

Data Mining on Attributed Relationship Graphs (ARGs)

SAND2007-5945C

LDRD Day

September 19, 2007

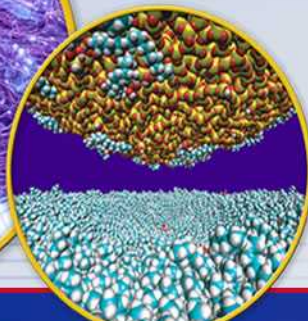
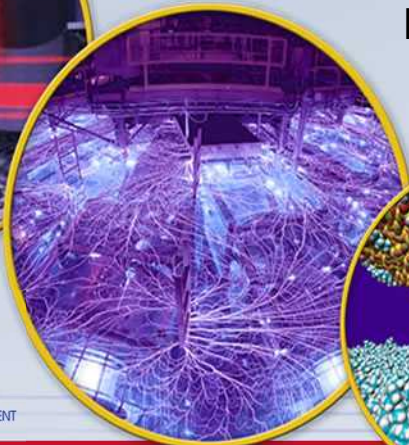
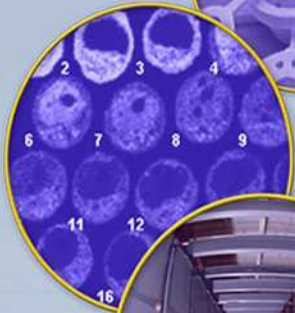
.....

Tamara G. Kolda

**Principal Investigator and PMTS
Informatics and Decision Sciences Dept. (Org. 8962)**

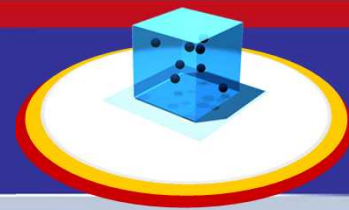
Team Members: Brett Bader (1416), Bruce Hendrickson (1415),
Ann Yoshimura (8116), Joe Kenny (8961), Travis Bauer (6341)

Plus collaborations with: Peter Chew (6341),
Danny Dunlavy (1411), Philip Kegelmeyer (8962)





Attributed Relationship Graphs



- **Graph**
 - Nodes represent entities such as people, places, or objects
 - Edges represent connections between entities
- **Attributed Relationship Graphs (ARGs)**
 - Nodes and edges have types, allowing multiple types in one graph
 - Both nodes and edges can have attributes (names, dates, etc.)
- **ARGs are used by intelligence analysts as a way to integrate data from disparate sources**

ARGs are also known as
Semantic Graphs



email



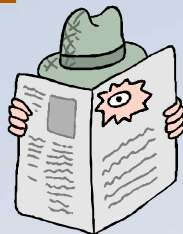
telephone



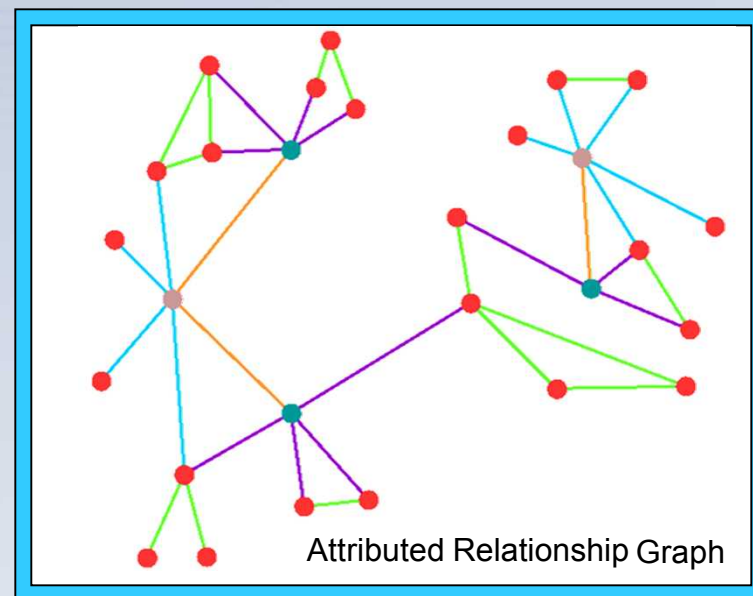
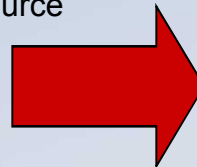
open source



