Marjorie M.K. Hlava, President
Jay Ven Eman, CEO
Access Innovations, Inc.
mhlava@accessinn.com
J_ven_eman@accessinn.com

# *TEXT MINING, TERM MINING, AND VISUALIZATION*

## *IMPROVING THE IMPACT OF SCHOLARLY PUBLISHING*

## MONDAY 16 APRIL 2012
## NICE, FRANCE

# What we will cover today

- Term and Text Mining

- The basics of visualization

- Case studies

- Using subject terms as metrics

- Applications

- Visualizing the results

# Definitions

- *Term Mining - a systematic comparison processing algorithmic method to find patterns in text*

- *Text Mining – using controlled vocabulary tags in text to find patterns and directions*

- Term & text mining
  - Many similarities
  - Can be complimentary; not mutually exclusive

# Term mining

- Precise

- Meaningful semantic relationships; contextual

- Replicable; repeatable; consistent

- Vetted; controlled

- Based on a controlled vocabulary

- Trends; gaps; relationship analysis; visualizations

- Less data processing load

# Text mining

- Algorithmic; formulaic

- Neural nets, statistical, latent semantic, co - occurrence

- Serendipitous relationships

- Sentiment; hot topics; trends

- False drops; noise;

- Misleading semantic relationships
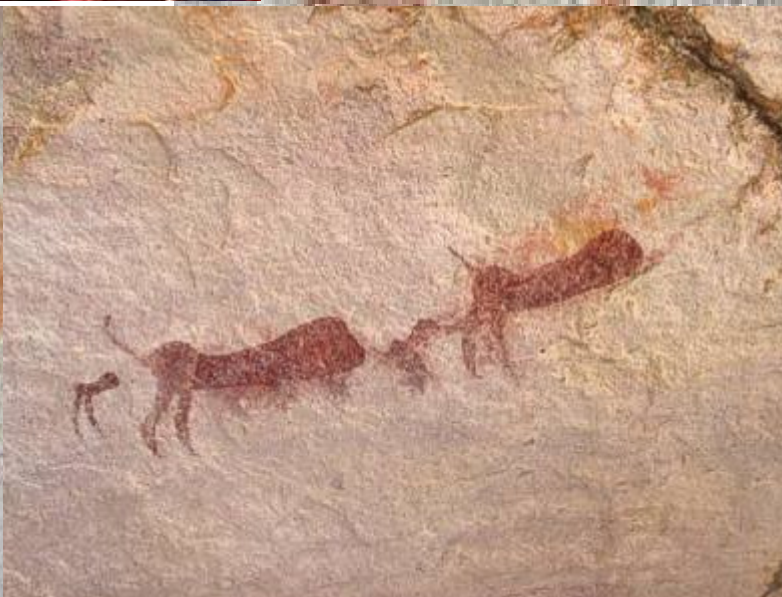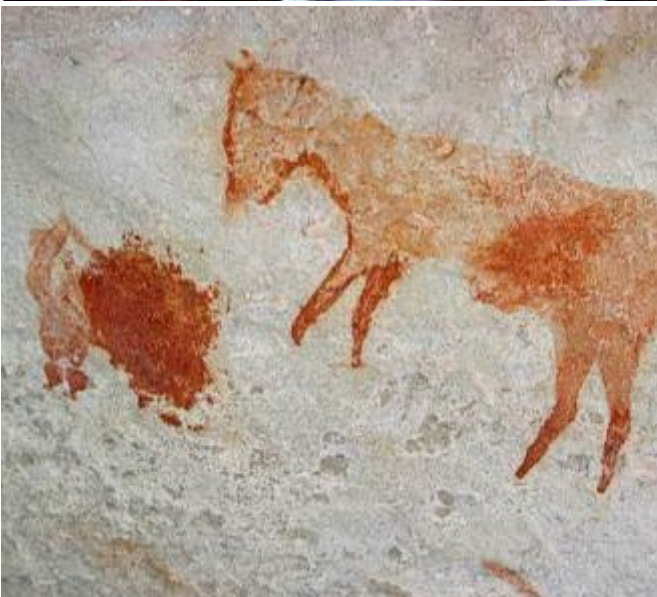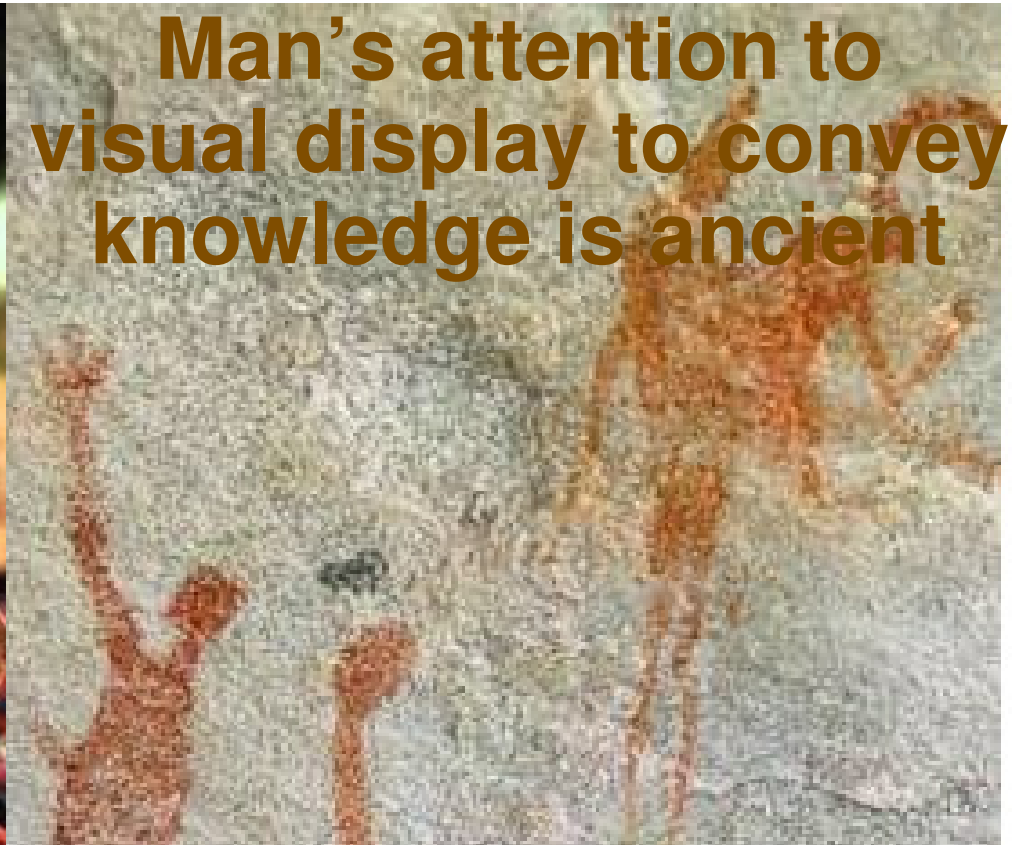
- Heavy processing load

# Why take a visual look?

- Humans can process information 17 times faster in visual presentations

- Now  data can be analyzed, manipulated and presented as visual displays.

-  To see the trends effectively we need to make the data into rich graph-able formats

# Visualization of data

- Needs
  - Measurement
  - Metrics
  - Numbers
- Shows
  - Adjacency
  - Relationships
  - Trends
  - Co – occurrence
  - Conceptual distance

- Is richer with
  - Linking
  - Semantic enrichment
  - Classification

- Supports
  - Forecasting
  - Trend analysis
  - Segmentation
  - Distribution

**Man's attention to visual display to convey knowledge is ancient**
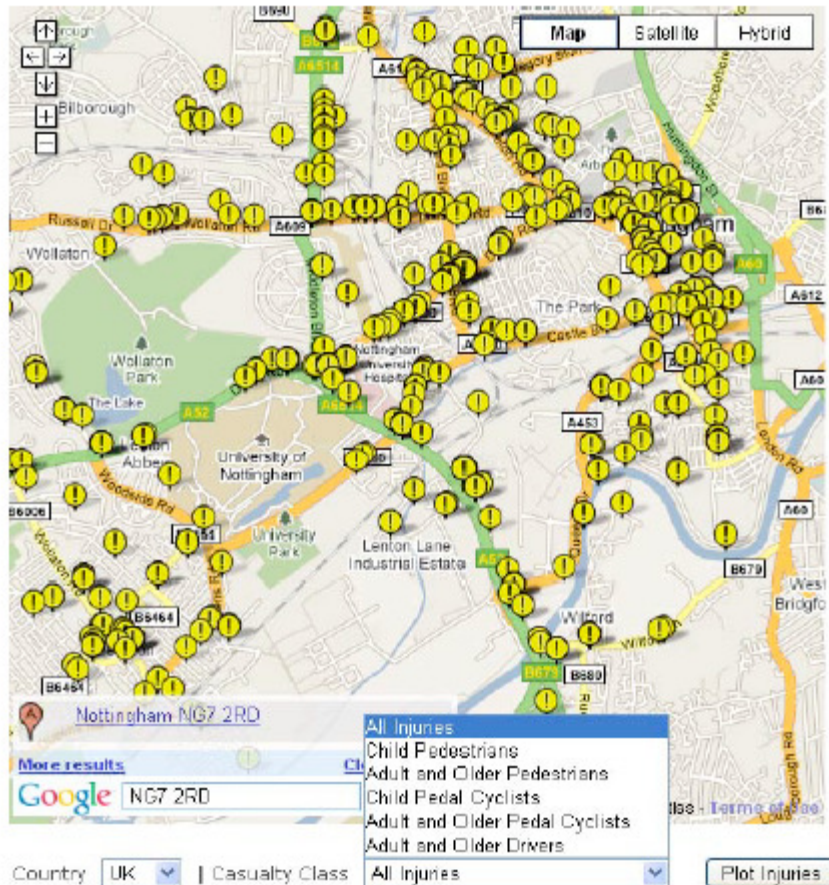
# The art in maps is a longstanding tradition



