

# Algebraic and Tensor Methods for Anomaly Detection

***Brett Bader***

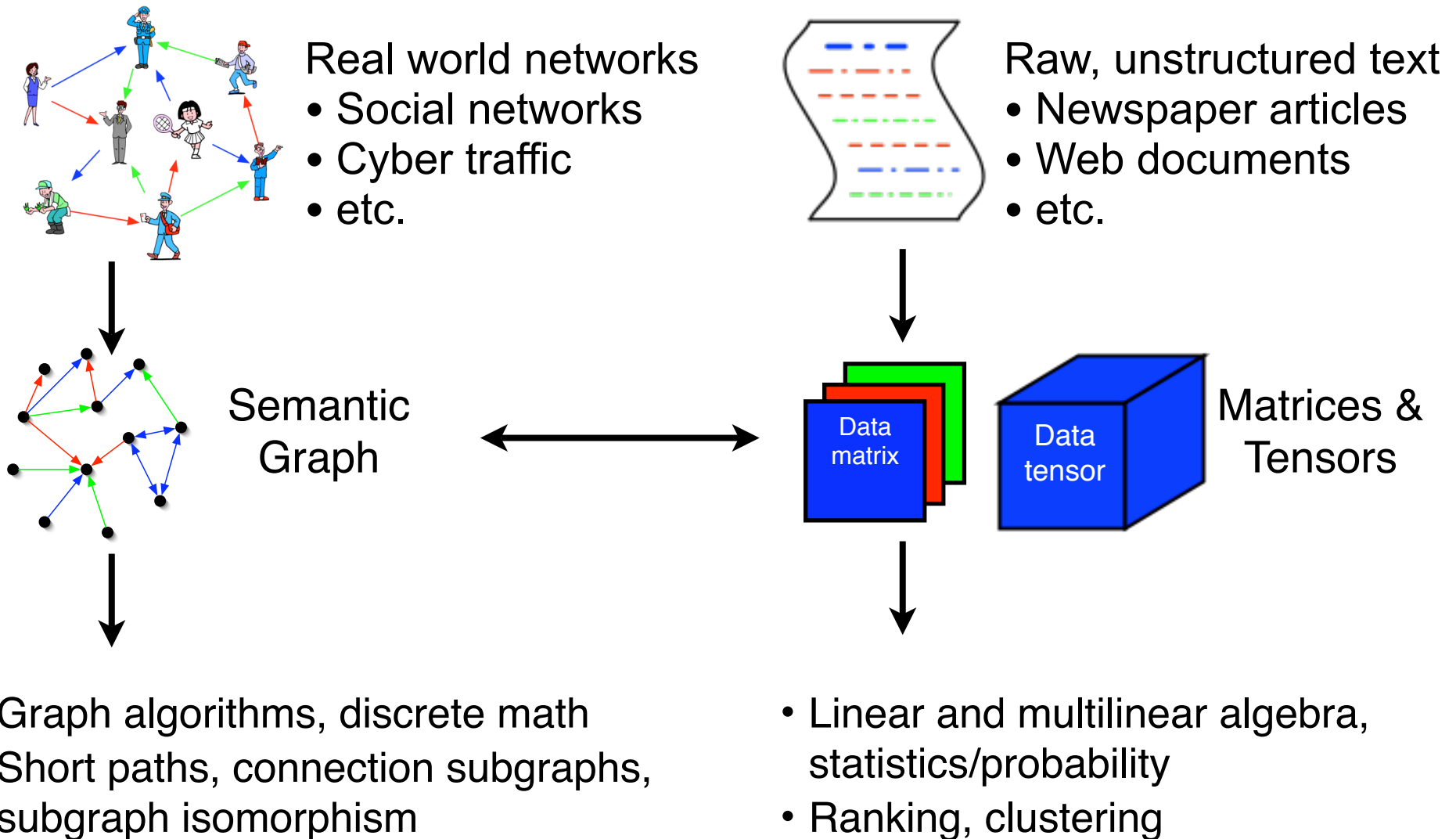
Senior Member of Technical Staff  
Computer Science and Informatics

June 17, 2010

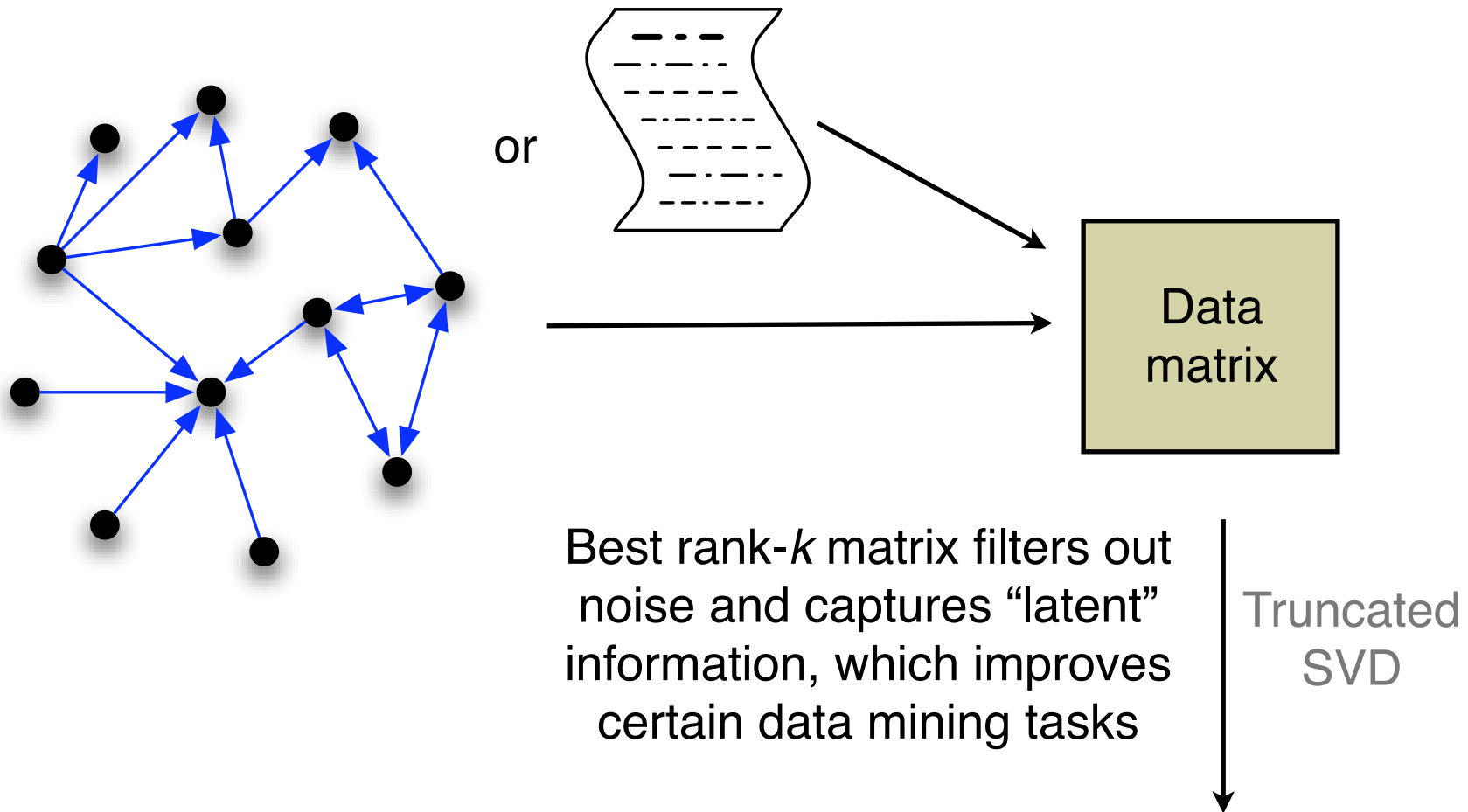
# Robust data analysis requires appropriate data abstractions and algorithms