cell phone

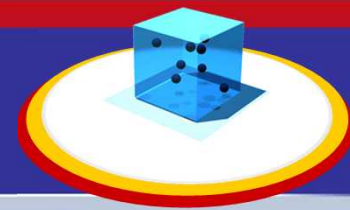


intelligence reports





LDRD Purpose: Develop Latent Semantic Analysis for ARGs

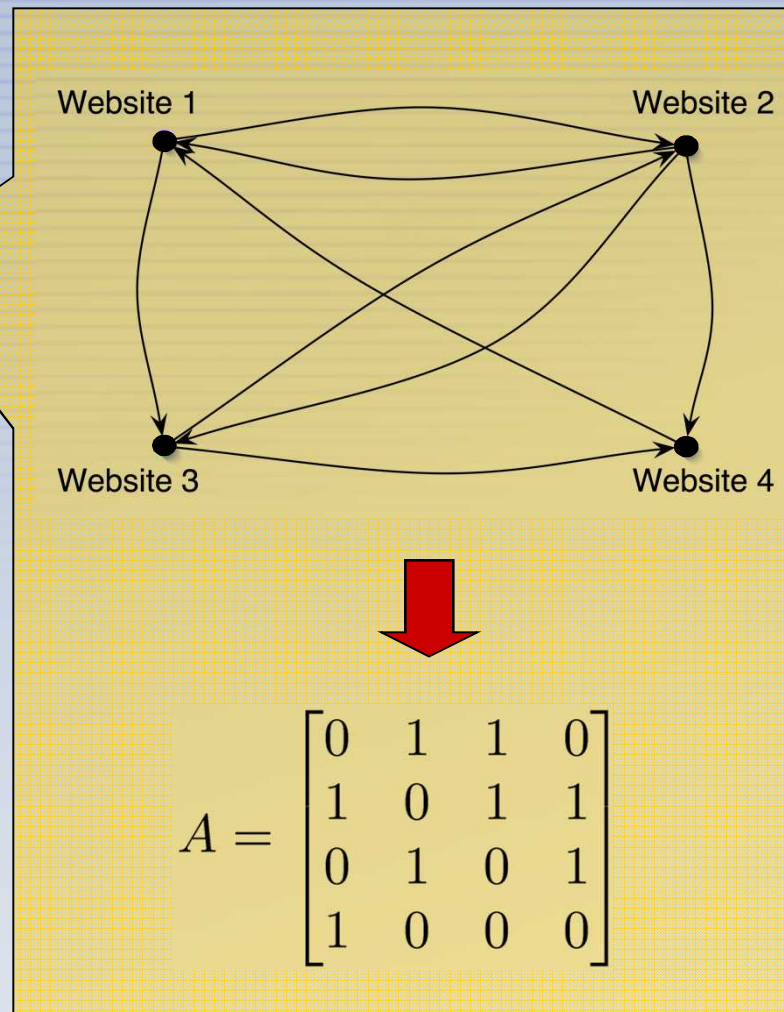


- **Matrix decompositions are a tool for graph analysis**

- Hubs & Authorities (HITS)
 - Matrix stores hyperlinks between web pages
 - Compute singular value decomposition (SVD)
 - Kleinberg (1998)
- Google PageRank
 - Stochastic matrix stores hyperlinks plus random jumps
 - Compute leading eigenvector
 - Page et al. (1998)
- Latent Semantic Indexing (LSI)
 - Matrix links documents and terms
 - Compute singular value decomposition (SVD)
 - Dumais et al. (1988)

- **Tensors provide a natural representation for ARGs**

- Tensor decompositions can be used for similar analyses!





LDRD Approach: Use Tensor Decompositions to Analyze ARGs

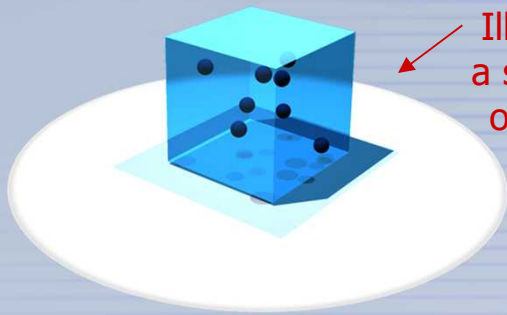
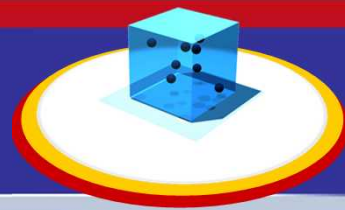


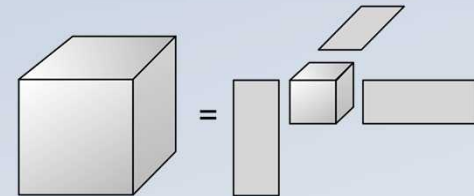
Illustration of
a sparse third-
order tensor

- A tensor is a multi-way array
- Not to be confused with tensor fields such as stress tensors.
- We have focused specially on sparse tensors!