Credits: Raremaps

# Super imposing data is now common
# A mash up example



Figure 1: *Traffic Injury Map showing incidents in the Nottingham area*

**Traffic Injury Map**
> UK Data Archive
> US National Highway
> > Safety Administration
> Google Maps Base

**Accident categories include**
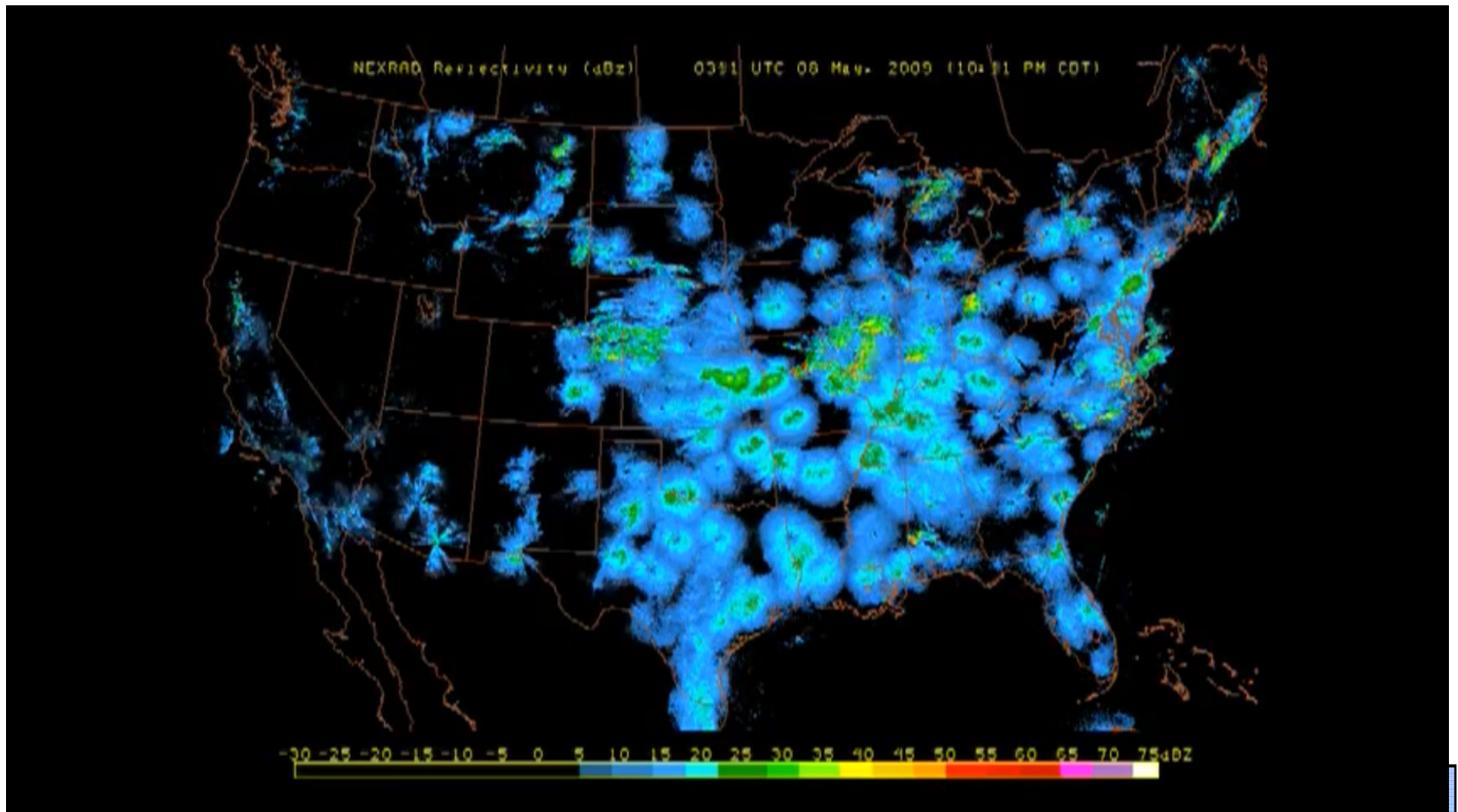> children
> automobile
> bicycle
> etc.

**Data**

> time
> place
> type

Source:
JISC TechWatch: Data Mash-ups September 2010

# Mash up of bird flight migrations and weather patterns

http://www.youtube.com/watch?v=uPff1t4pXiI&feature=youtu.be

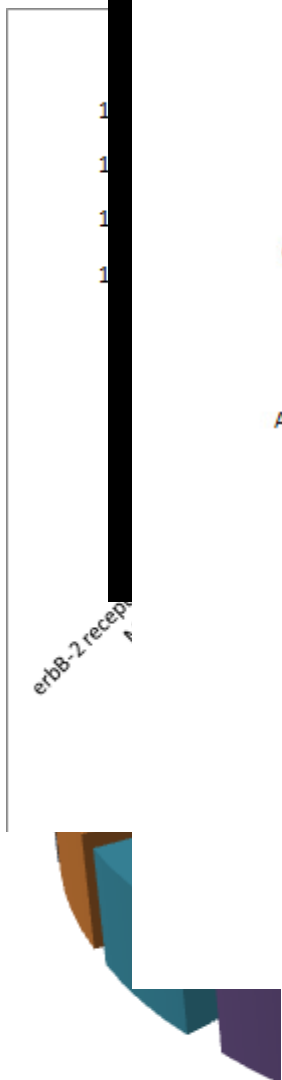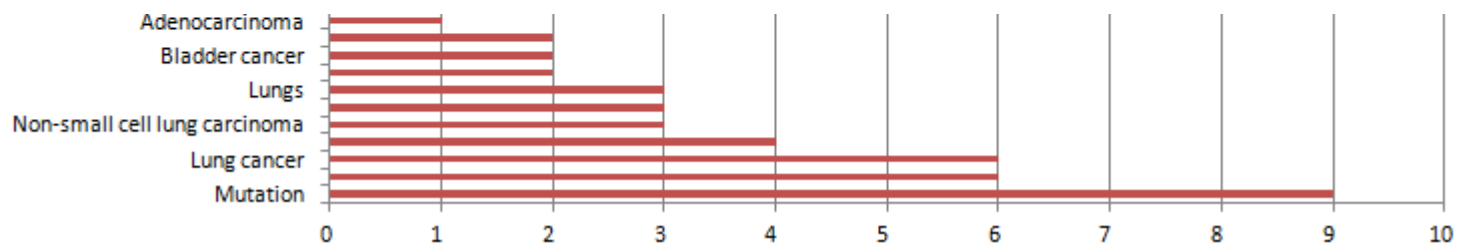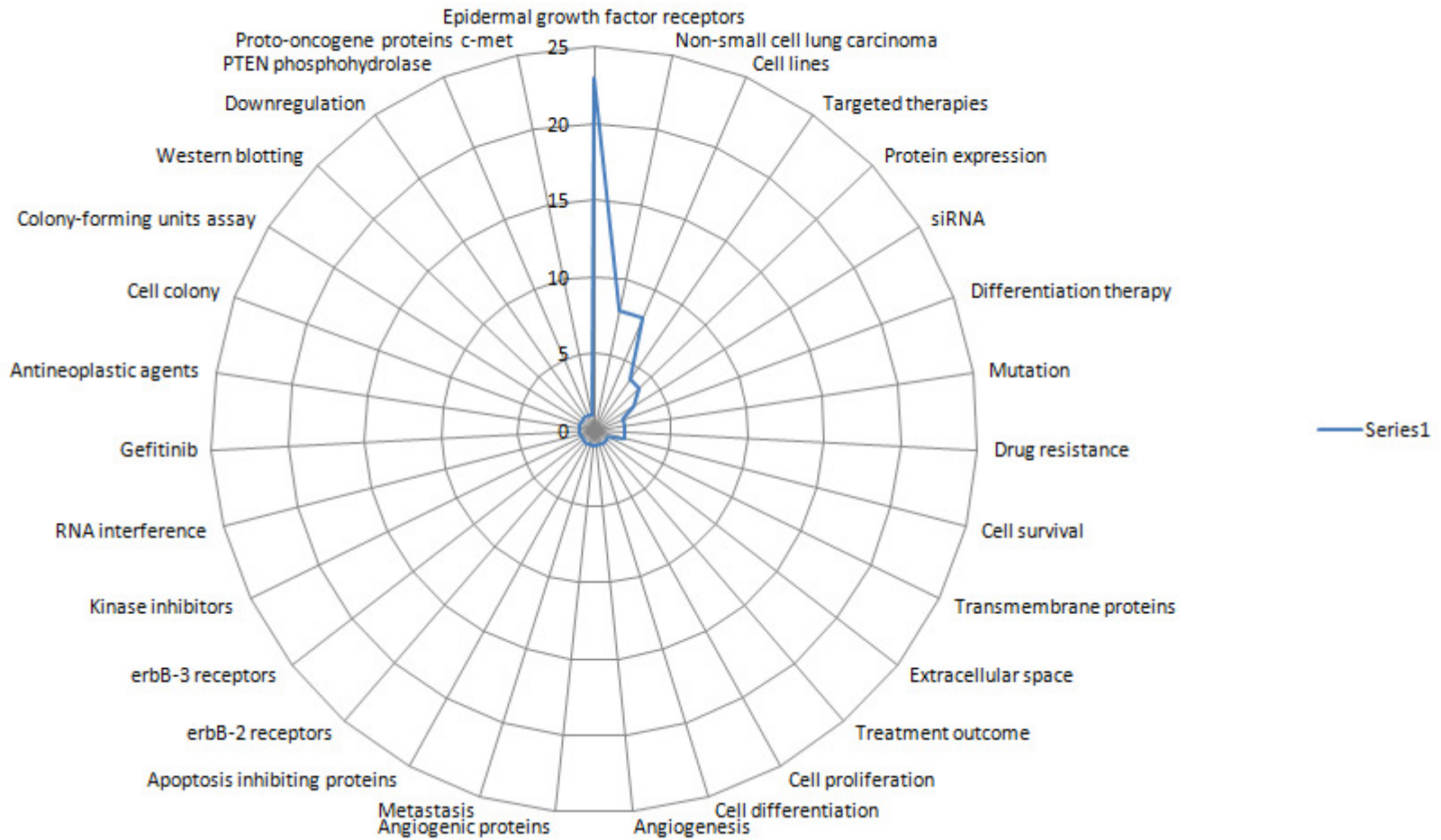http://www.youtube.com/watch?v=nokQBjk1s_8&feature=player_embedded

