# Vision, Resources, and Opportunities for Next Generation of Data Mining and Cyber-Enabled Discovery and Innovation

**Maria Zemankova, Program Director**
**mzemanko@nsf.gov**

**Information Integration & Informatics Cluster (III)**
**Sylvia Spengler, Cluster Lead**

**Information & Intelligent Information Division (IIS)**
**Haym Hirsh, Division Director**

**Computer & Information Science & Engineering Directorate (CISE)**
**Jeannette Wing, Assistant Director**

**National Science Foundation (NSF)**
**Arden Bement, Director**

**National Science Foundation**
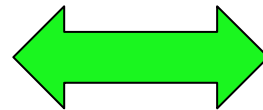WHERE DISCOVERIES BEGIN

# "We"

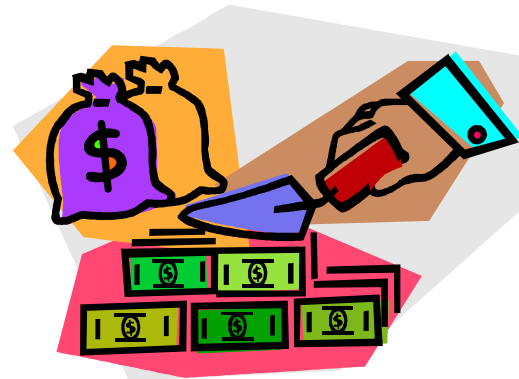**You: researchers**
Academia
Industry
Government labs
International

**Us: supporters**
NSF, NASA, NIH, DOE,…
Industry
Private funders
International

# What "us" need from you

- Research challenges
- Infrastructure needs
- Recommendations for partnerships
- Innovative proposals
  - Research
  - Workshops
  - Educational activities
- Dedication to education, outreach
- Taking broader impacts seriously
- Participation in review process
- Examples of successful research
- Come to NSF as program director, division director,…