Sandia uses *semantic graphs* and *tensors* as unifying data abstractions

- Supports rich relationship-centered analysis
- Combines large, heterogeneous data corpora
- Different abstractions support different analytics



# Traditional Analysis



Examples:

- Latent Semantic Analysis
- Text Analysis (LSI)
- Web search (HITS)
- Clustering

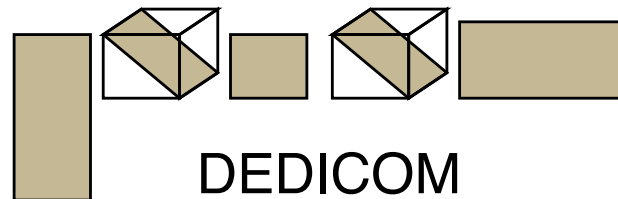
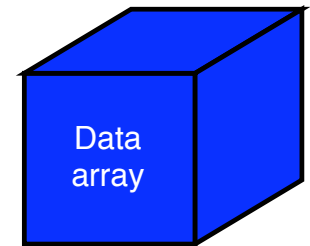
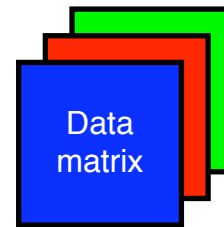
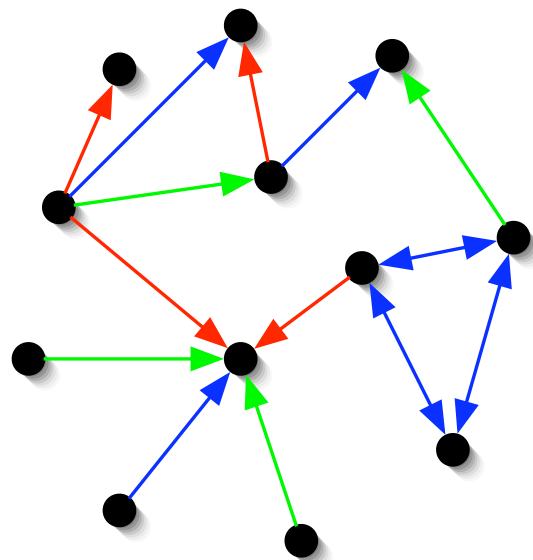
The diagram shows the SVD decomposition of a data matrix  $A$  into three components:  $U_k$ ,  $\Sigma_k$ , and  $V_k^T$ . Each component is shown in a box.  $\Sigma_k$  is a square box with a diagonal line from the top-left to the bottom-right, indicating it is a diagonal matrix. Below these boxes, the equation  $A_k = U_k \Sigma_k V_k^T = \sum_{i=1}^k \sigma_i u_i v_i^T$  is written, with the index  $k$  in red.

$$A_k = U_k \Sigma_k V_k^T = \sum_{i=1}^k \sigma_i u_i v_i^T$$

But there may be more useful information in the data!

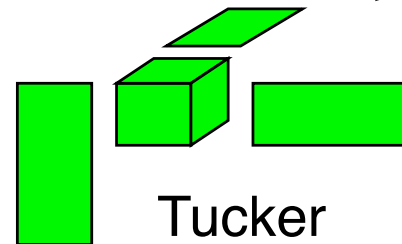
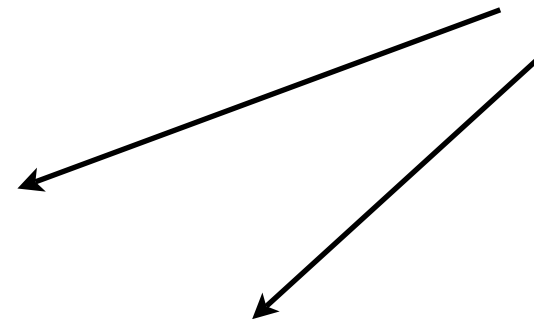
# New Paradigm: “Multidimensional Data Mining”

Build a “data array” such that there is a data matrix for each link type.



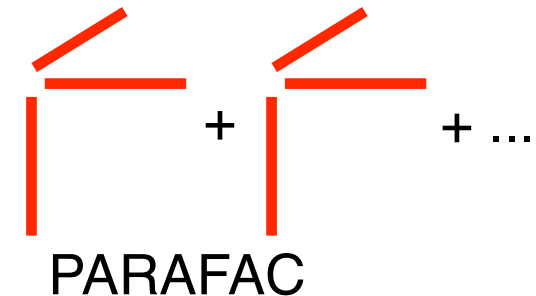
DEDICOM

Third dimension offers more explanatory power: uncovers new latent information and reveals subtle relationships



Tucker

Multilinear algebra



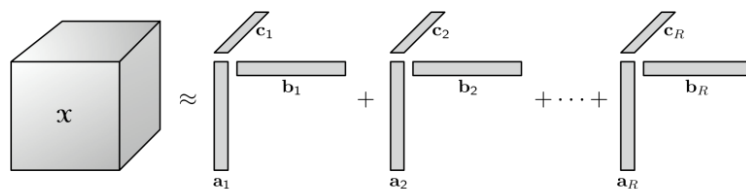
PARAFAC



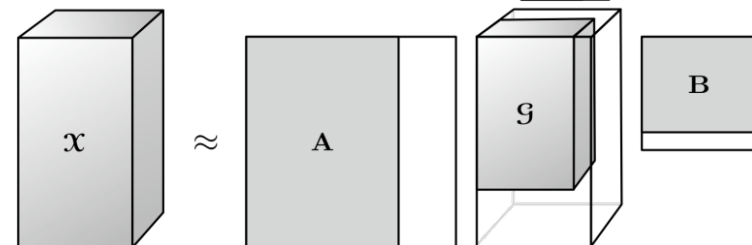
Unique data mining capability  
developed at Sandia