# How does it work?

- Develop controlled vocabulary
  » Prefer one with hierarchy

- Apply to full text
  » Or to the "heads"

- Decide on data points to convey information

- Divide the XML into graphable sections

# Start with data – like this XML file

# Index or tag using subject terms from thesaurus or taxonomy

- date,        category,        taxonomy term, frequency

| | | | | |
|---|---|---|---|---|
| 765 | SunMon2012_#18 | Molecular and Cellular Biology 1 | Histone deacetylase | 9 |
| 766 | SunMon2012_#2 | Molecular and Cellular Biology 1 | Mutation | 9 |
| 767 | SunMon2012_#24 | Molecular and Cellular Biology 1 | Neoplasm invasiveness | 9 |
| 768 | SunMon2012_#17 | Molecular and Cellular Biology 1 | Pyruvates | 9 |
| 769 | SunMon2012_#27 | Molecular and Cellular Biology 1 | Chemotherapy | 10 |
| 770 | SunMon2012_#4 | Molecular and Cellular Biology 1 | Epidermal growth factor re | 10 |
| 771 | SunMon2012_#8 | Molecular and Cellular Biology 1 | erbB-2 receptors | 10 |
| 772 | SunMon2012_#14 | Molecular and Cellular Biology 1 | Neoplasm invasiveness | 10 |
| 773 | SunMon2012_#1 | Molecular and Cellular Biology 1 | Phosphorylation | 10 |
| 774 | SunMon2012_#23 | Molecular and Cellular Biology 1 | Non-small cell lung carcino | 11 |
| 775 | SunMon2012_#11 | Molecular and Cellular Biology 1 | Phosphorylation | 12 |
| 776 | SunMon2012_#14 | Molecular and Cellular Biology 1 | Epithelial mesenchymal tr | 13 |
| 777 | SunMon2012_#3 | Molecular and Cellular Biology 1 | erbB-2 receptors | 14 |
| 778 | SunMon2012_#17 | Molecular and Cellular Biology 1 | Renal cell carcinoma | 14 |
| 779 | SunMon2012_#23 | Molecular and Cellular Biology 1 | Exons | 15 |
| 780 | SunMon2012_#26 | Molecular and Cellular Biology 1 | Mice | 17 |
| 781 | SunMon2012_#22 | Molecular and Cellular Biology 1 | Prolactin | 17 |
| 782 | SunMon2012_#9 | Molecular and Cellular Biology 1 | Epidermal growth factor re | 18 |
| 783 | SunMon2012_#14 | Molecular and Cellular Biology 1 | Intercellular signaling pep | 18 |
| 784 | SunMon2012_#23 | Molecular and Cellular Biology 1 | Mutation | 19 |
| 785 | SunMon2012_#10 | Molecular and Cellular Biology 1 | Epidermal growth factor re | 20 |
| 786 | SunMon2012_#18 | Molecular and Cellular Biology 1 | Mammary glands | 20 |
| 787 | SunMon2012_#7 | Molecular and Cellular Biology 1 | Epidermal growth factor re | 23 |
| 788 | SunMon2012_#26 | Molecular and Cellular Biology 1 | Colon cancer | 24 |
| 789 | SunMon2012_#23 | Molecular and Cellular Biology 1 | Epidermal growth factor re | 25 |

# SunMon2012_#7



Epidermal growth factor receptors
Proto-oncogene proteins c-met
PTEN phosphohydrolase
Non-small cell lung carcinoma
Cell lines
Downregulation
Targeted therapies
Western blotting
Protein expression
Colony-forming units assay
siRNA
Cell colony
Differentiation therapy
Antineoplastic agents
Mutation
Gefitinib
Drug resistance
RNA interference
Cell survival
Kinase inhibitors
Transmembrane proteins
erbB-3 receptors
Extracellular space
erbB-2 receptors
Treatment outcome
Apoptosis inhibiting proteins
Cell proliferation
Metastasis
Cell differentiation
Angiogenic proteins
Angiogenesis

Series1

Adenocarcinoma
Bladder cancer
Lungs
Non-small cell lung carcinoma
Lung cancer
Mutation