# Examples of Successful Research

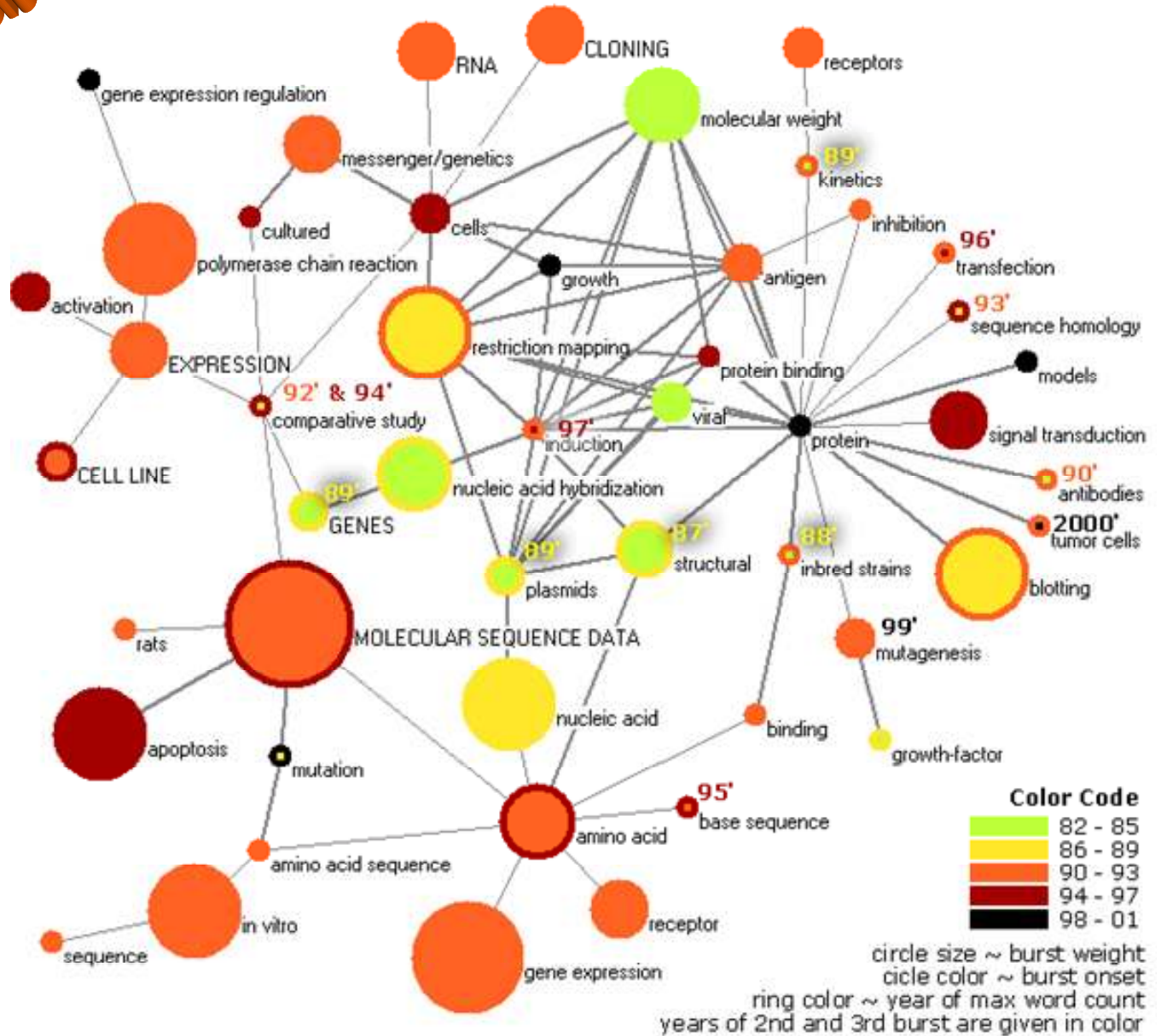- Inform the public

- Justify funds spent

- Request more funding

Katy Borner, Indiana University
*CAREER: Visualizing Knowledge Domains*

**Topic Bursts**

**Mapping**

Visualization of keywords appearing in
PNAS 1982 -- 2001.

The top 10% of highly cited PNAS papers
were analyzed to reveal keywords listed
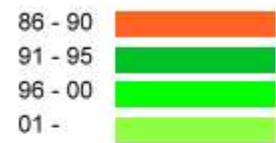with high frequency or rapid initial
appearance (Topic burst).



gene expression regulation
RNA
CLONING
receptors
molecular weight
messenger/genetics
89' kinetics
cultured
cells
inhibition
polymerase chain reaction
growth
antigen
96' transfection
activation
restriction mapping
93' sequence homology
EXPRESSION
protein binding
models
92' & 94'
comparative study
97' induction
viral
protein
signal transduction
CELL LINE
nucleic acid hybridization
90' antibodies
89' GENES
89'
97' structural
88' inbred strains
2000' tumor cells
plasmids
blotting
rats
MOLECULAR SEQUENCE DATA
99' mutagenesis
apoptosis
nucleic acid
binding
growth-factor
mutation
95' base sequence
amino acid
amino acid sequence
in vitro
receptor
sequence
gene expression

**Color Code**
82 - 85
86 - 89
90 - 93
94 - 97
98 - 01

circle size ~ burst weight
cicle color ~ burst onset
ring color ~ year of max word count
years of 2nd and 3rd burst are given in color

Ke, Visvanath & Börner
*Mapping the Evolution of Co-Authorship Networks*
(Won 1st price at the IEEE InfoVis Contest, 2004)