# Many Types of Tensor Decompositions

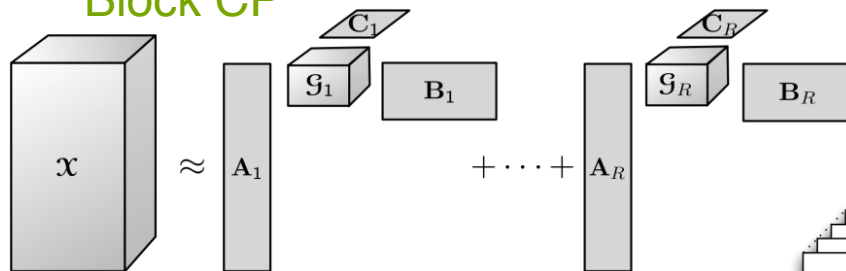
## CANDECOMP/PARAFAC (CP)



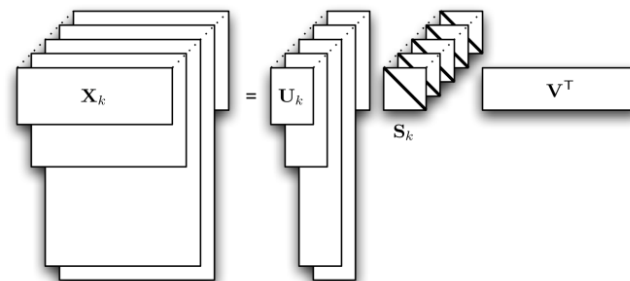
## Tucker



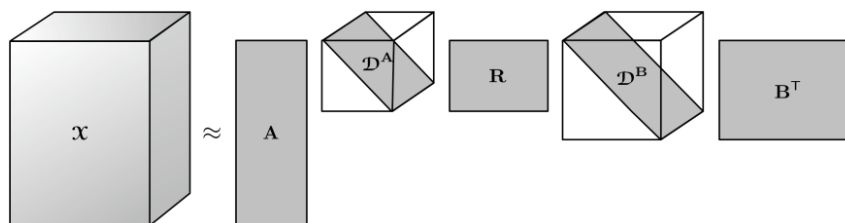
## Block CP



## PARAFAC2



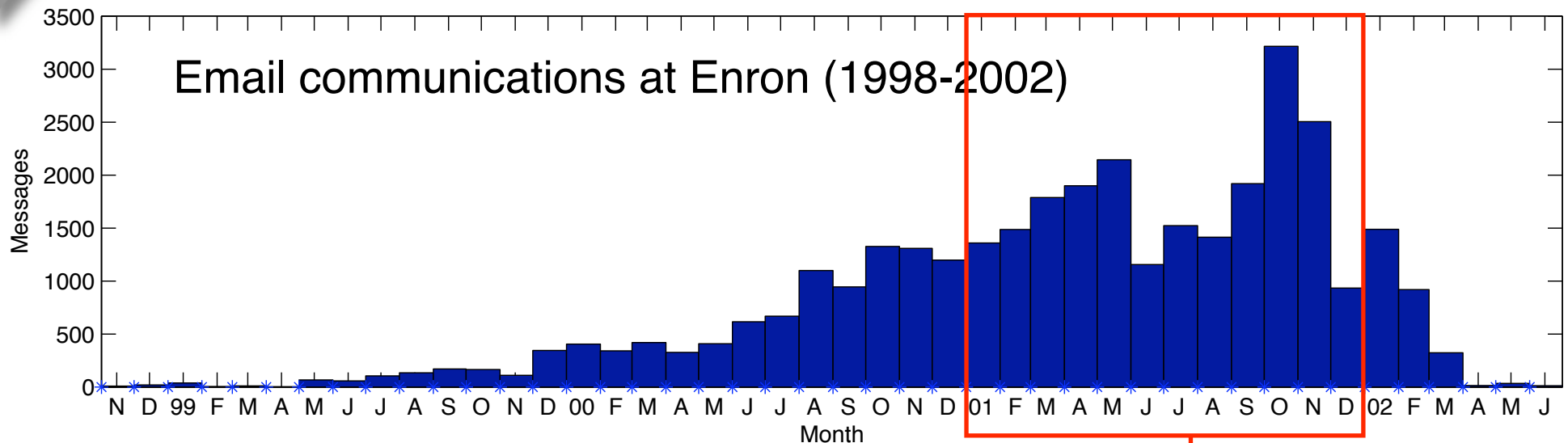
## PARATUCK2 (DEDICOM3)



Kolda & Bader, Tensor  
Decompositions and Applications,  
SIAM Review, 2009

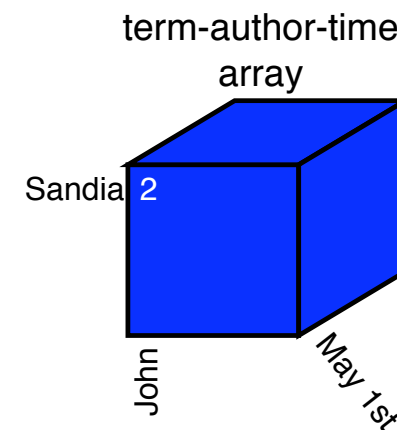
Each decomposition provides a  
different interpretation of the data