0 1 2 3 4 5 6 7 8 9 10

# Load to a visualization program
# Like Prefuse

# Or Pajek



# Networks / Pajek

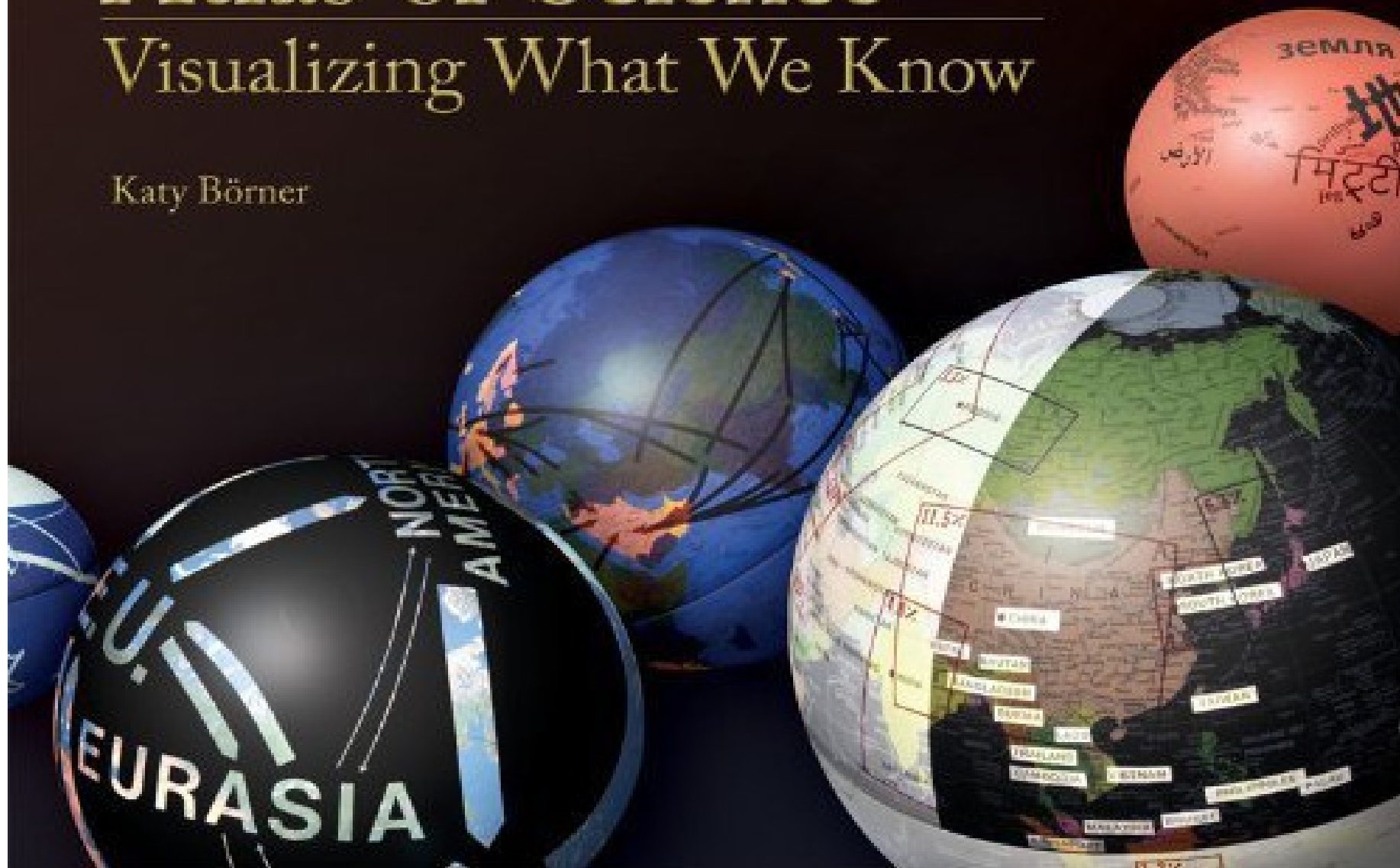**Package for Large Network Analysis**

# How to

**In construction**

1. First steps to visualization of networks with **Pajek** (May 2002);
   Data to play with (ZIP)
2. Tutorial and course material at Faculty of Social Sciences, Ljubljana
3. Printing pictures created with **Pajek**
4. Converting Excel datasets into **Pajek** format
5. Converting text file datasets into **Pajek** format
6. Run **Pajek** from command line or other programs
7. Analysis of networks and vectors sent from **Pajek** to program **R**
8. Analysis and visualization of genealogies with **Pajek**; see also Toolkit and Notes
9. Kinship macros
10. Generating random genealogies
11. Some hints on working with extremely large networks in **Pajek**
12. Runing **Pajek** from read only location
13. Runing **Pajek** on Linux
14. Combining SVG with Background Picture

**Pajek**; Vlado/Networks

# Atlas of Science
## Visualizing What We Know

Katy Börner

# National Information Center for Educational Media

- Albuquerque's own
  - » Sandia developed VxInsight
  - » Access Innovations = NICEM
- Same data – several views
- Primary and Secondary Education in US
- Shows the US Valley of Science
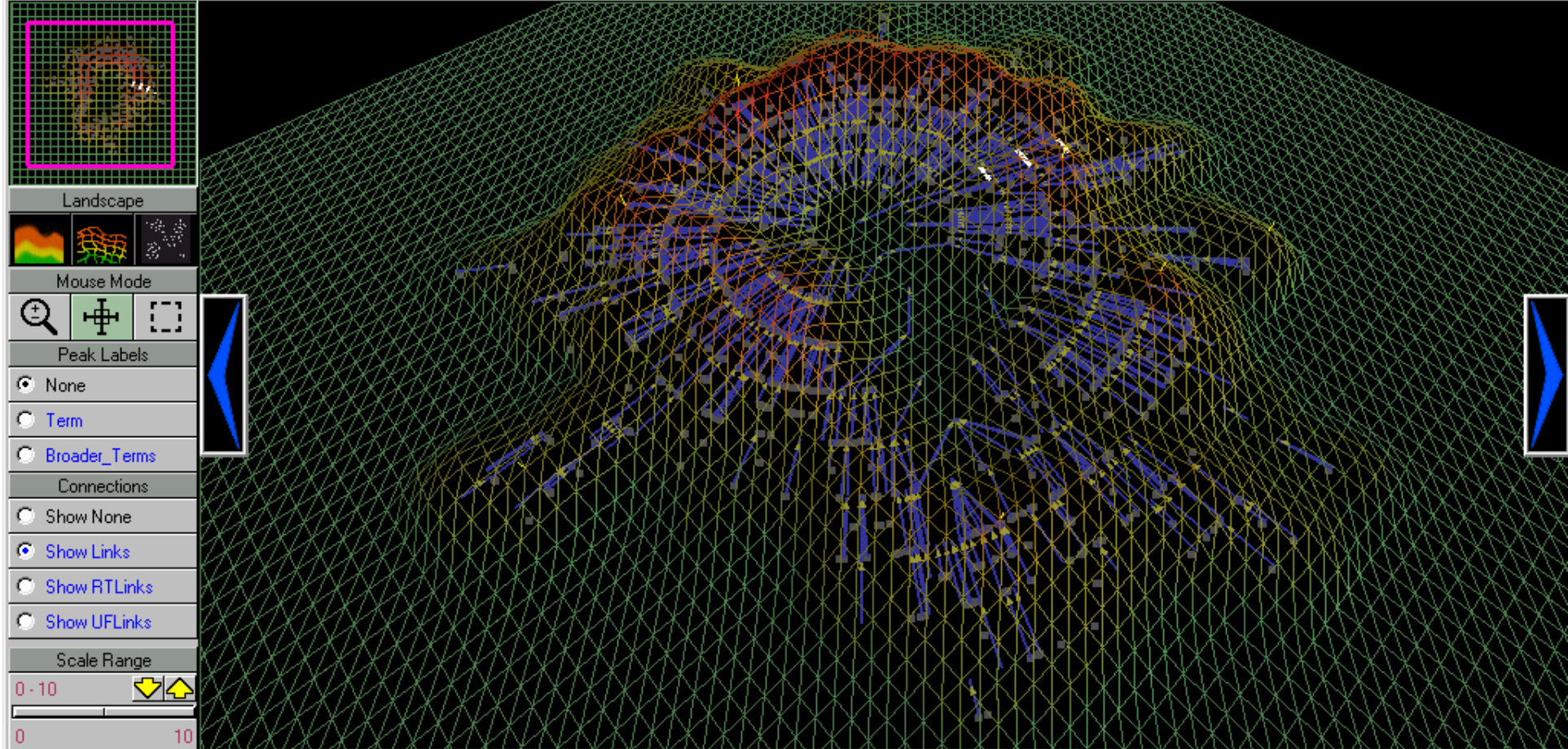- Little Science taught in elementary years

VxInsight: 2.163

File  Options  Tools  Plugins  Help

Term
Broader_Terms
Narrower_Terms
Related_Terms
Rank    UsedFor_Terms    History    Status

Landscape

Mouse Mode

Peak Labels
○ None
○ Term
○ Broader_Terms

Connections
○ Show None
● Show Links
○ Show RTLinks
○ Show UFLinks

Scale Range
0 - 10
0        10

Database Query
Show objects where  broader_terms  IS_LIKE  music  Send

Legend
● broader_terms IS_LIKE %music% --- 87

Advanced Query    Send to Excel    Clear Legend    Clear Selected

Slide 7 of 9    Default Design

Start    Eudora - [J...    contentwe...    Data Harm...    VxInsigh...    2:41 PM

# Using visualization to show

- From a society / publisher perspective
  - » Identify Core, Boundary and Cross Border
  - » Provides Indicators
    - Activity
    - Growth
    - Relatedness
    - Centrality
  - » Locates Journal domains

