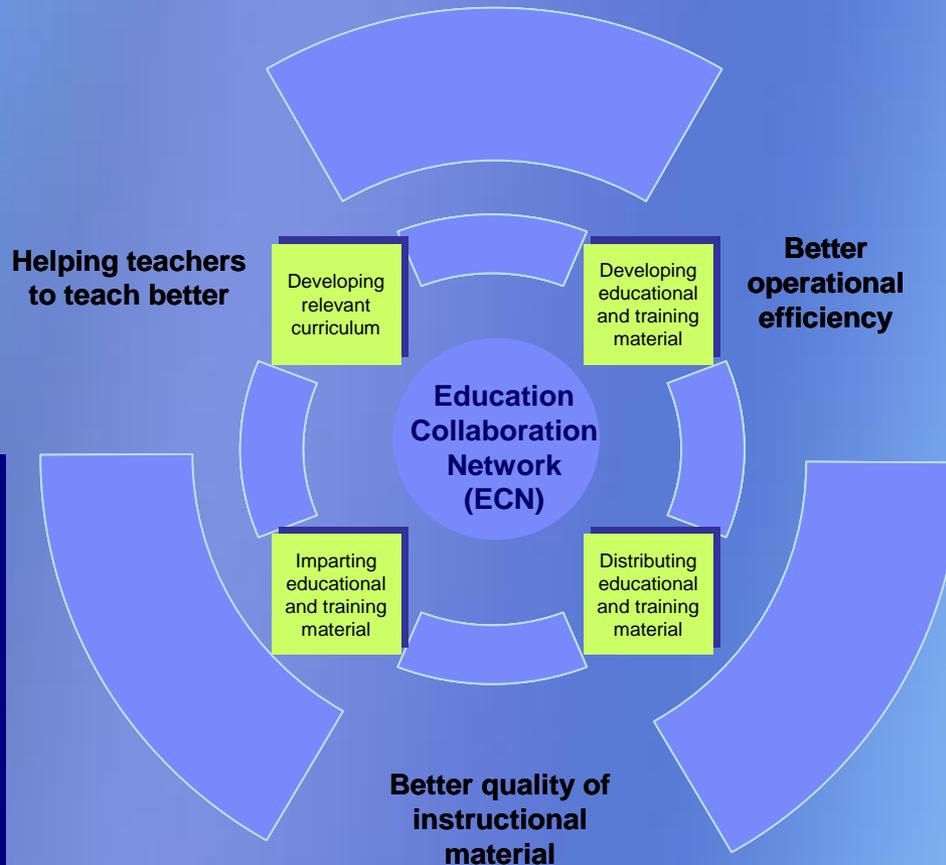


# Education Collaboration Network



- Low teacher-student ratios
- Instruction material poor and often out-of-date
- Poorly trained teachers
- High student drop-out rates

- A hardware and a software infrastructure built on industry standards that empower teachers, educators, and administrators to collectively create, manage, and access educational material, impart education, and increase their skills



- Accumulation and re-use of teaching material
- Distributed, evolutionary content creation
- New pedagogy: teacher as discussant
- Multi-lingual

- Teachers are able to find material that help them understand the subject matter and obtain access to teaching aids that others have found useful.
- Teachers also enhance the material with their own contributions that are then available to others on the network.
- Experts come to the classroom virtually



Improving India's Education System through Information Technology.  
 IBM Report to the President of India. 2005.