# Case Study: Discussion Tracking in Email

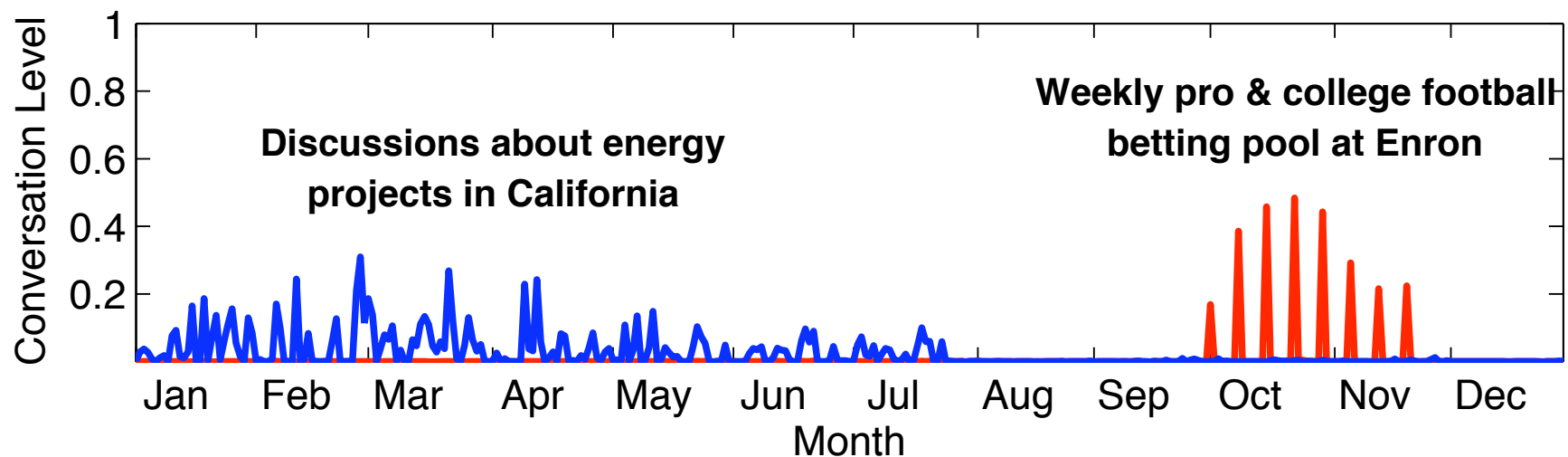
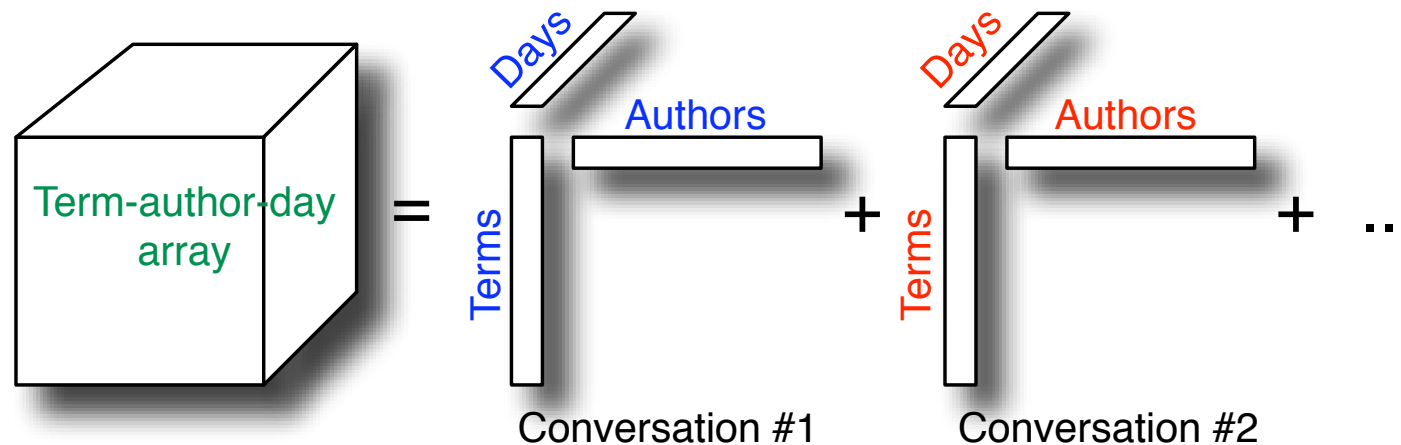


53,733 messages  
from 184 employees

- Situational awareness
- What can we learn from these email conversations?
  - **What** are the major topics of conversations?
  - **Who** are the major participants?
  - **When** are they taking place?



# Tensor analysis finds unusual activity by associating terms with people over time



Key terms: California, power, utilities, energy, utility, governor, market

games, week, missed, picked, prize, wins, scored, upsets

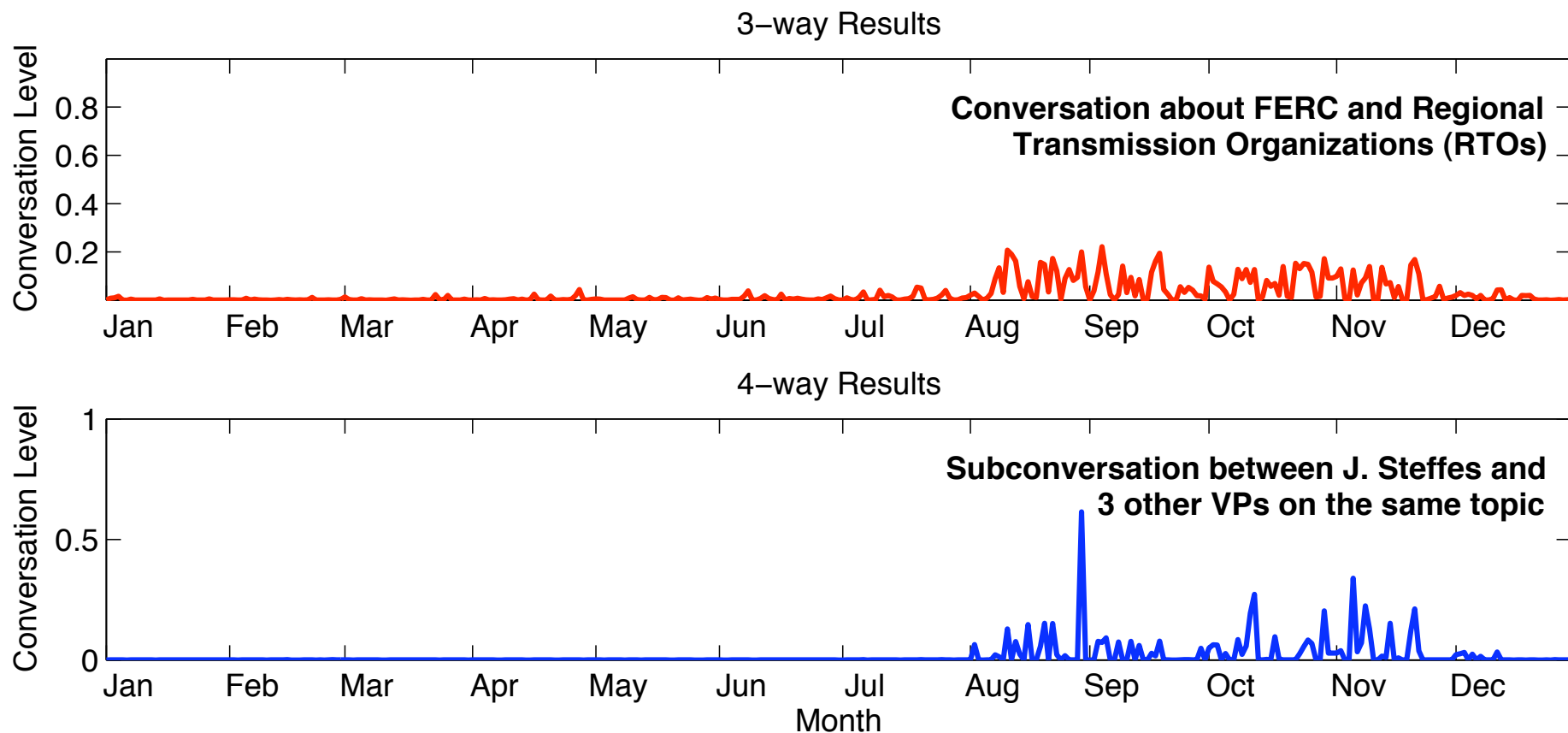
Key authors: J. Steffes, S. Kean, J. Dasovich, R. Shapiro, P. Allen, ...