- From a thesaurus perspective
  - » Identifies terms that are too broadly defined
  - » Potential Improvements in thesaurus structure using topic structures
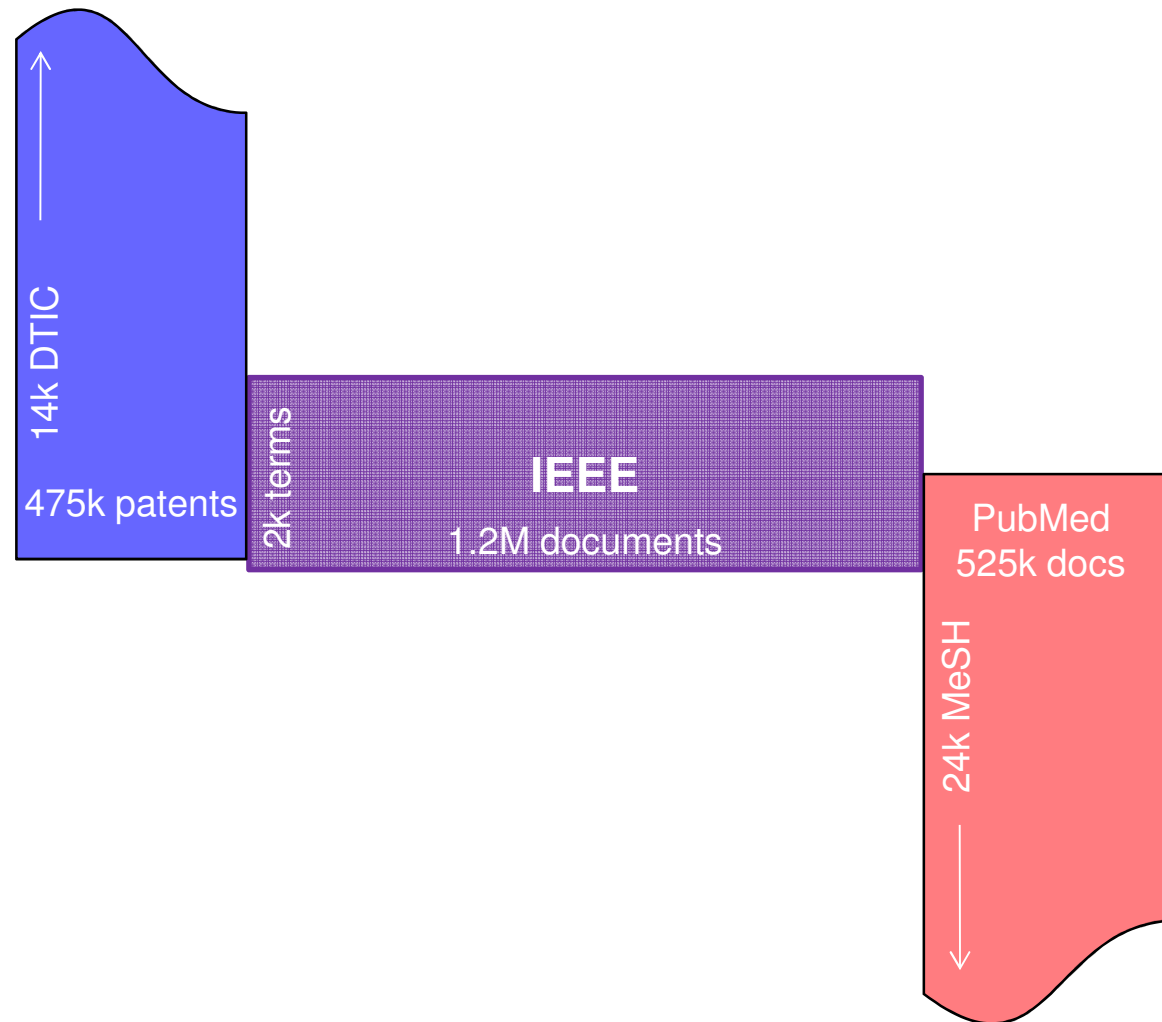
# Case Study:
# Mapping IEEE thesaurus space

- We are interested in an expanded map that includes adjacencies to the IEEE data
    - » Expanded term set shows adjacent white space; opportunities for expansion

- Overlaps and edges of the science
    - » We need comparison data

- Learn the directions in the field
    - » Low occurrence rate in IEEE documents?
    - » Linkage to terms in IEEE documents?

- Where do we find these terms? How can we add them?

ACCESS Innovations.

# The process

- Built a rule base to auto index IEEE content
  - » "90 % accuracy out of the box on journal data"*
  - » "80% out of the box on proceedings data"*
- **The overlapping data sets**
  - » Auto indexed 1.2 million Xplore records
  - » Auto indexed 10 years of US Patent data
  - » Auto indexed 10 years of Medline
- Term sets used
  - » IEEE thesaurus terms rule base
  - » Medical Subject Headings (MeSH) (and simple rule base)
  - » Defense Technical Information Center (DTIC) Thesaurus ( and simple rule base)
  - » Similar level of detail to current IEEE thesaurus terms
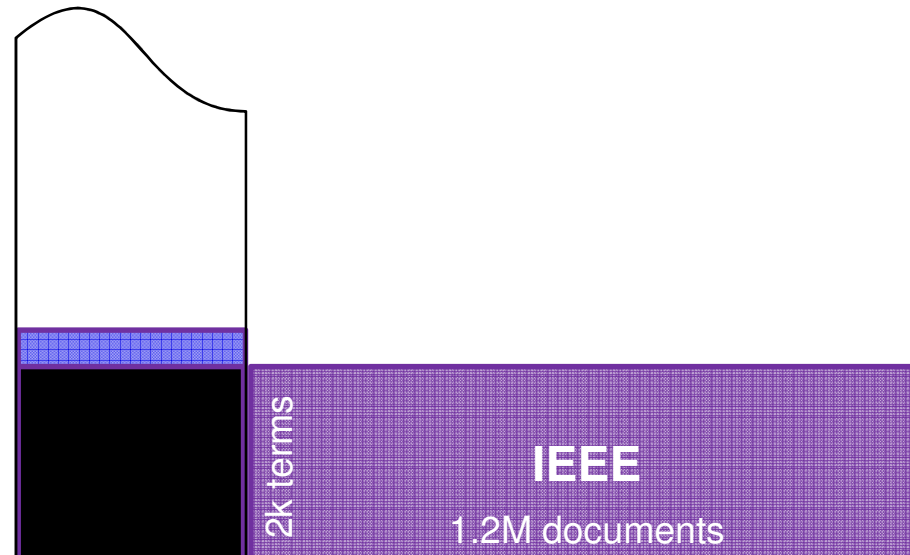
# Defining expanded term space