1990

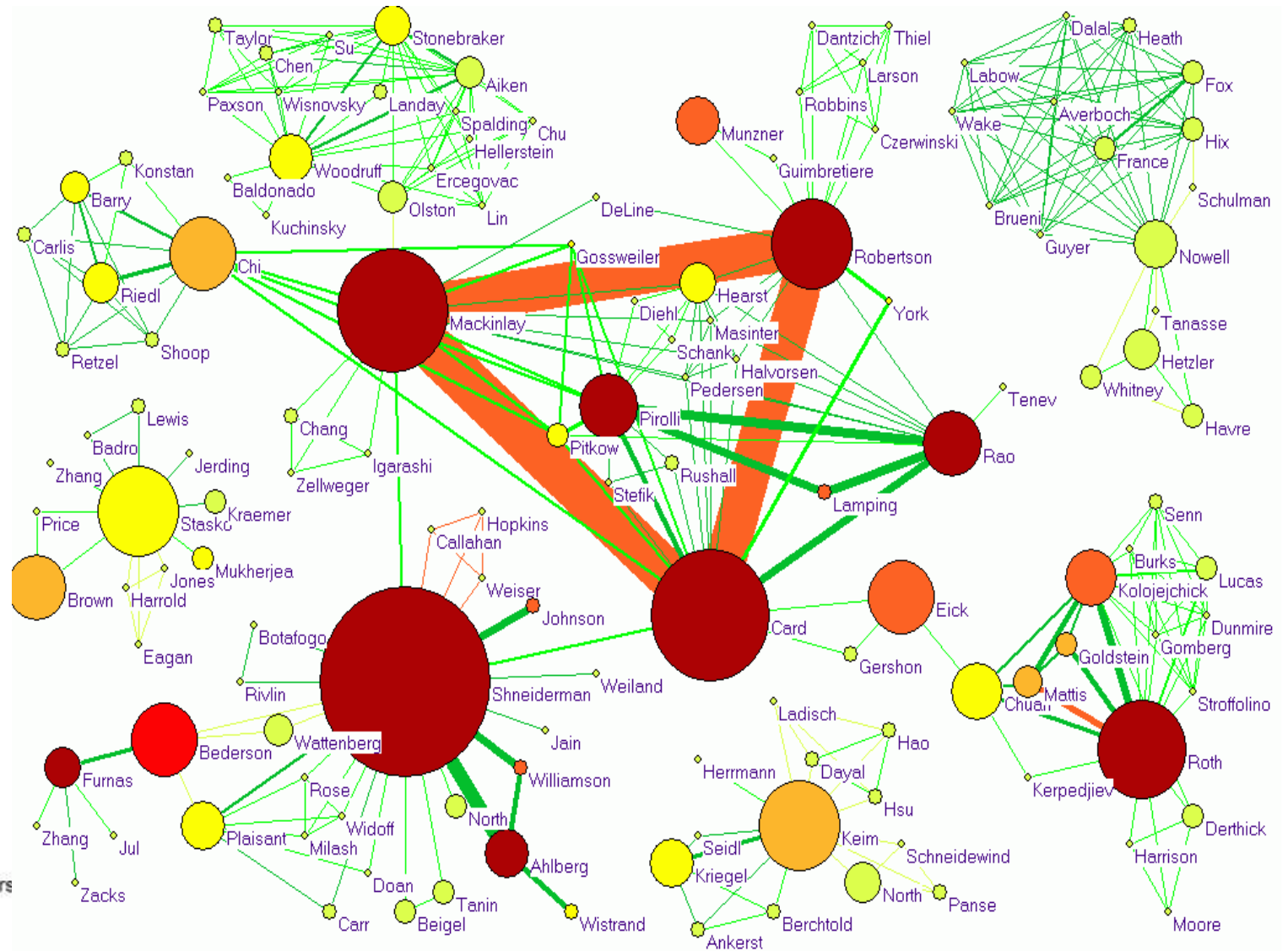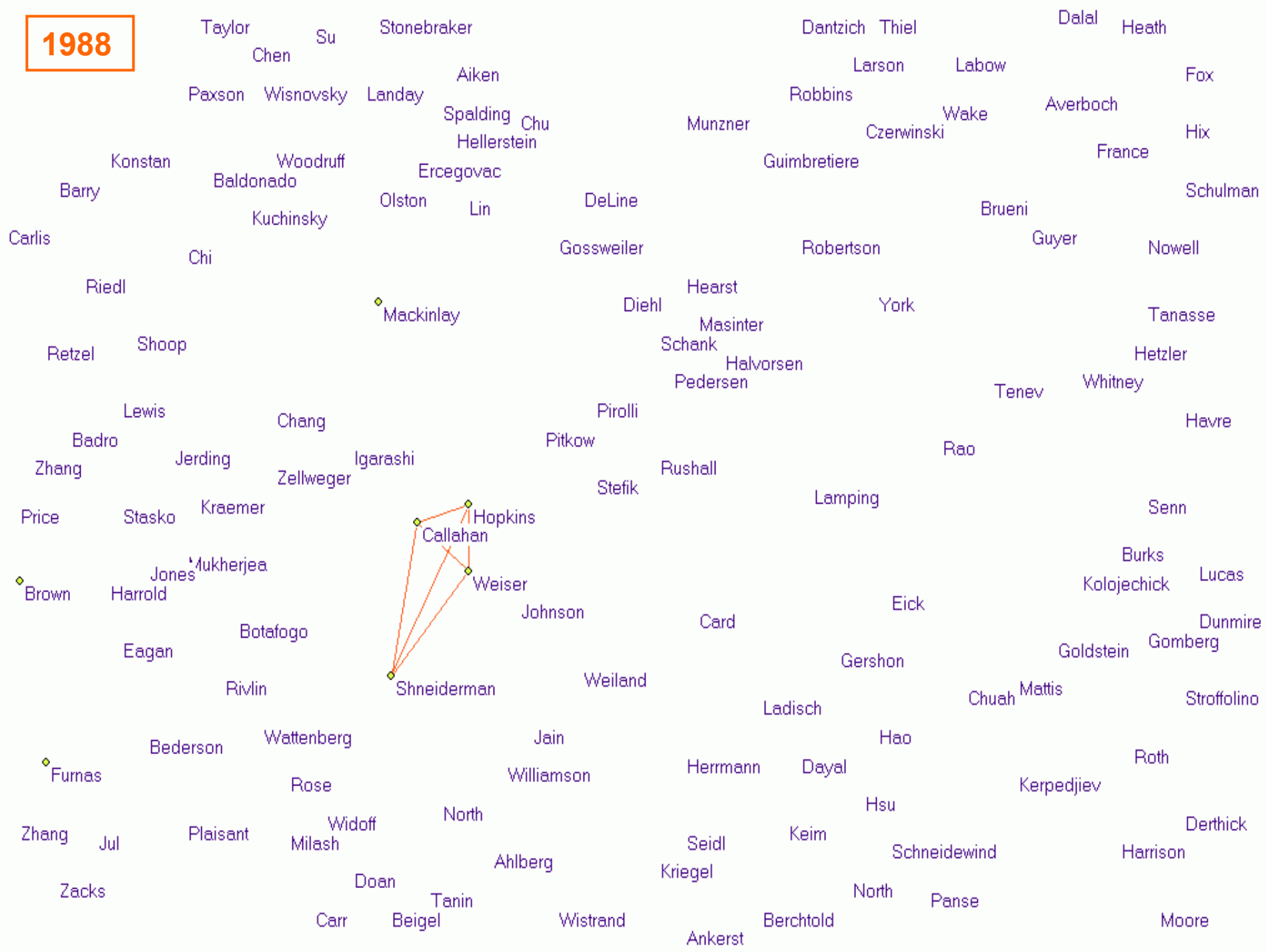Taylor   Chen   Su   Stonebraker         Dantzich   Thiel      Dalal   Heath