A. Pace, L. Campbell, C. Dean

# Four-way analysis shows deeper relationships

4-way array: Author x Recipient x Terms x Time

- 4-way analysis may track subconversation already found by 3-way analysis
- Provides context and temporal patterns of social network

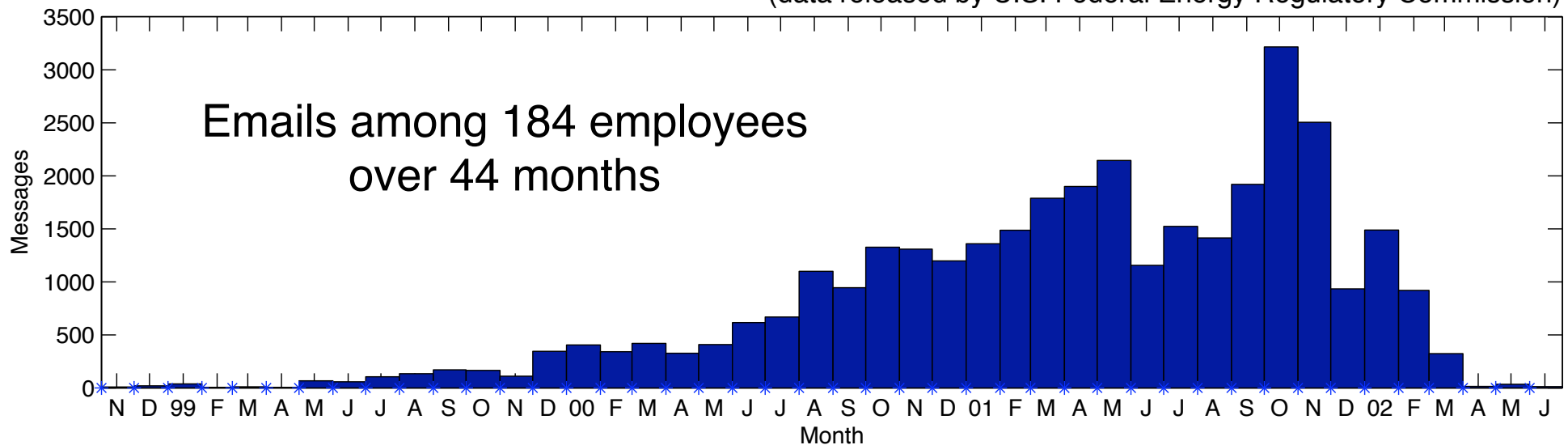




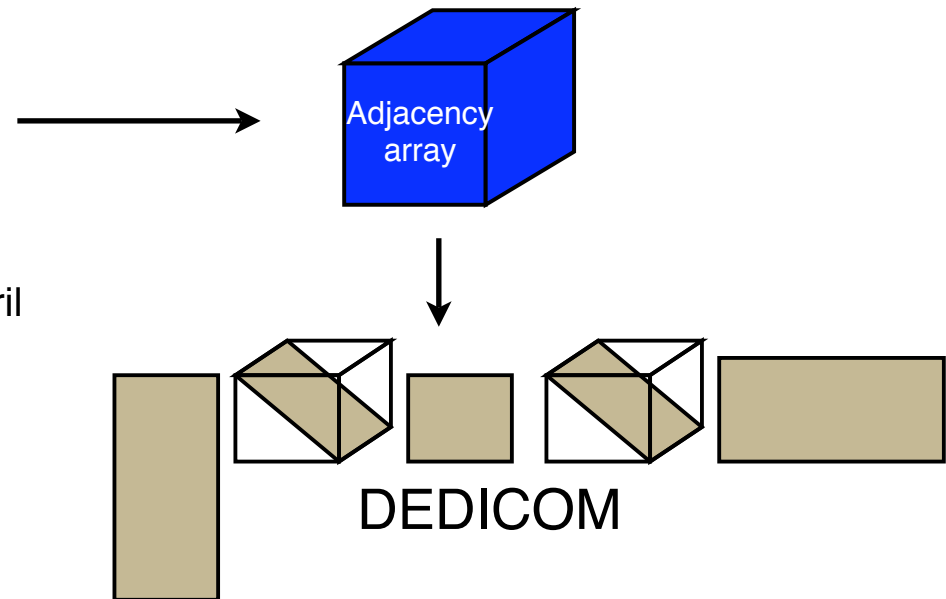
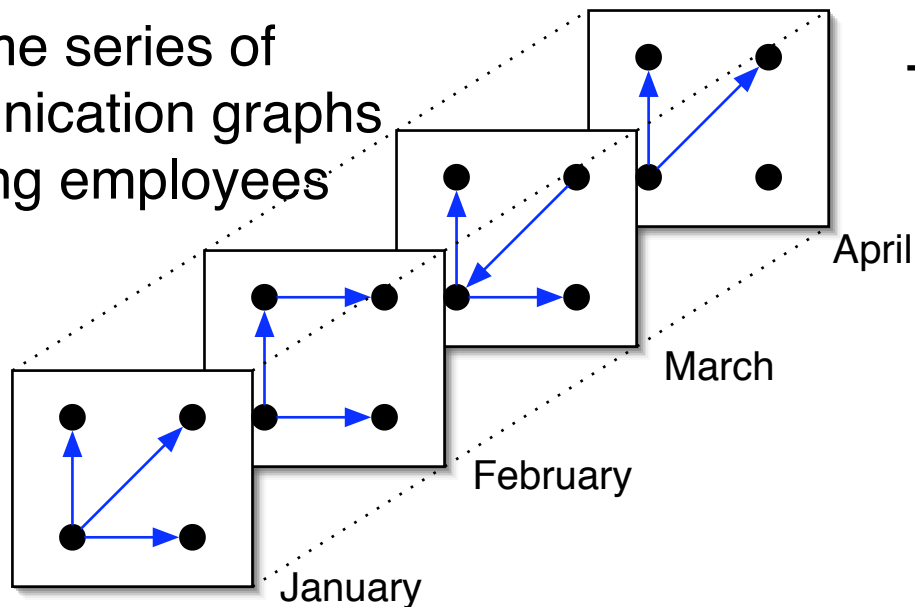
# Case Study: Pattern Analysis in Email Networks

Email communications at Enron (1998-2002)