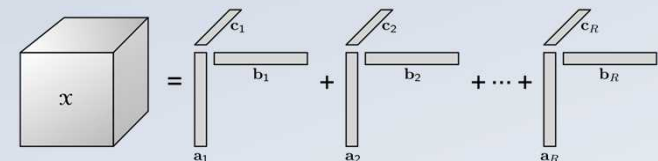
- Tensor decompositions reveal latent structure in the entries of a tensor

- Higher-order analogues of matrix SVD/PCA
- Date back to Hitchcock (1927)
- Popularized
 - 1970s in psychometrics
 - 1980s in chemometrics
 - 1990s in signal processing
 - 2000s in data mining, etc.

Tucker (1966)

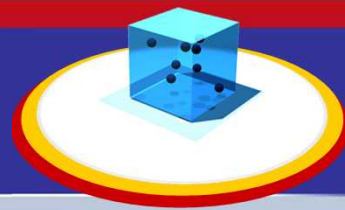


CANDECOMP/
PARAFAC (1970)





TOPHITS – A Three-Dimensional View of the Web



Endangered Species
Animals today are being threatened by a variety of environmental pressures. For example, the jaguar is losing prime habitat in the world. Zoos are trying to raise awareness of their plight.

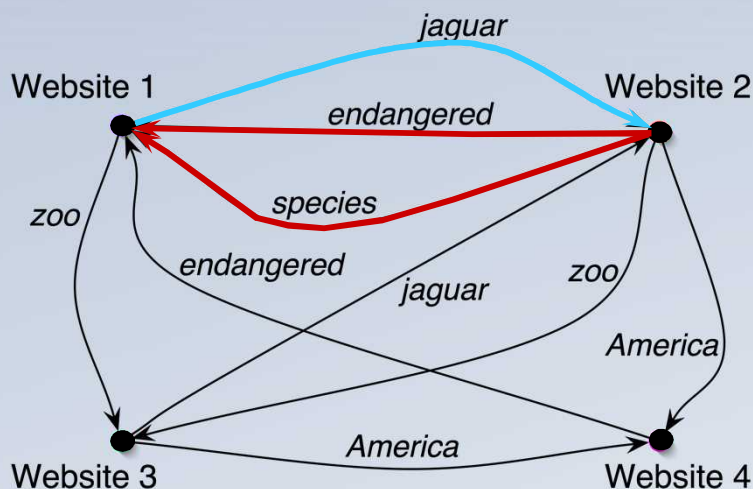
Jaguar FAQ
Jaguars are an endangered species that live in the tropical rain forests of Central and South America. They live about 11 years in the wild and up to 22 years at a zoo.

Rain Forest Zoo
We have a new exhibit opening next month highlighting the endangered species of the Americas, including the jaguar.

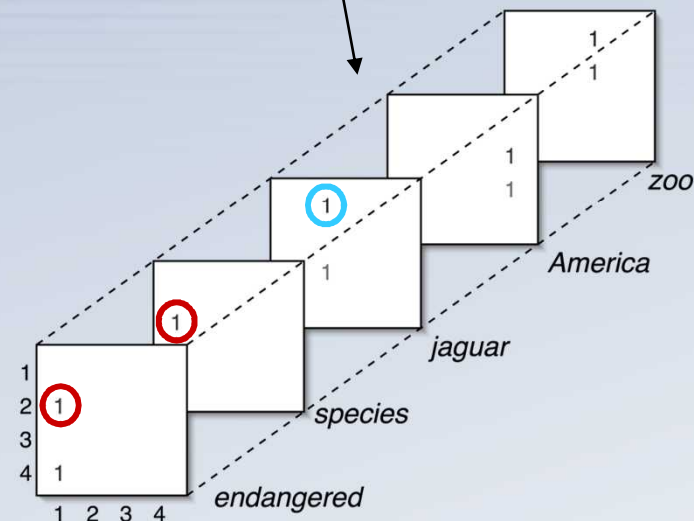
Online Atlas
View maps of animal habitats from around the world, including those of endangered animals in North, South, and Central America.