# Collaborative Knowledge Creation (Educational Material)

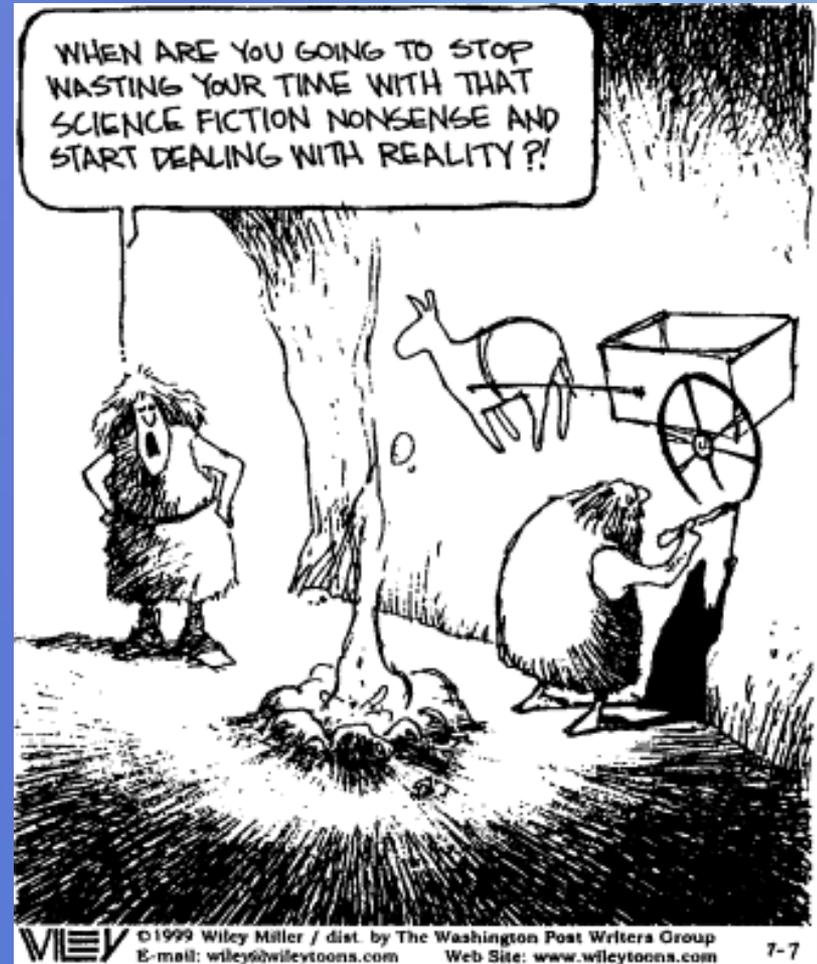


- Inspired by Wikipedia
- But multiple viewpoints rather than one consensus version!
- How to personalize search to find the material suitable for one's own style of teaching?
- Management of trust and authoritativeness?

- More than 3.5 million articles in 75 languages
- Fashioned by more than 25,000 writers
- 1 million articles in English (80,000 in Encyclopedia Britannica)

# Some Ideas

- Personal data mining
- Enable people to get a grip on their world
- Enable people to become creative
- Enable people to make contributions to society
- **Data-driven science**

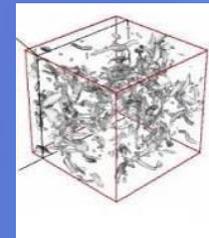


# Science Paradigms

- Thousand years ago:  
science was **empirical**  
describing natural phenomena
- Last few hundred years:  
**theoretical** branch  
using models, generalizations
- Last few decades:  
a **computational** branch  
simulating complex phenomena
- Today:  
**data exploration** (eScience)  
unify theory, experiment, and simulation  
using data management and statistics
  - Data captured by instruments  
Or generated by simulator
  - Processed by software
  - Scientist analyzes database / files



$$\left(\frac{\dot{a}}{a}\right)^2 = \frac{4\pi G\rho}{3} - K\frac{c^2}{a^2}$$



- Historically,  
Computational Science  
= simulation.
- New emphasis on  
informatics:
  - Capturing,
  - Organizing,
  - Summarizing,
  - Analyzing,
  - Visualizing