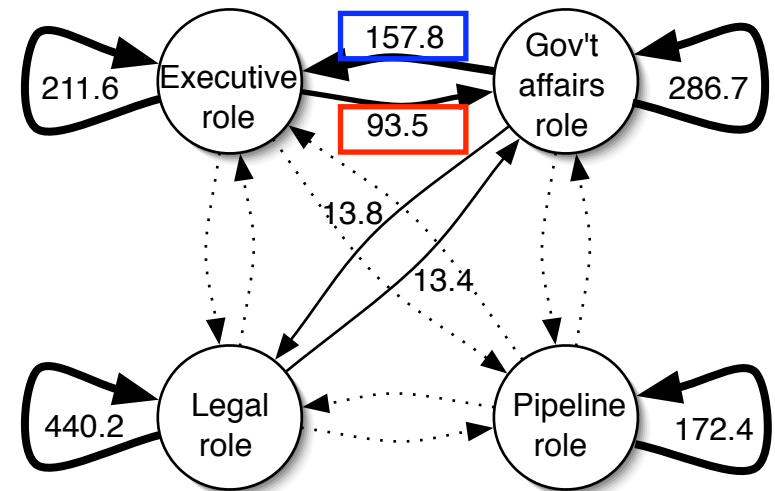
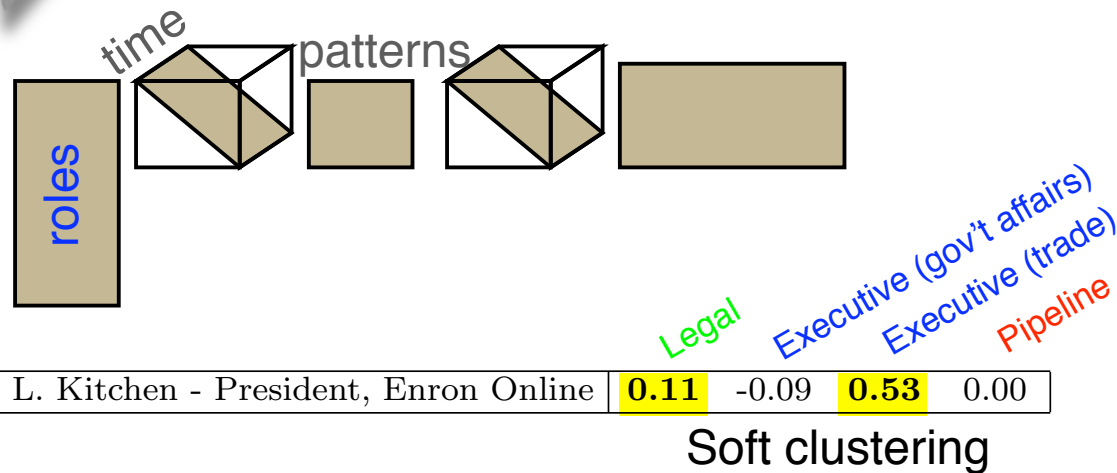
(data released by U.S. Federal Energy Regulatory Commission)



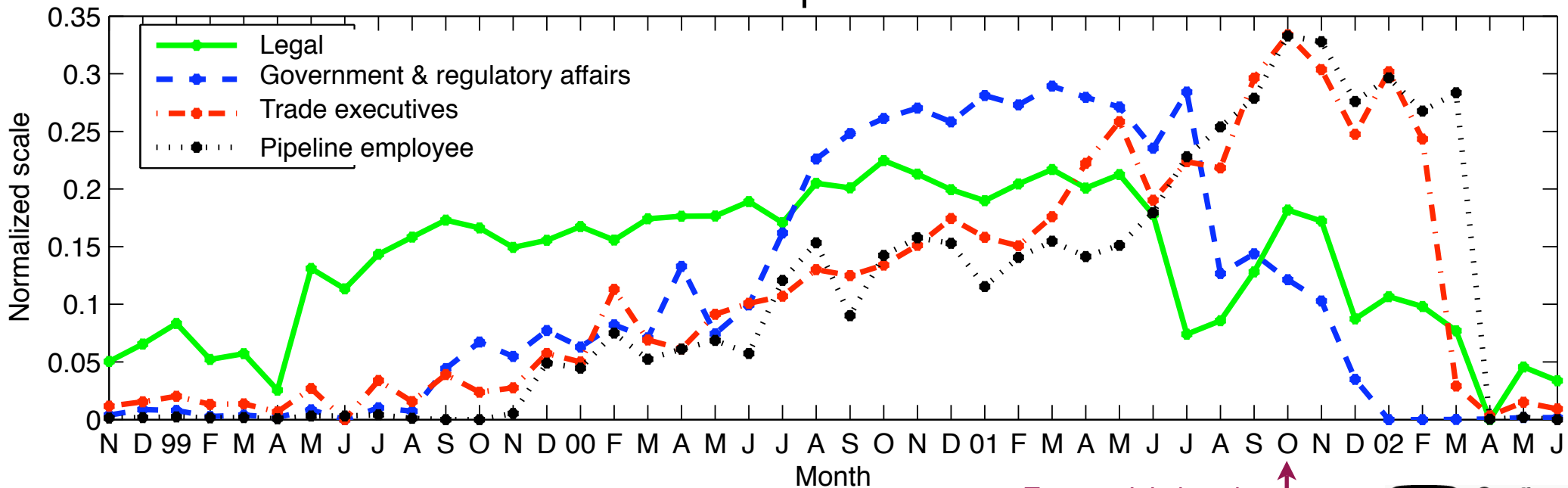
Time series of  
communication graphs  
among employees



# Analysis shows employee roles and communication patterns



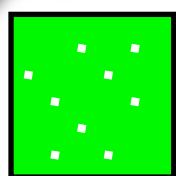
## Communication patterns over time



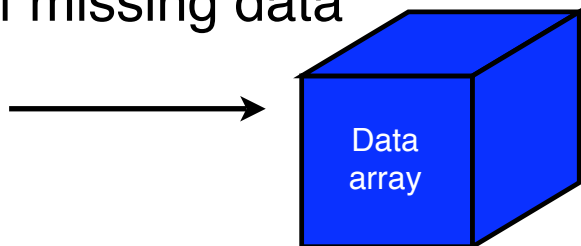
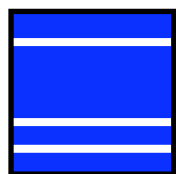
Bader, Harshman, Kolda, Temporal analysis of semantic graphs using ASALSAN, in ICDM 2008.

Enron crisis breaks; investigation begins

# Our algorithms can handle missing data

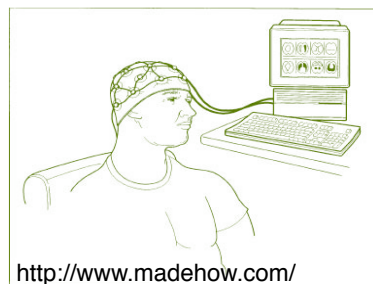


Random or  
systematic patterns  
of missing data

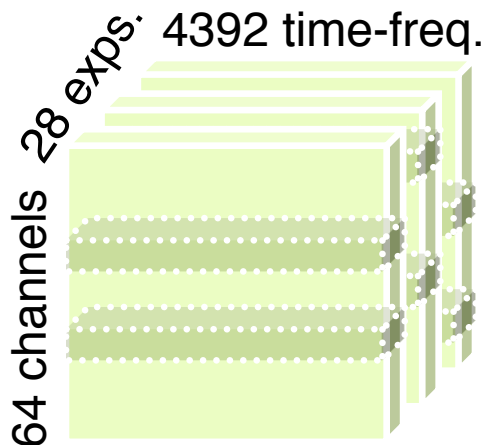


Fit model using  
derivative-based  
algorithms

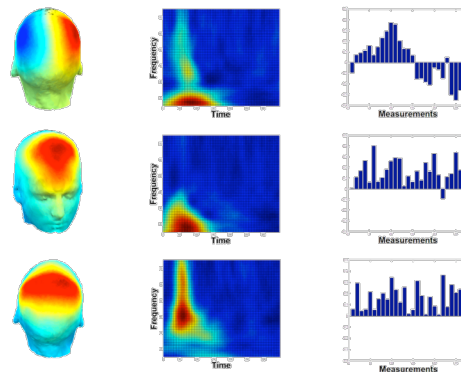
- Simultaneous analysis in 3 ways will fill in the gaps
- Our approach is faster than alternatives
- Specialized algorithm for large-scale problems
  - 500 x 500 x 500 with 99% missing data (1.25M nonzeros)
  - 1000 x 1000 x 1000 with 99.5% missing data (5M nonzeros)