Tensor Definition

$$x_{ijk} = \begin{cases} 1 & \text{if page } i \rightarrow \text{page } j \\ & \text{with term } k \\ 0 & \text{otherwise} \end{cases}$$

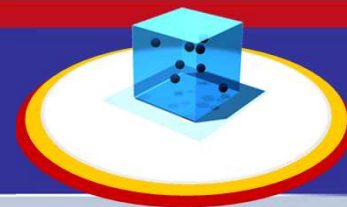


Observe that this tensor is very sparse!





TOPHITS Terms & Authorities on Sample Data



Main Idea: Extend the idea behind the HITS model to incorporate term (i.e., topical) information.

$$x_{ijk} = \begin{cases} \frac{1}{\log(w_k)+1} & \text{if } i \rightarrow j \text{ with term } k \\ 0 & \text{otherwise} \end{cases}$$

w_k = # unique links using term k

1st Principal Factor			
.23	JAVA	.86	java.sun.com
.18	SUN	.38	developers.sun.com

2nd Principal Factor			
.17	PLATF	.99	www.lehigh.edu
.16	SOLAR	.20	NO-READABLE-TEXT

3rd Principal Factor			
.15	EDITIO	.16	SEAR
.15	DOWN	.16	NEWS

4th Principal Factor			
.12	SOFTV	.16	COMF
.12	NO-RE	.12	LEHIC
.12	WEB	.23	CITIZ

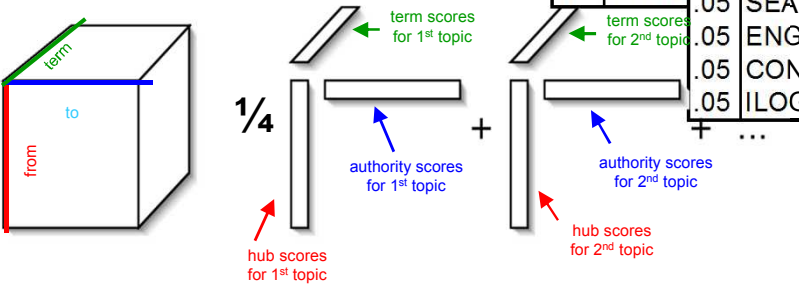
6th Principal Factor			
.11	DEVEL	.26	PRESIDENT
.11	LINUX	.22	OTHE
.11	RESO	.19	CENT

12th Principal Factor			
.10	DOWN	.15	PUBLI
.10	TECH	.15	U.S
.10	TECH	.15	WELC

13th Principal Factor			
.13	FREE	.13	PRES
.13	BUDG	.06	TREE
.13	HOUS	.08	DECIS

16th Principal Factor			
.13	OFFIC	.05	SEAR
.13	OFFIC	.05	ENGIN
.13	OFFIC	.05	CONT

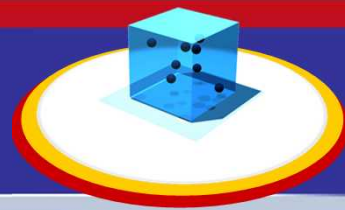
19th Principal Factor			
.22	TAX	.73	www.irs.gov
.17	TAXES	.43	travel.state.gov
.15	CHILD	.22	www.ssa.gov
.15	RETIREMENT	.08	www.govbenefits.gov
.15	BENEFITS	.06	www.usdoj.gov
.15	STATE	.03	www.census.gov
.15	POLICY	.03	www.usmint.gov
.14	INCOME	.02	www.nws.noaa.gov
.13	SERVICE	.02	www.gsa.gov
.13	REVENUE	.02	www.gsa.gov
.12	CREDIT	.01	www.annualcreditreport.com