Aiken                Larson   Labow        Fox

Paxson   Wisnovsky   Landay        Robbins              Averboch

Spalding   Chu                    Wake              Hix

Hellerstein              Munzner   Czerwinski        France

Konstan   Woodruff        Ercegovac         Guimbretiere           Schulman

Baldonado                                           Brueni

Barry                Olston   Lin      DeLine

Kuchinsky                                      Guyer      Nowell

Carlis

Chi              Gossweiler              Robertson

Riedl                         Hearst         York

Mackinlay       Diehl              Tanasse

Masinter

Retzel   Shoop                Schank              Hetzler

Halvorsen

Pedersen        Whitney

Lewis                         Tenev

Chang              Pirolli           Rao           Havre

Badro         Jerding   Igarashi   Pitkow

Zhang         Zellweger              Rushall

Stefik                    Lamping           Senn

Price   Stasko   Kraemer                        Burks

Hopkins                              Kolojechick   Lucas

Callahan

Jones Mukherjea   Weiser                    Eick         Dunmire

Brown                Johnson                         Gomberg

Harrold                                   Goldstein

Card         Gershon

Eagan                                          Mattis

Botafogo                                  Chuah        Stroffolino

Rivlin   Shneiderman   Weiland            Ladisch            Roth

Hao              Kerpedjiev

Bederson   Wattenberg      Jain

Furnas         Williamson   Herrmann   Dayal

Hsu                Derthick

Rose            North              Keim

Zhang   Jul   Plaisant   Widoff              Seidl         Schneidewind   Harrison

Milash              Ahlberg      Kriegel      North   Panse

Doan                                          Moore

Zacks      Tanin

Carr   Beigel   Wistrand      Berchtold

Ankerst

1994

1996

Taylor  Stonebraker
Su
Chen  Aiken
Paxson  Wisnovsky  Landay
Woodruff
Spalding  Chu
Hellerstein
Ercegovac
Olston  Lin

Dantzich  Thiel
Larson
Robbins
Munzner
Guimbretiere

Dalal  Heath
Labow  Fox
Wake  Averboch
Czerwinski  Hix
France
Brueni
Guyer  Nowell
Schulman

Konstan
Barry
Carlis  Chi
Riedl
Retzel  Shoop

Baldonado
Kuchinsky

DeLine
Gossweiler
Mackinlay
Hearst
Diehl
Masinter
Schank
Halvorsen
Pedersen
Pirolli
Pitkow
Rushall
Stefik
Lamping

Robertson
York

Rao

Tanasse
Hetzler
Whitney
Tenev
Havre

Lewis
Badro  Jerding
Zhang
Stasko  Kraemer
Price
Jones  Mukherjea
Brown  Harrold
Eagan

Chang
Zellweger  Igarashi

Hopkins
Callahan
Weiser
Johnson
Botafogo  Weiland
Rivlin  Shneiderman
Jain
Bederson  Williamson
Furnas  Wattenberg
Rose  North
Zhang  Widoff  Ahlberg
Jul  Plaisant  Milash
Zacks  Doan  Tanin
Carr  Beigel  Wistrand

Card

Eick
Gershon
Ladisch
Hao
Herrmann  Dayal
Hsu
Seidl  Keim
Kriegel
Ankerst  Berchtold
North  Panse

Senn
Burks  Lucas
Kolojechick
Goldstein  Dunmire
Gomberg
Chuah  Mattis  Stroffolino
Roth
Kerpedjiev
Derthick
Harrison
Schneidewind
Moore

1997

1998

1999

Taylor  Stonebraker
Su
Chen
Aiken
Paxson  Wisnovsky  Landay
Spalding  Chu
Hellerstein
Woodruff
Baldonado  Ercegovac
Olston  Lin

Konstan
Barry
Kuchinsky
Carlis
Chi
Riedl
Retzel  Shoop

Dantzich  Thiel
Munzner  Larson
Robbins
Czerwinski
Guimbretiere

Dalal  Heath
Labow  Fox
Wake  Averboch
Brueni  Hix
Guyer  Nowell
France
Schulman

DeLine
Gossweiler  Robertson
Mackinlay  Hearst  York
Diehl
Masinter
Schank
Halvorsen
Pedersen
Pirolli  Rao
Pitkow  Tenev
Rushall
Stefik
Lamping