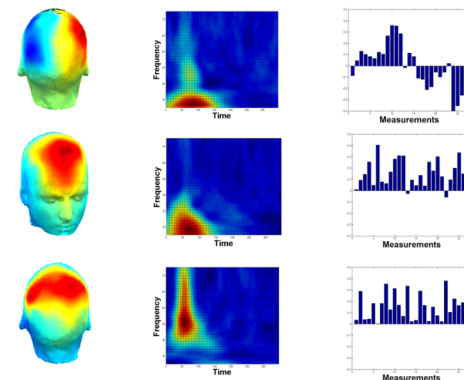
<http://www.madehow.com/>



No Missing Data

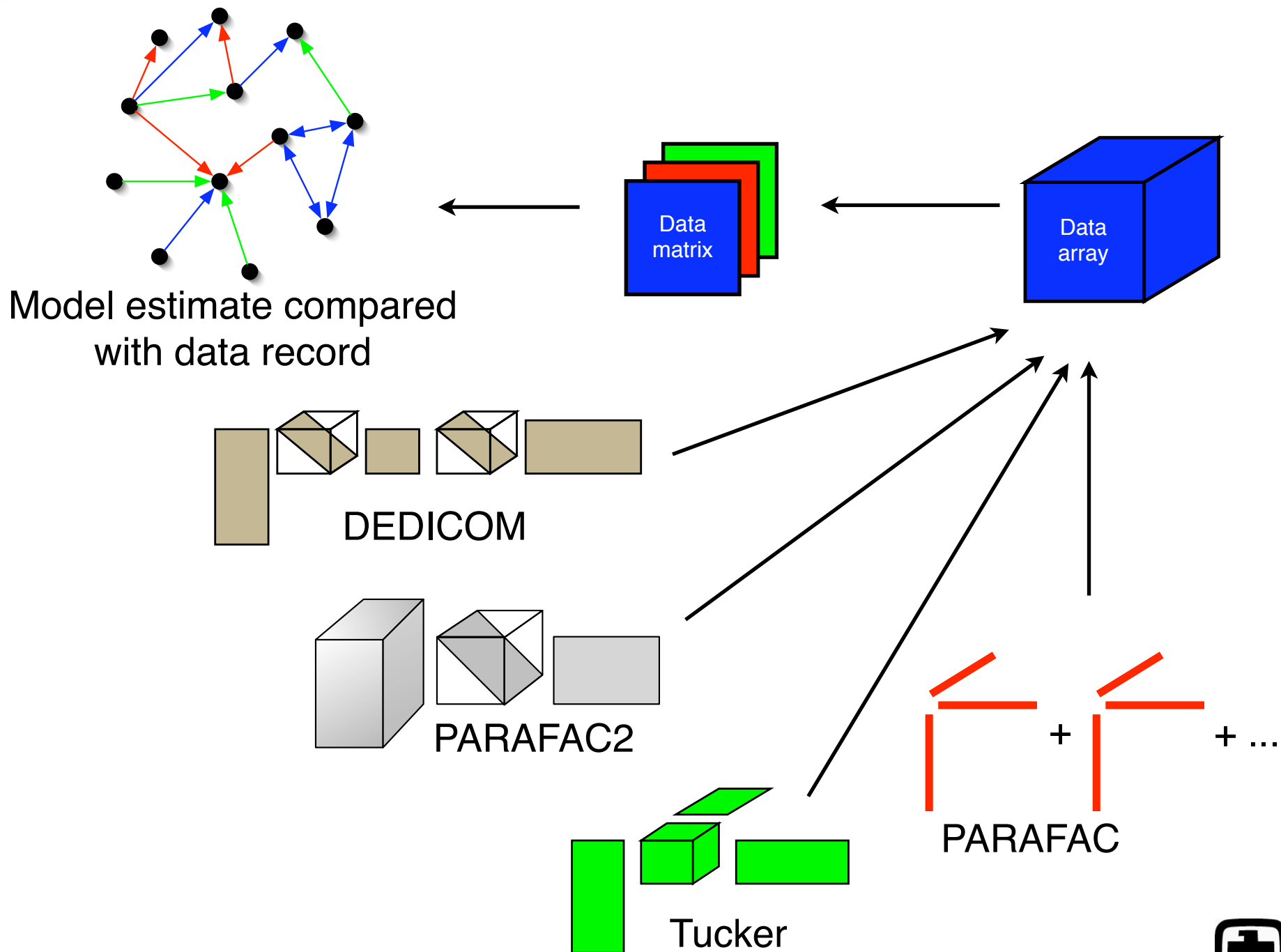


46% Missing Data



- Acar, Dunlavy, Kolda, Mørup, Scalable Tensor Factorization with Missing Data, SDM2010.
- Acar, Dunlavy, Kolda, Mørup, Scalable Tensor Factorization with Incomplete Data, in revision, 2010.

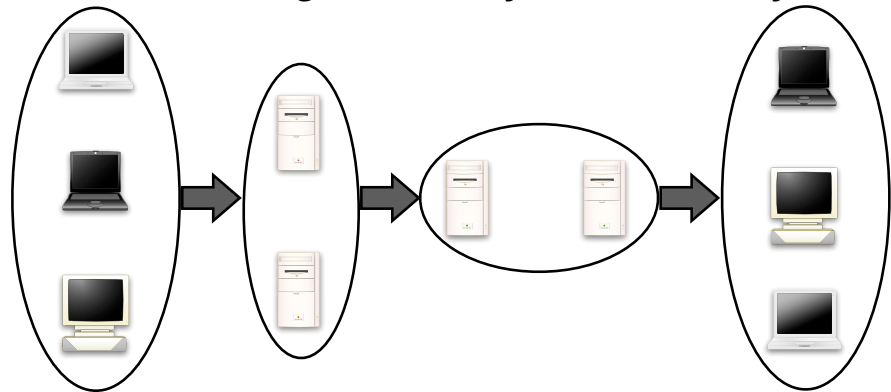
# Missing data facilitates another approach to anomaly detection