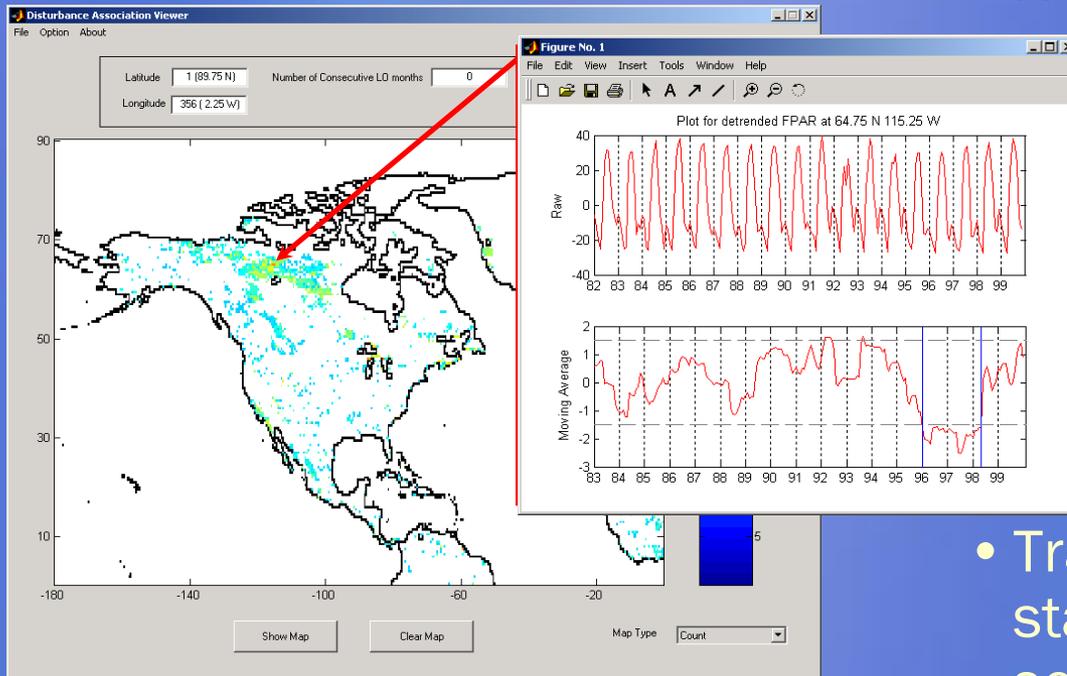
Microsoft  
**Research**

Courtesy Jim Gray, Microsoft Research.

# Understanding Ecosystem Disturbances



Vipin Kumar  
U. Minnesota



- Watch for changes in the amount of absorption of sunlight by green plants to look for ecological disasters

- NASA satellite data to study
  - How is the global Earth system changing?
  - How does Earth system respond to natural & human-induced changes?
  - What are the consequences of changes in the Earth system?
- Transformation of a non-stationary time series to a sequence of disturbance events; association analysis of disturbance regimes



Potter et al. "Recent History of Large-Scale Ecosystem Disturbances in North America Derived from the AVHRR Satellite Record", *Ecosystems*, 2005.

# Some Other Data-Driven Science Efforts

- **Bioinformatics Research Network**



- Study brain disorders and obtain better statistics on the morphology of disease processes by standardizing and cross-correlating data from many different imaging systems
- 100 TB/year

- **Earthscope**



- Study the structure and ongoing deformation of the North American continent by obtaining data from a network of multi-purpose geophysical instruments and observatories
- 40 TB/year



# Call to Action



- Move the focus of future work towards humane data mining (applications to benefit individuals):
  - Personal data mining (e.g. personal health)
  - Enable people to get a grip on their world (e.g. dealing with the long tail of search)
  - Enable people to become creative (e.g. inventions arising from linking non-interacting scientific literature)
  - Enable people to make contributions to society (e.g. education collaboration networks)
  - Data-driven science (e.g. study ecological disasters, brain disorders)
- Rooting our work in these (and similar) applications, will lead to new abstractions, algorithms, and systems

# Thank you!



Search Labs' mission is to invent next in Internet search and applications.