## 1. The data - Select related corpus

14k DTIC

475k patents

2k terms

**IEEE**

1.2M documents

PubMed
525k docs

24k MeSH

Access Innovations

# Defining expanded term space

## 2. Identify related terms
## Use the IEEE Thesaurus to index the three collections

2k terms

**IEEE**

1.2M documents

# Defining expanded term space

## 2. Identify related terms
## Use MESH and DTIC to also index the three collections

IEEE

1.2M documents

2k terms

# Defining expanded term space

## 3. Resulting term set
## The co-indexed items from the three collections


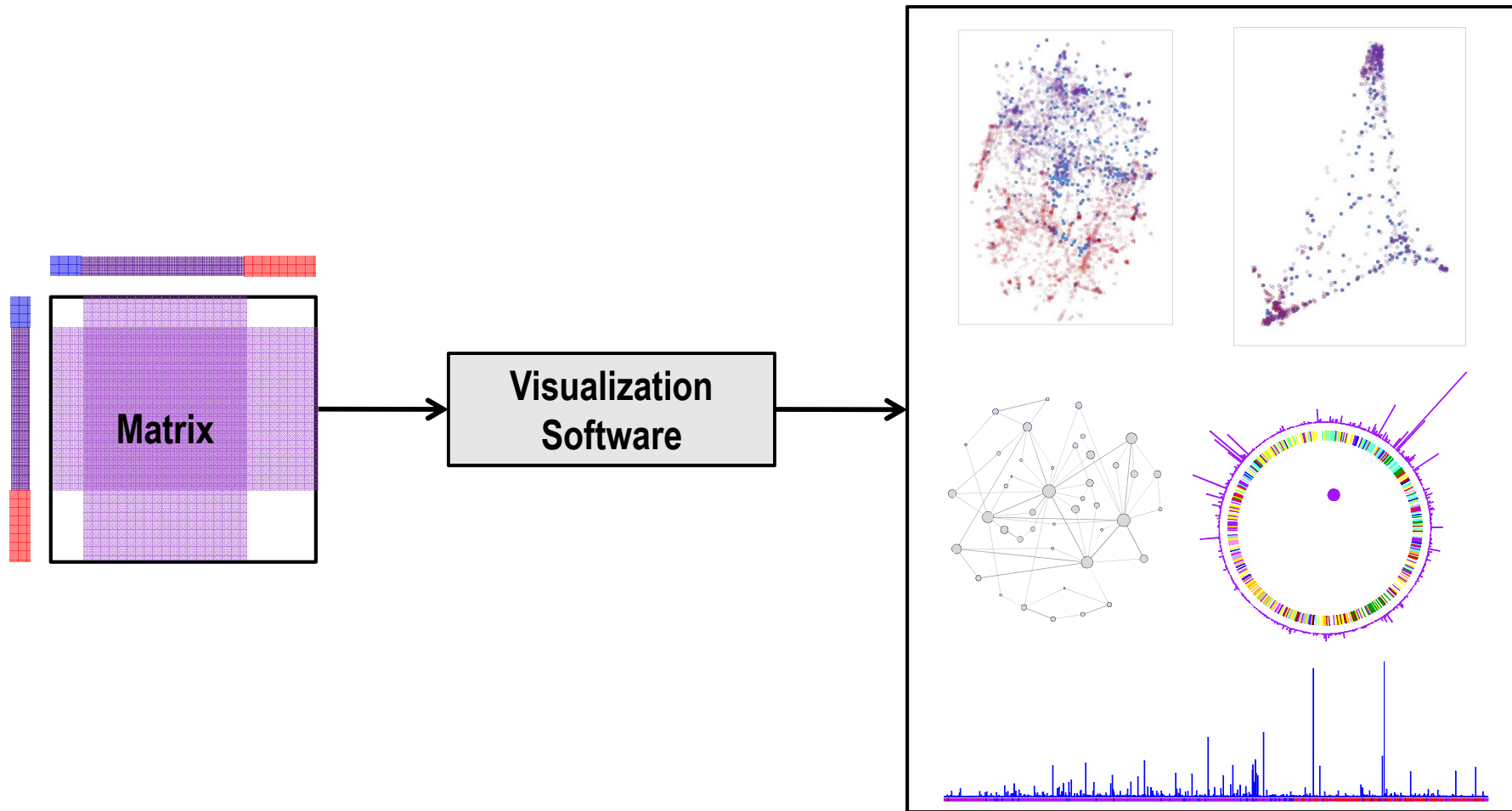
2k terms

**IEEE**
1.2M documents

# Defining expanded term space
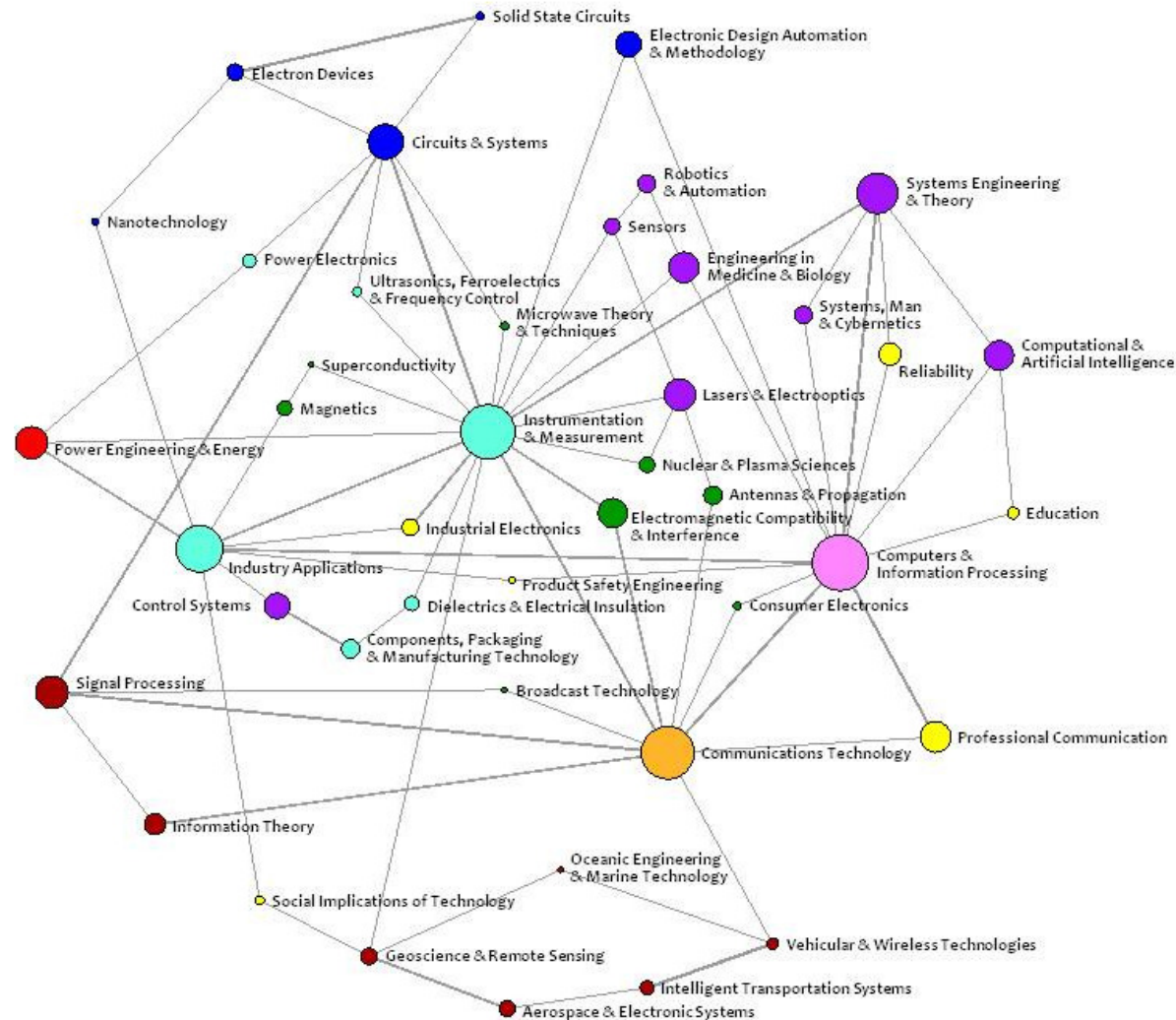
## 4. Term:Term Matrix
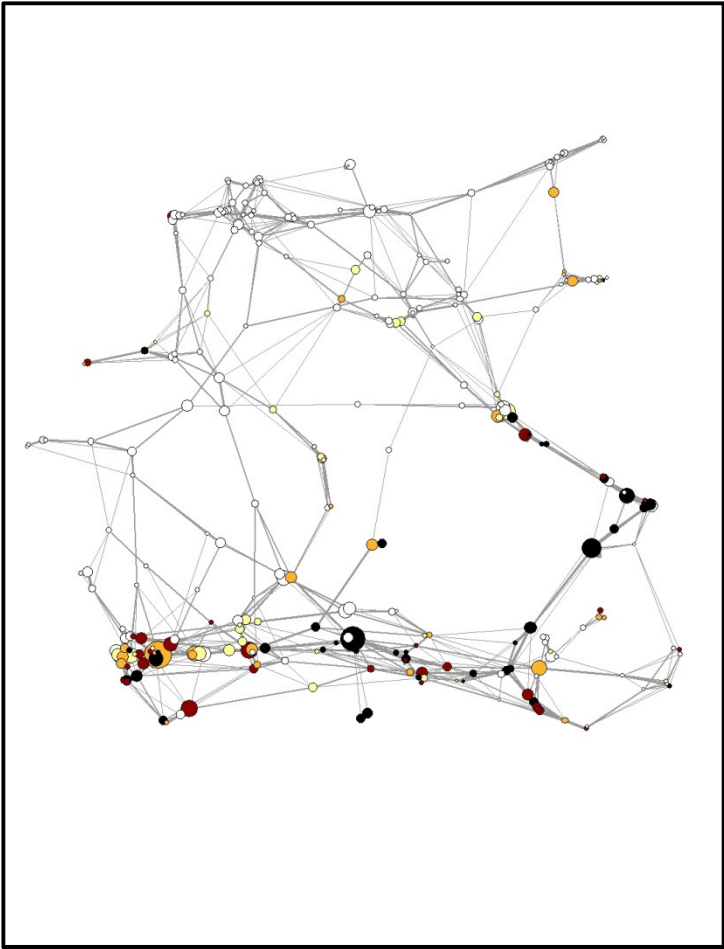## Where do the articles and their indexing intersect?