Tanasse
Hetzler
Whitney
Havre

Lewis
Badro  Jerding
Zhang
Price  Stasko  Kraemer
Jones  Mukherjea
Brown  Harrold

Chang
Zellweger  Igarashi

Hopkins
Callahan
Weiser
Johnson
Weiland
Shneiderman
Jain

Eick
Gershon
Card

Senn
Burks
Kolojechick  Lucas
Goldstein  Dunmire
Chuah  Mattis  Gomberg
Stroffolino

Botafogo
Eagan
Rivlin

Bederson
Furnas
Zhang  Jul
Zacks

Wattenberg
Rose
Plaisant  Widoff
Milash
Doan  Beigel
Carr  Tanin

North
Williamson
Ahlberg
Wistrand

Ladisch

Roth
Kerpedjiev
Derthick
Harrison
Moore

Hao
Herrmann  Dayal
Hsu
Keim
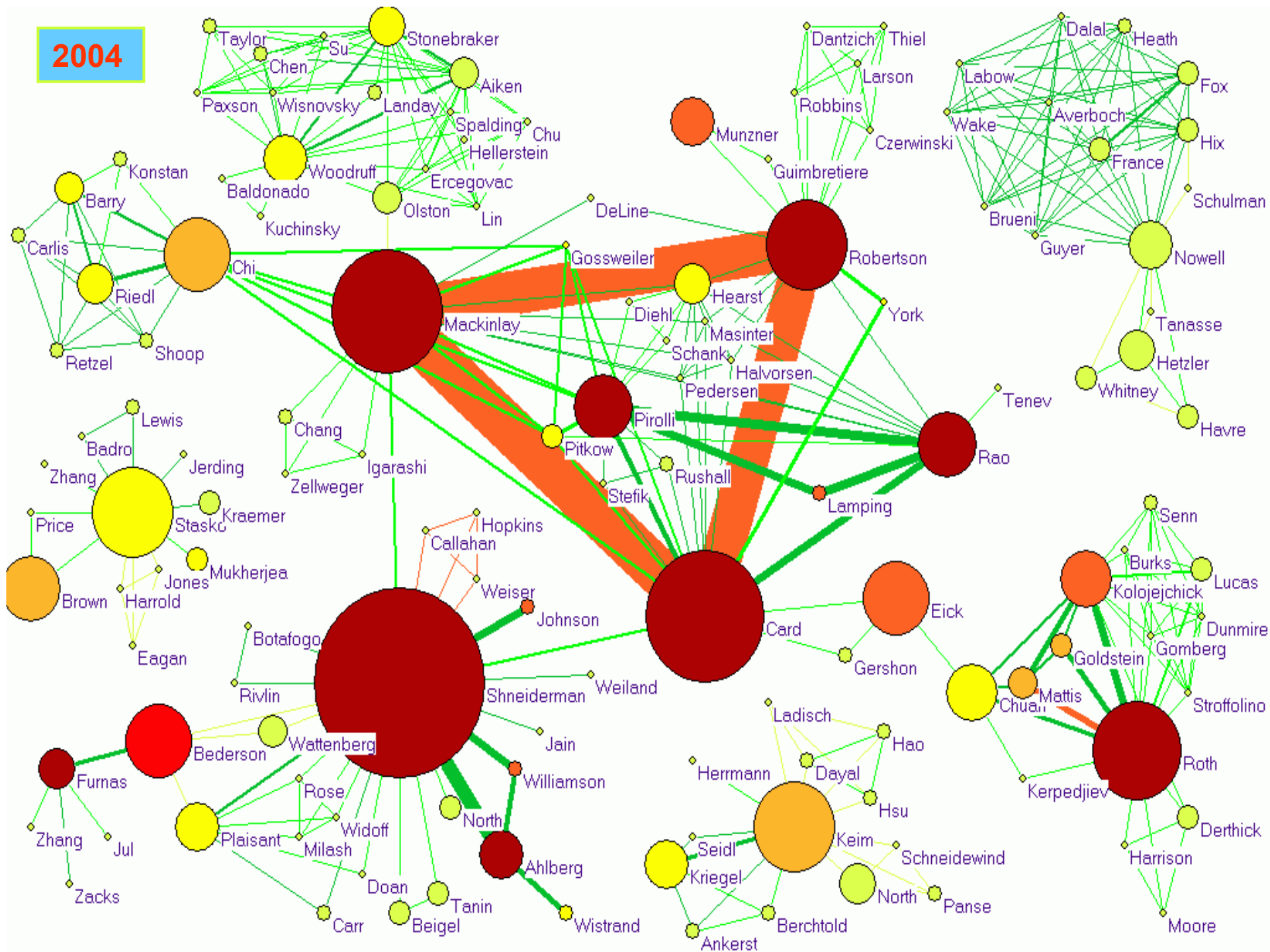Seidl
Kriegel  Schneidewind
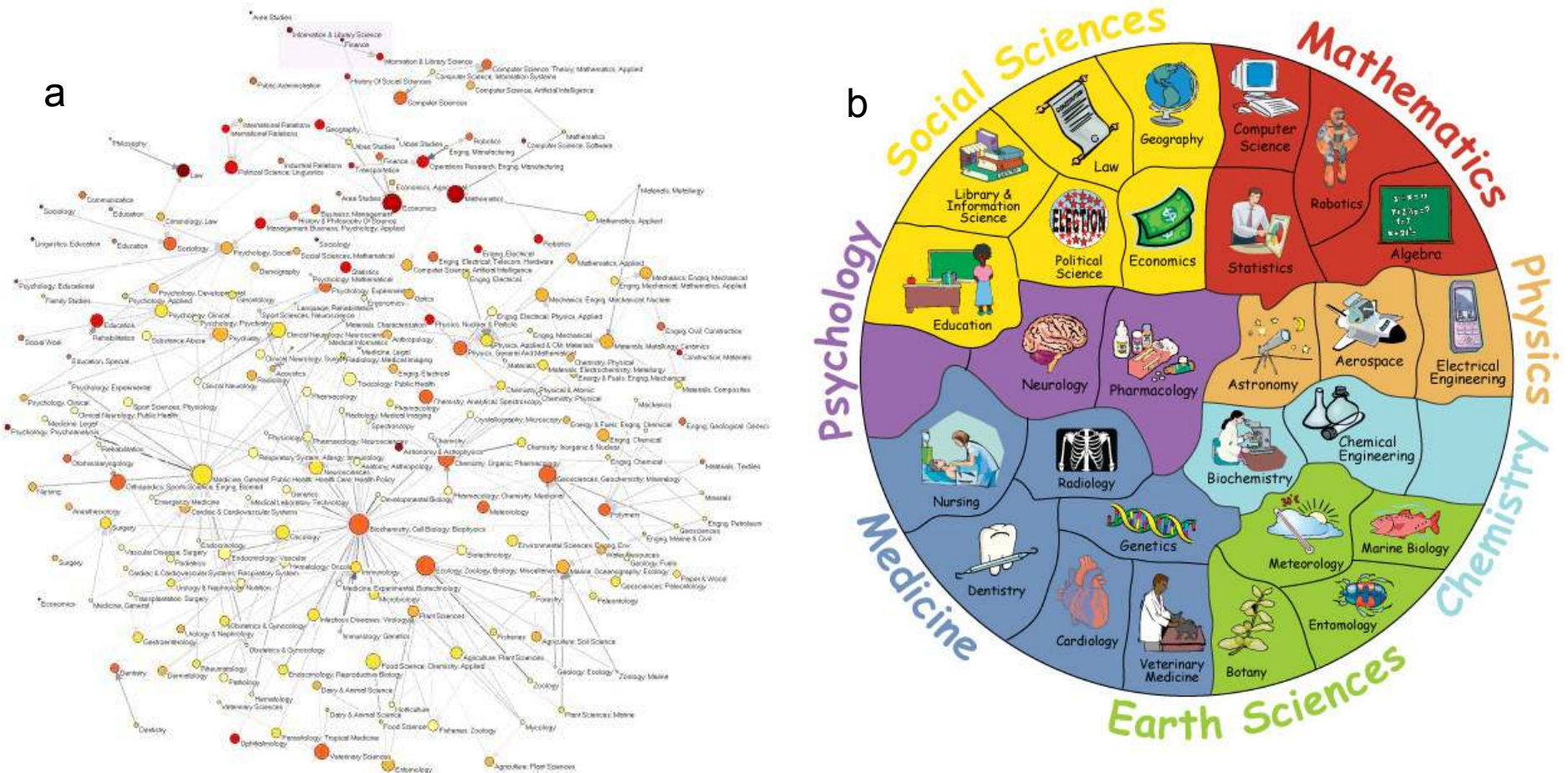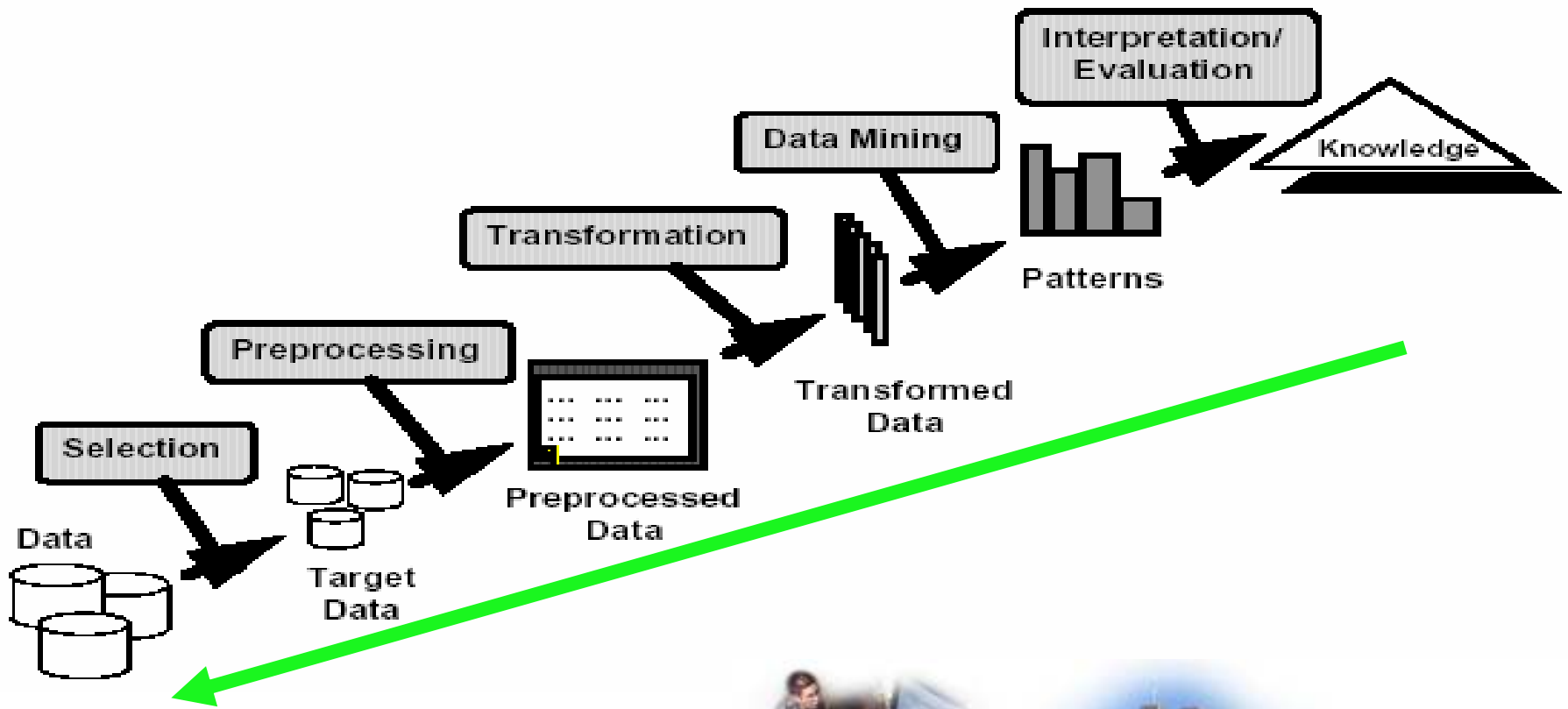Ankerst  Berchtold  North  Panse