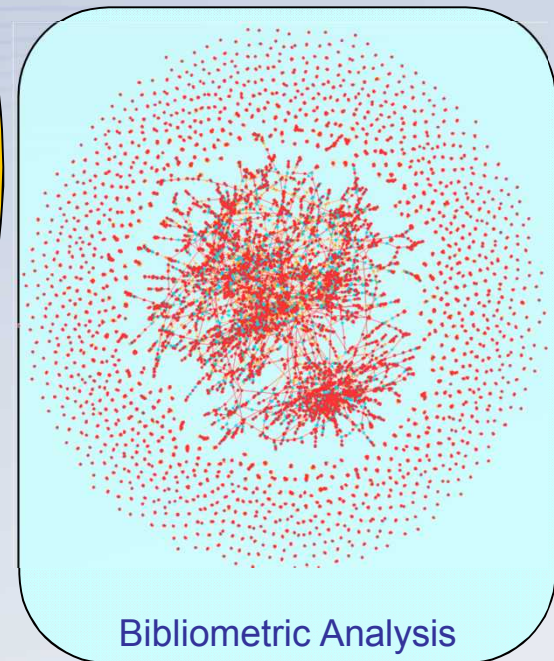
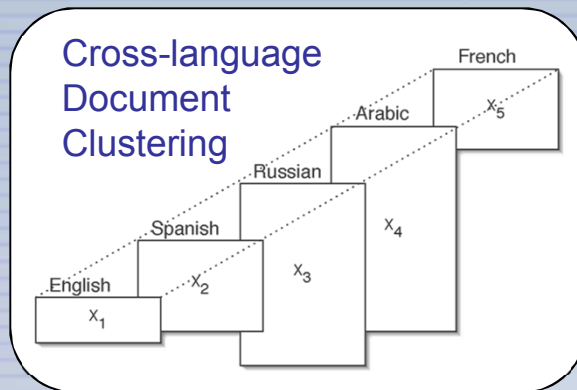
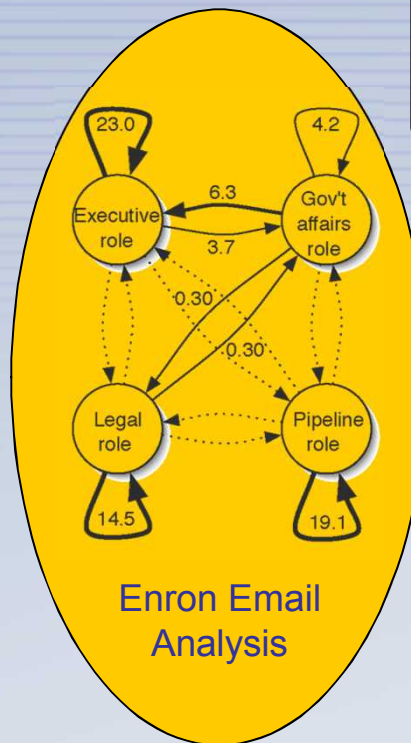
Tensor PARAFAC



Significance: Tensors are a New Approach to Graph & Data Analysis

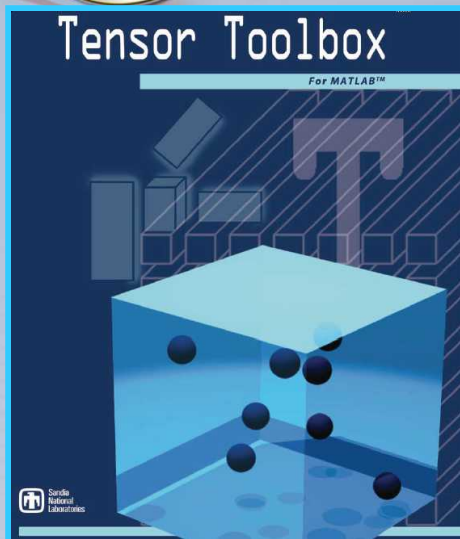
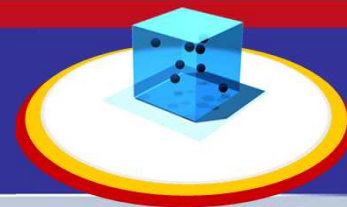


- **TOPHITS for higher-order web link analysis**
- **Enron email text analysis**
 - With Mike Berry, U. Tennessee
- **Tutorial: Mining large time-evolving data using matrix and tensor tools**
 - With C. Faloutsos, J. Sun (CMU)
- **Cross-language document clustering**
- **Enron email temporal pattern analysis**
 - With R. Harshman, U. Western Ontario
- **Bibliometric analysis**
- **Social network analysis**