# Visualization Strategies
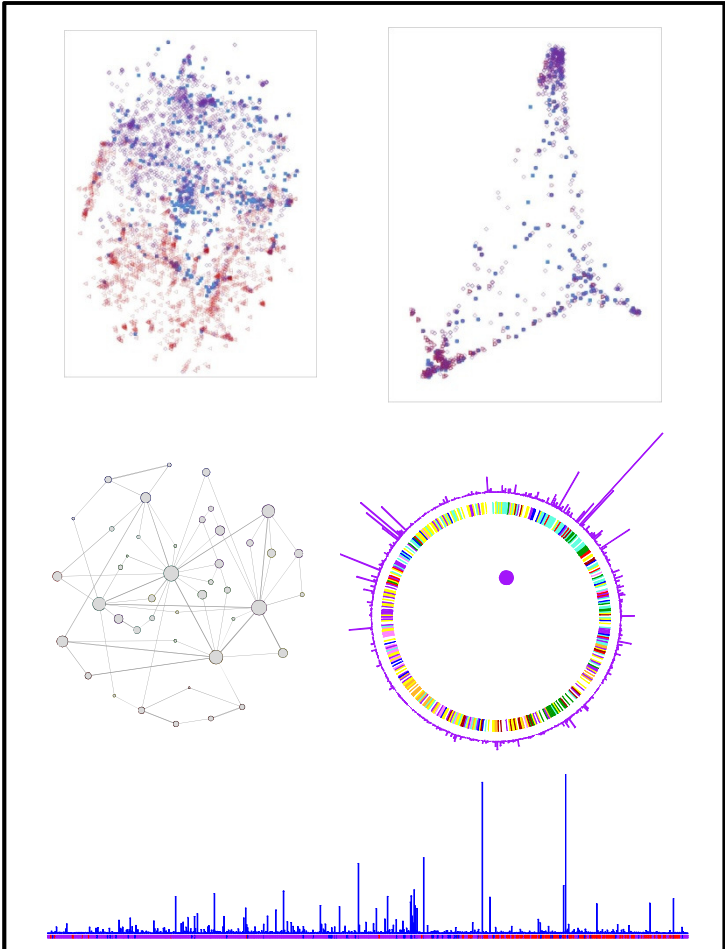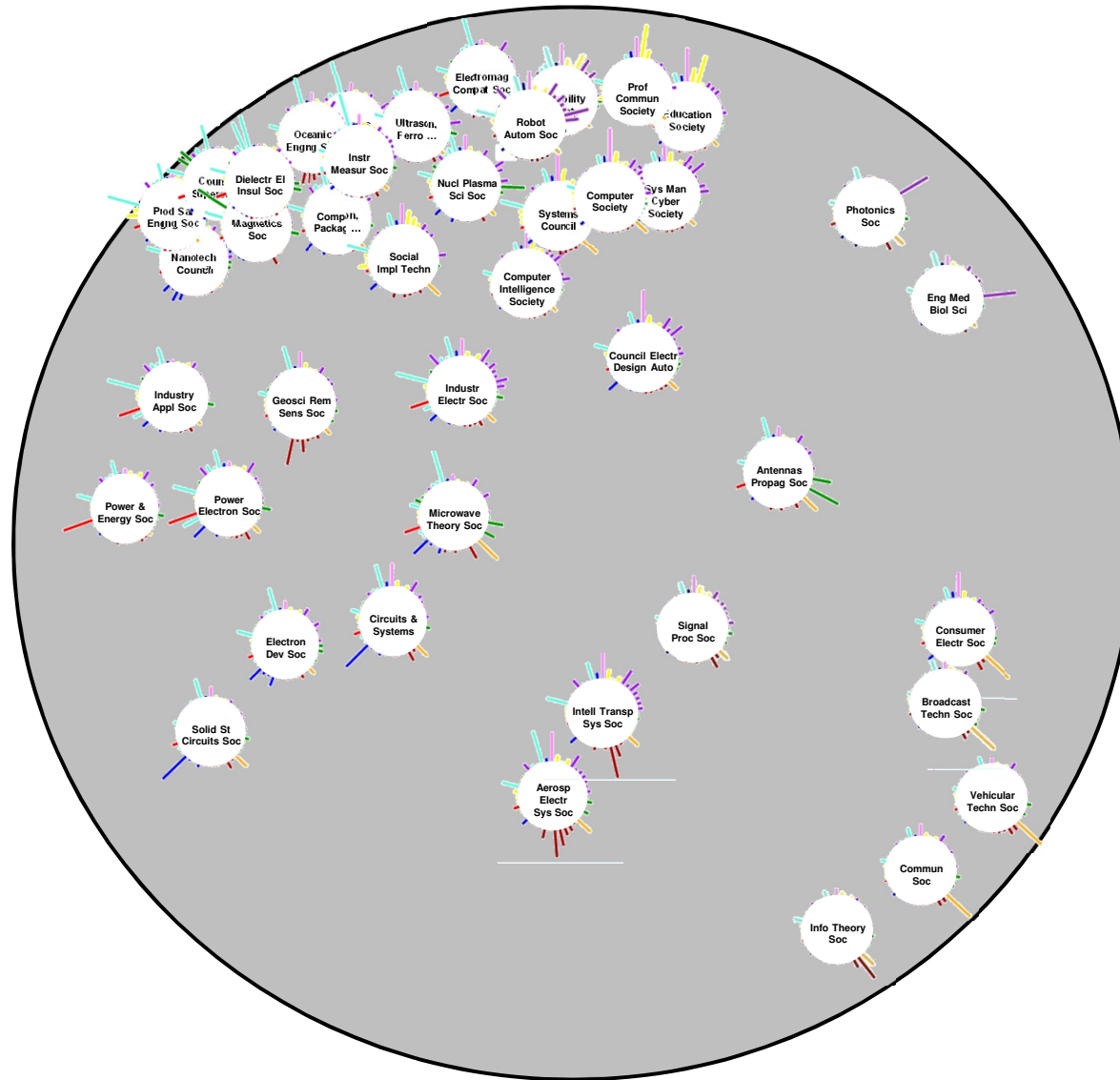
# All data up-posted to the top level