2002

2003

Map of Science derived from text mining of volumes of scientific literature.
(a) Visualization for researchers where 2D proximity represents similarity of research areas and size and color of research areas indicate activity level. Interactive map allows drilling to research articles.
(b) Puzzle for elementary school students that represents similarity of research areas.

0308264
**Vipin Kumar and Jaideep Srivastava, University of Minnesota**
**Data Mining for Rare Class Analysis**

- A precursor to many attacks on networks is often a reconnaissance operation, more commonly referred to as a scan. These computers, once compromised, are used to send spam, serve pornography, or launch large scale attacks on the Internet.

- This project developed a new, effective method for scan detection by employing a data mining approach, which makes it possible to use automated techniques for knowledge discovery from massive data.

- Extensive experiments on real network traffic data have shown that the new methods have substantially better performance than the state of the art methods, in terms of coverage, false alarm rate and speed of detection.

- With millions of compromised computers constantly scanning the Internet, such improved scan detection techniques are critical to providing security analysts with the information they need to take preventive measures.

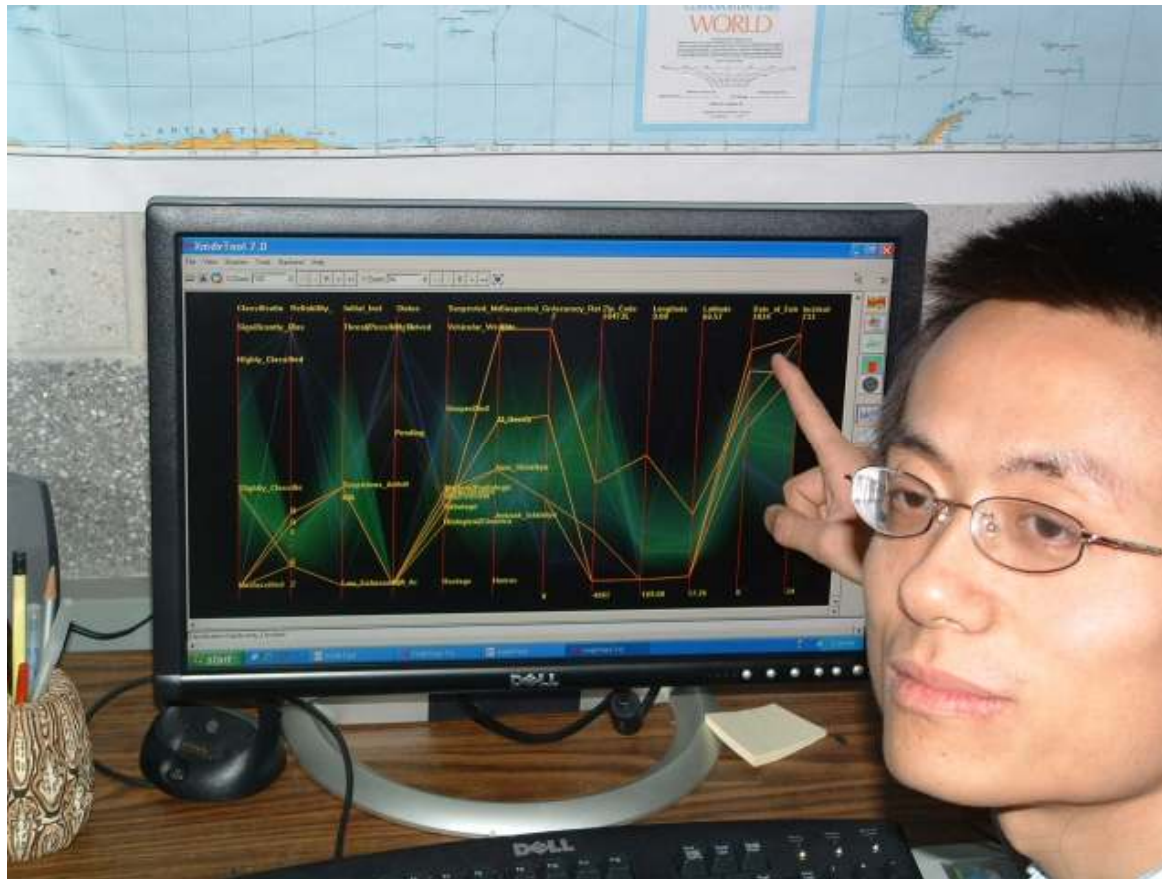**Matthew O. Ward and Elke A. Rundensteiner, WPI**

**Quality-Aware Visual Exploration Tools**

In general, exploratory data analysis assumes data is of high quality and that data transformations, either computational (e.g., data mining) or visualization, do not cause loss of information and result in new knowledge of high quality.

In reality, data is rarely of uniformly high quality or certainty, and every transformation performed distorts or loses information. This disparity leads to results (new knowledge) that may be misleading or incorrect, with potentially serious consequences.

This project adds a new and important dimension to knowledge discovery – making the researcher aware of the quality of the data they are analyzing as well as the information loss resulting from filtering, sampling, clustering, visualizing, and all other transformations applied to the data, resulting in more accurate and reliable knowledge and decisions can be drawn from data.

A data analyst points out high confidence clusters in a simulated homeland security dataset.
Lower quality information has been automatically deemphasized.

# Opportunities at NSF: http://www.nsf.gov

- **MyNSF**
  - Create your profile
  - Receive
    - Announcement on funding opportunities
      - **Explosives and Related Threats: Frontiers in Prediction and Detection (EXP)**
    - Relevant News, Discoveries, reports, …
- **Find Funding**
  - Search by keywords, dates, directorates, …
  - Browse the NSF Web:
    - **NSF-wide/Cross-cutting Programs**
      - Integrative Graduate Education and Research Traineeship Program (IGERT)
      - NSF Graduate Teaching Fellowship in K-12 Education (GK-12)
      - Grant Opportunities for Academic Liaison with Industry (GOALI)
    - **Office of International Science & Engineering**
      - Partnerships for International Research and Education (PIRE)
- **Awards Search**
  - Search by keywords
    - Find researchers, relavant NSF programs, program directors, …
  - Search by NSF Units (divisions, programs, …)
    - Get an idea about the scope

# CISE++ Opportunities:
http://www.nsf.gov/dir/index.jsp?org=CISE

Cyber-Enabled Discovery and Innovation (CDI)

Information & Intelligent Systems Program Solicitation

Expeditions in Computing

Foundations of Data and Visual Analytics

CreativeIT

CyberTrust

Community-Based Data Interoperability Networks

Sustainable Digital Data Preservation and Access Network Partners (DataNet)

Mathematical Sciences: Innovations at the Interface with Computer Sciences

Industry/University Cooperative Research Centers Program (I/UCRC)

CISE Pathways to Revitalized Undergraduate Computing Education (CPATH )

Research Experience for Undergraduates (REU) Sites & Supplements