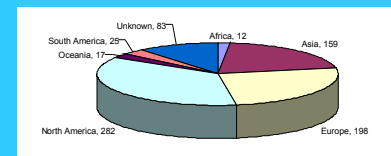




Results: New Algorithms & Software for Data Analysis

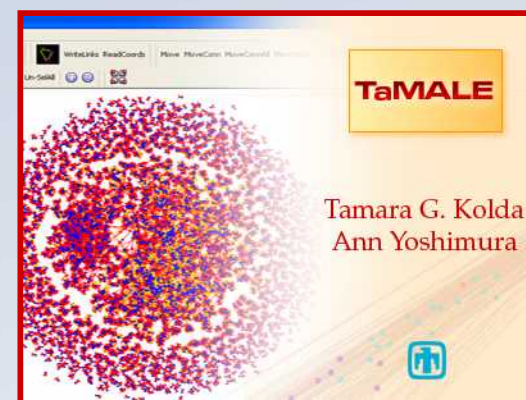


- **Tensor Toolbox for MATLAB**
- **B. W. Bader (1416) & T. G. Kolda (8962)**
- **Version 2 released externally Sept. 2006; 775+ Downloads from all over the world**
- **Publications on the toolbox in ACM Transactions on Mathematics Software and SIAM J. Scientific Computing**
- **Unique capability: Methods for large-scale, sparse and structured tensors**



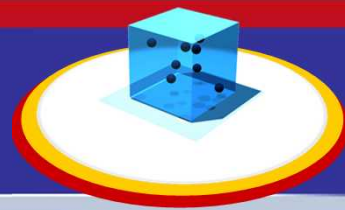
Downloads by Continent

- **TaMALE = Tensor Multi-Attribute Link Explorer**
- **T. G. Kolda (8962) and A. Yoshimura (8116)**
- **Import, visualization, and analysis of ARGs**
- **Released internally in March 2007**
- **Version 2.0 now available**

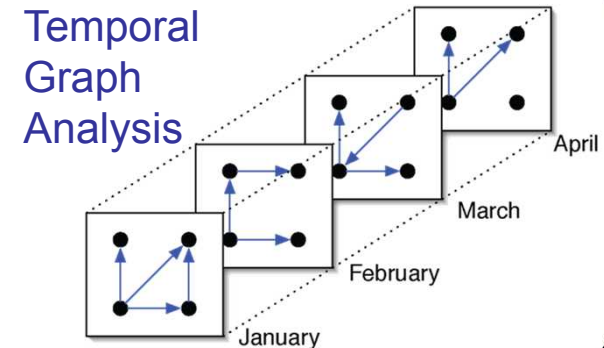




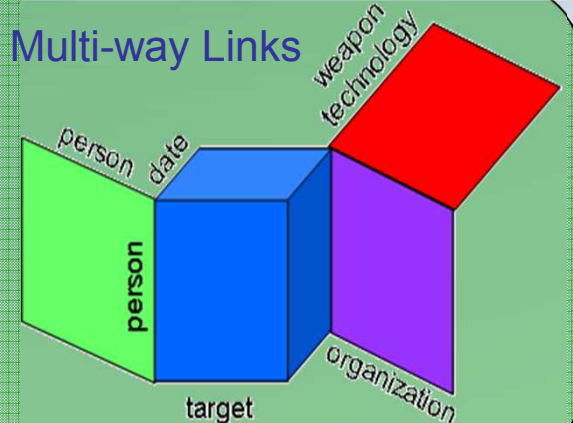
This project is over, but there is more graph analysis on the horizon



- **Enabling All-Threat Analysis Through Intelligent Filtering of Network Traffic**
 - DHS LDRD FY07-09
 - PI: Jamie Van Randwyk (8965)
- **Network Discovery, Prediction and Disruption**
 - Grand Challenge LDRD FY08-10
 - PI: Bruce Hendrickson (1415)
- **Leveraging Multi-way Linkages on Heterogeneous Data**
 - EPS LDRD FY08-10
 - PI: Kolda