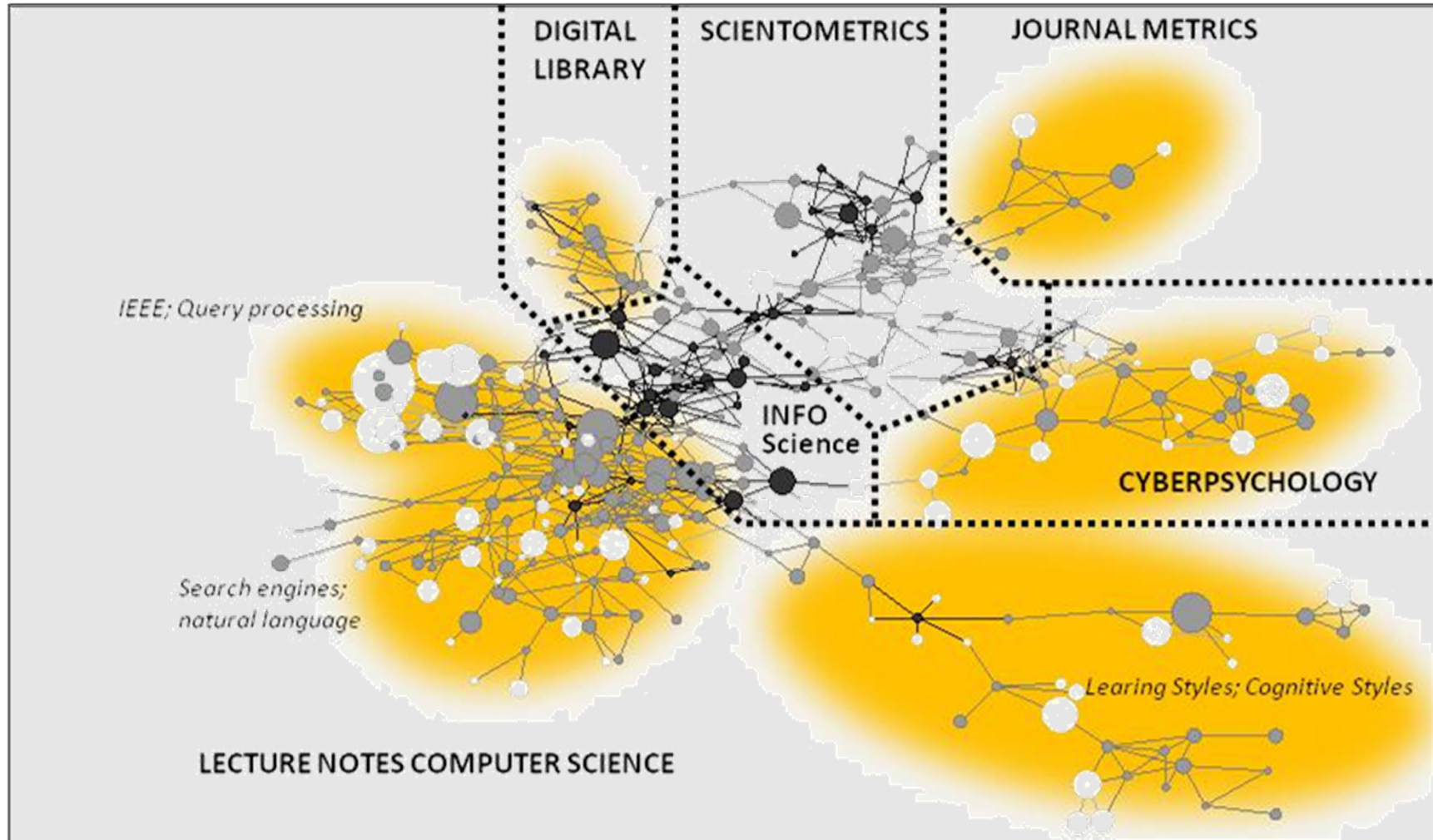# Many map options

**Previous Experience**

**IEEE Experience**

# IEEE Portfolio
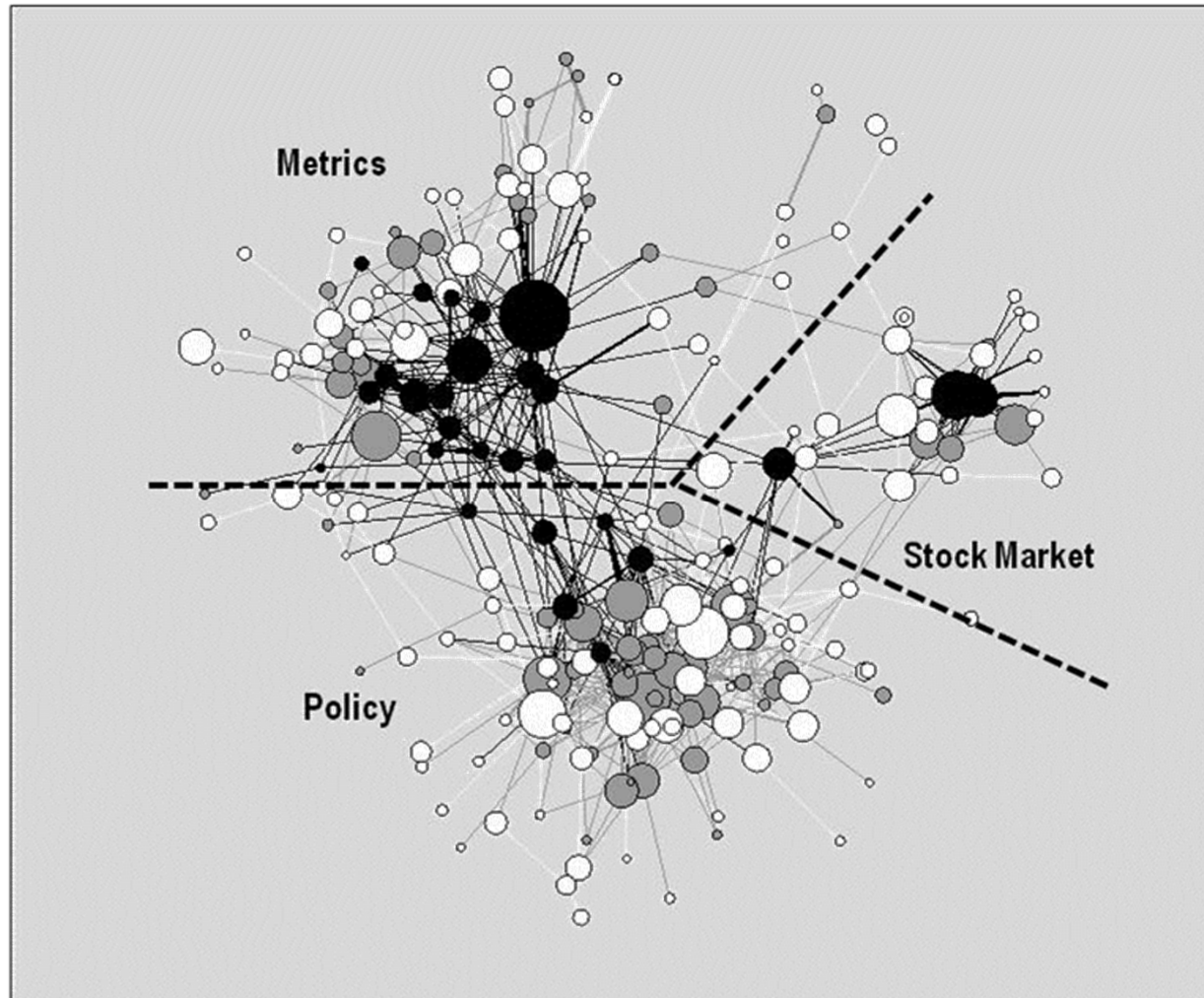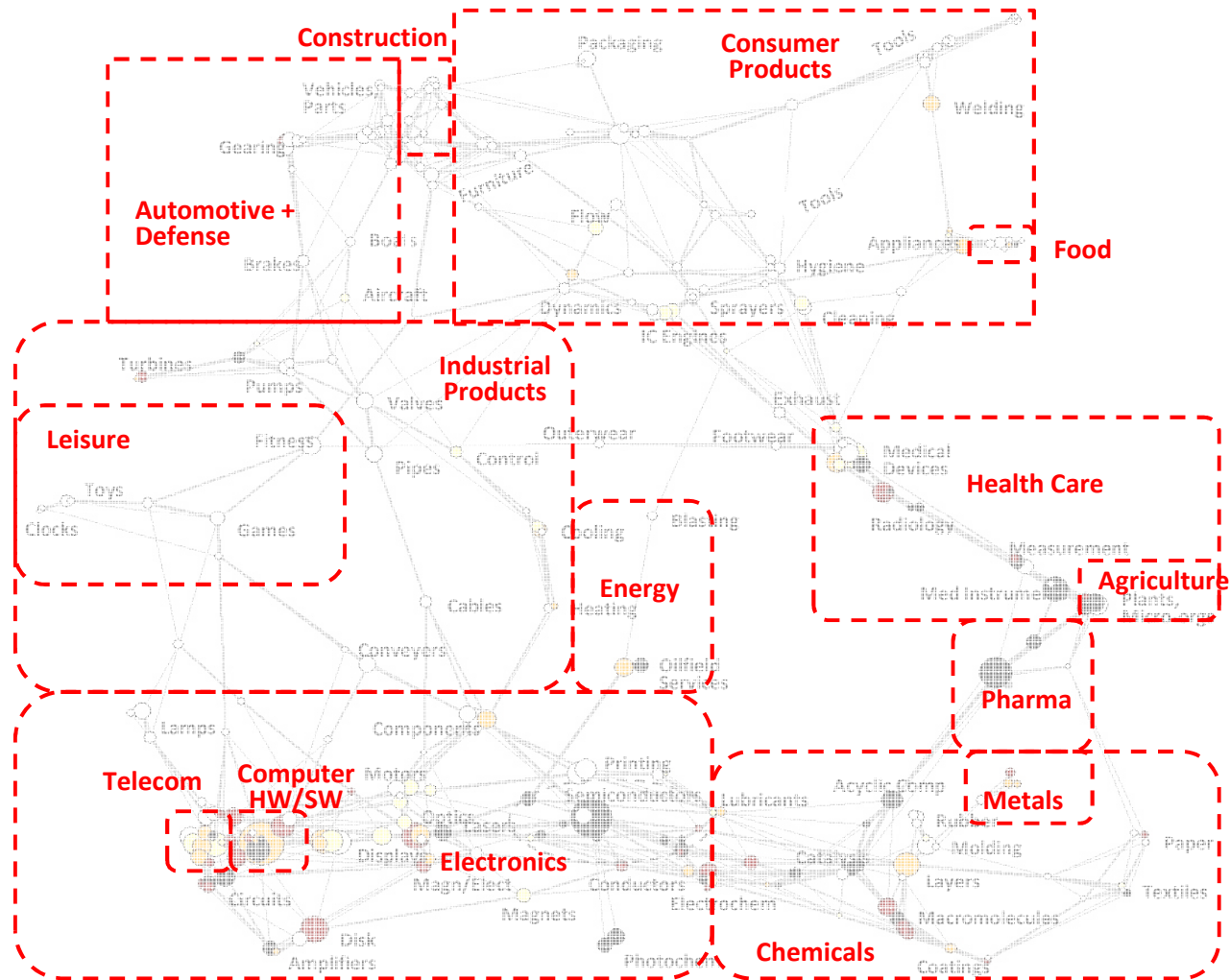
# Radial Visualization

# Publication Strategy



JASIST reference

# Conference Strategy

# Use a Thesaurus to Label Maps

# Questions Answered

- Is there a way, using our own information, to forecast our direction?

- Where is the industry headed?  What about by technology sector?

- Does our coverage match our mission and vision?

- Can we become smarter about our data and potential markets using our collection in new ways?
  Are the societies publishing and talking about what their charter indicates they cover?

- What are the trends – are topics emerging/cooling?

- Can we use technology and our own data to explore these questions while enhancing our data?

# The research team

- Access Innovations / Data Harmony
  - » Founded in 1978
  - » Data enrichment and normalization
  - » Suite of Semantic Enrichment tools

- SciTechStrategies
  - » Understanding data through visualization

- IEEE Indexing & Abstracting Group

# We looked at visualization of data

- Finding the Metrics
  - » Measurement
  - » Numbers
  - » Terms as indicators
- Ways to show
  - » Adjacency
  - » Relationships
  - » Trends
  - » Co – occurrence
  - » Conceptual distance

- How to enrich with
  - » Linking
  - » Semantic enrichment
  - » Classification
- Maps supporting
  - » Forecasting
  - » Trend analysis
  - » Segmentation
  - » Distribution

# Effective maps require

- Contextual data
- Detailed data
- Classification methods
- At least two directions in the matrix
- A little art for fun

# It just takes a little imagination

## Thank you

Marjorie M.K. Hlava
President
mhlava@accessinn.com

Jay Ven Eman, CEO
J_ven_eman@accessinn.com

, Access Innovations
505-998-0800 43