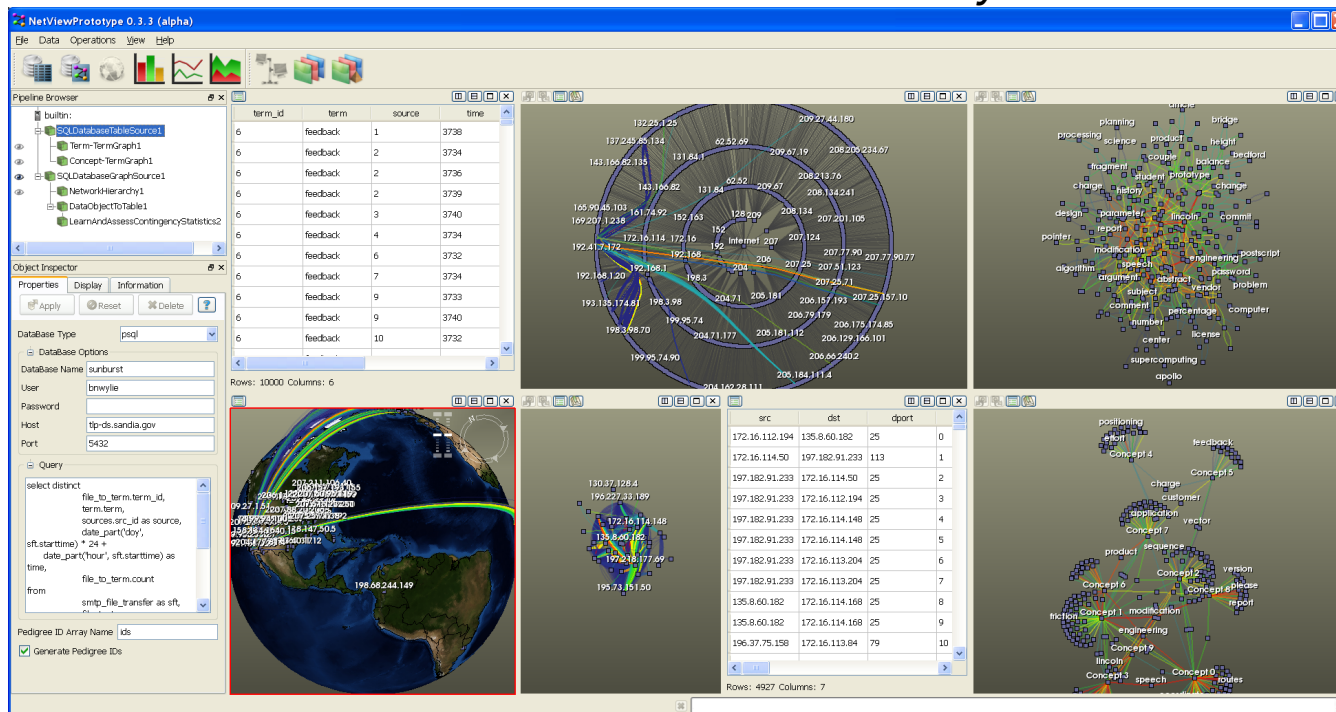
# Related Analysis Projects

- Multilingual document analysis and classification
- Uncovering plots buried in text (scenario discovery)
- IP address characterization (trace route analysis)
- Network traffic analysis (cyber, phone)
- Cyber data exfiltration analysis
- Link prediction
- Higher-order web link analysis

*Clustering nodes by their activity*

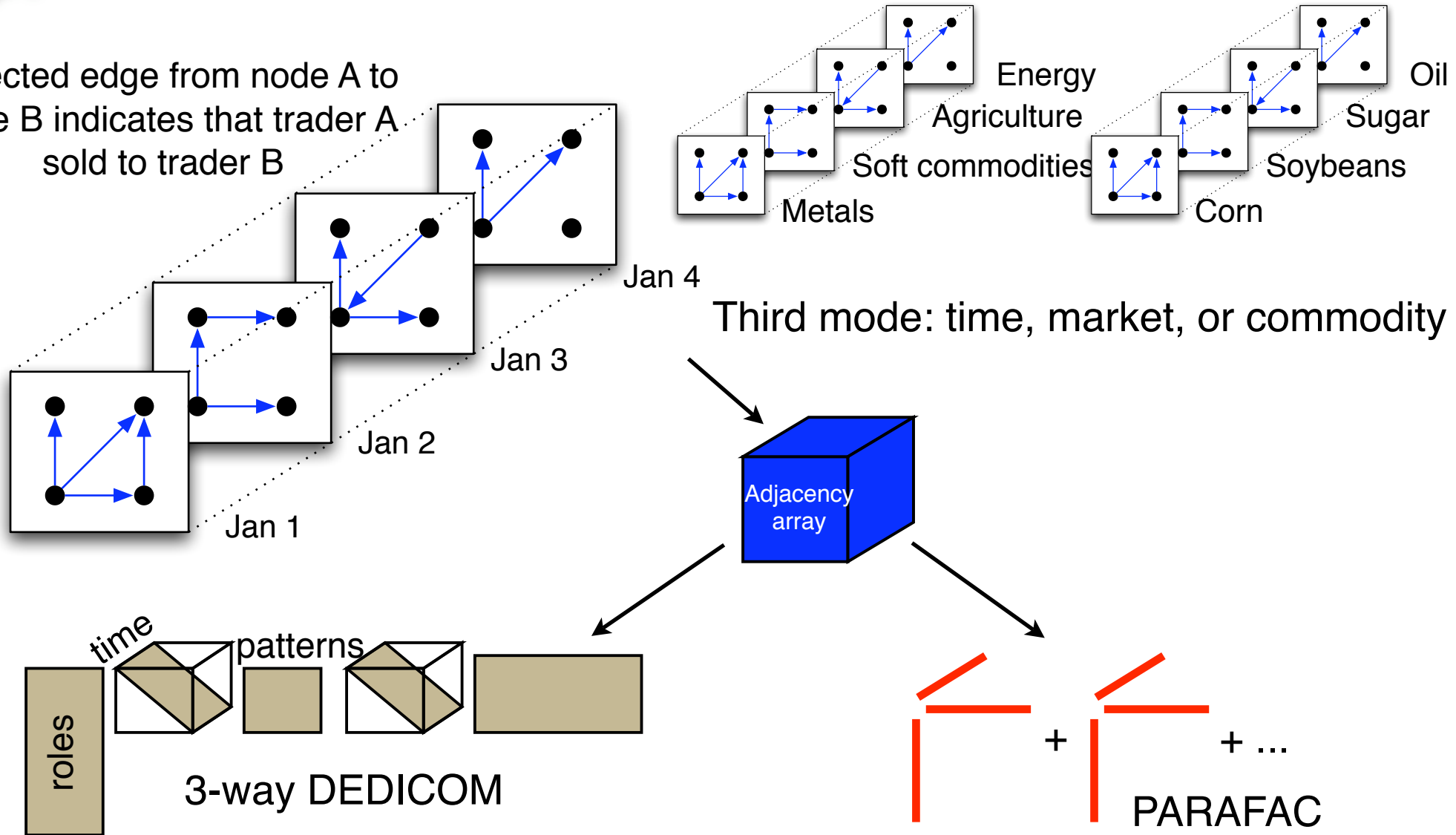


*Network traffic and content analysis*



# Preliminary Ideas for Trading Networks

Directed edge from node A to node B indicates that trader A sold to trader B



- Soft clustering of traders by their activity
- Aggregate trading patterns among clusters
- Behavior over time or by market
- Traders characterized by their “authority”
- Patterns in time or by market