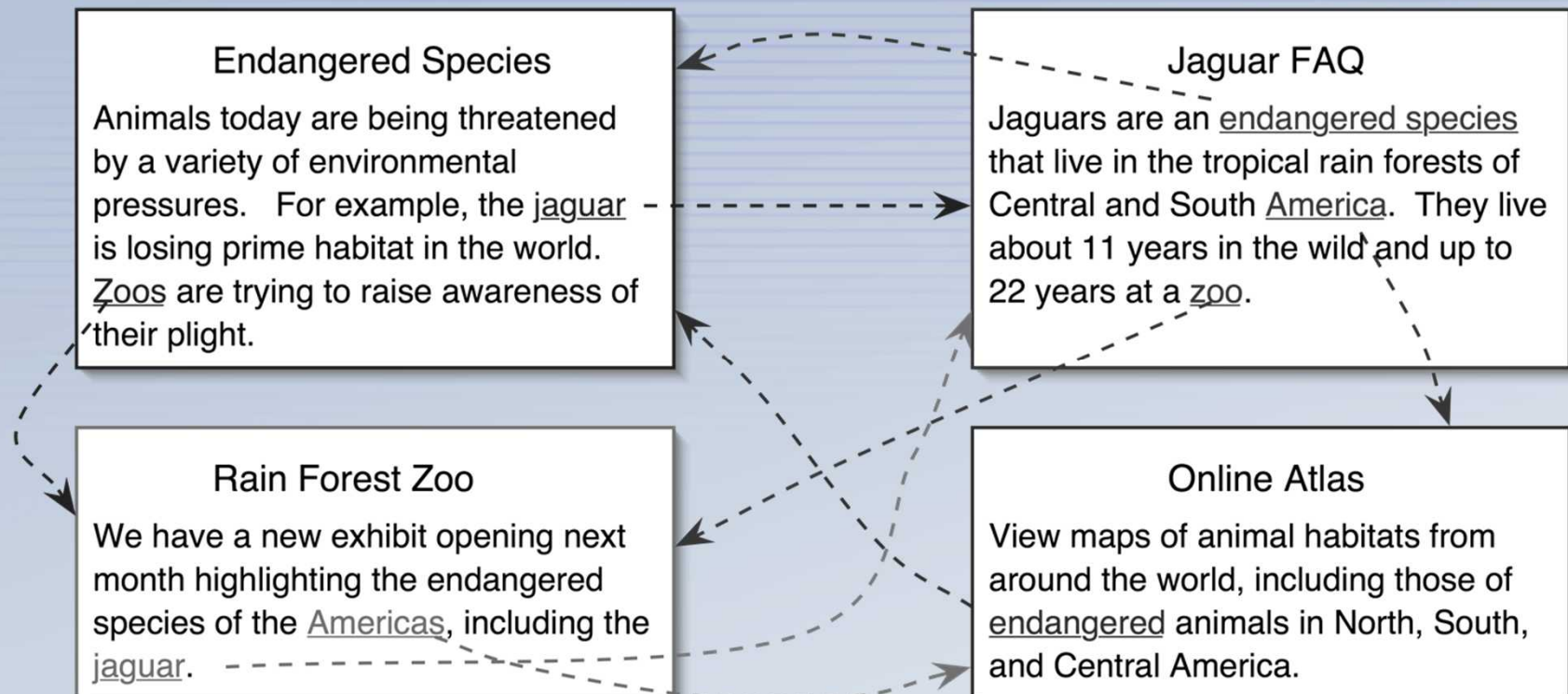
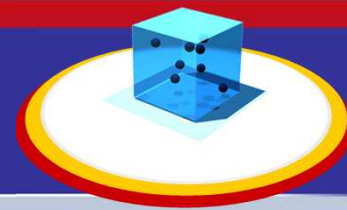
Multi-way Links



For more information, contact:
Tammy Kolda
294-4769, tgkolda@sandia.gov

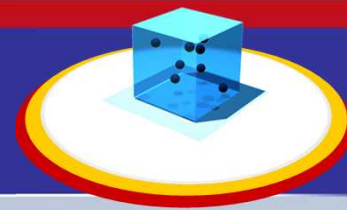


Detail: Example Web Pages



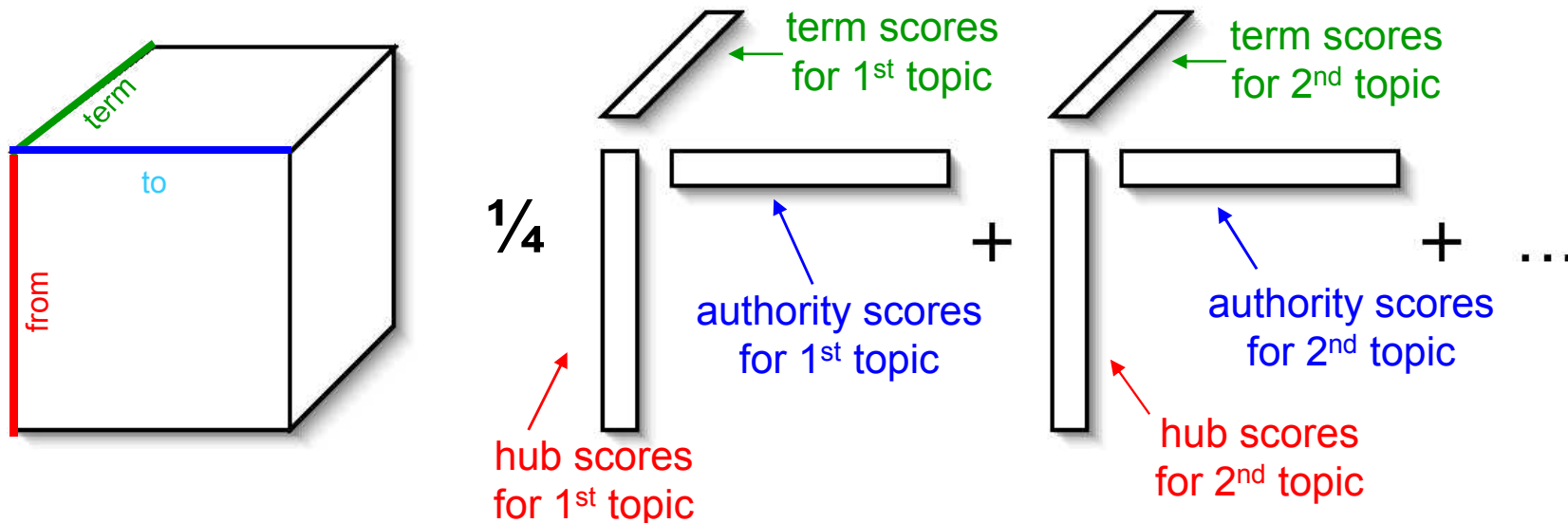


DETAIL: Topical HITS (TOPHITS) Tensor Decomposition



Main Idea: Extend the idea behind the HITS model to incorporate term (i.e., topical) information.

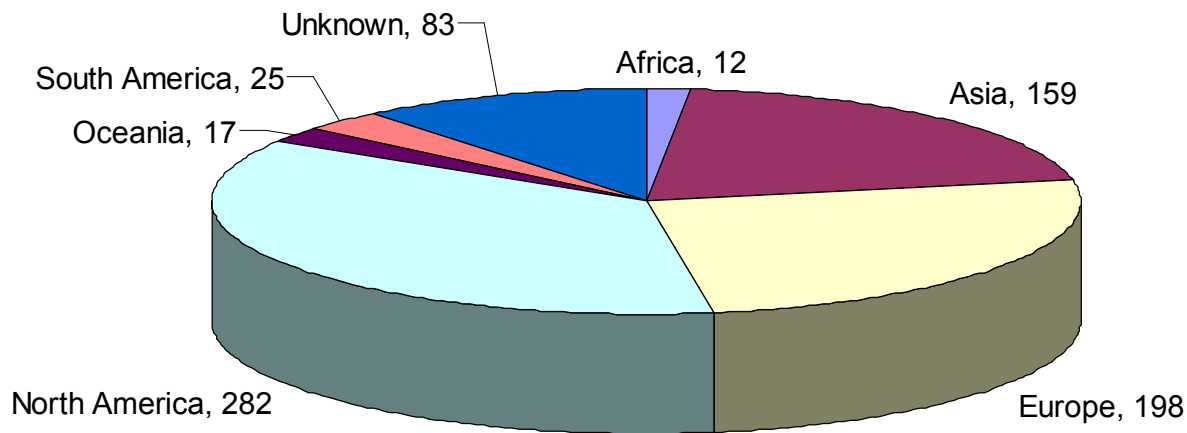
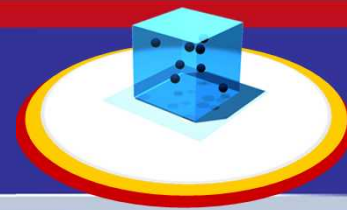
$$\mathcal{X} \approx [\lambda ; \mathbf{H}, \mathbf{A}, \mathbf{T}] = \sum_{r=1}^R \lambda_r \mathbf{h}_r \circ \mathbf{a}_r \circ \mathbf{t}_r$$



T. G. Kolda, B. W. Bader, and J. P. Kenny. Higher-order web link analysis using multilinear algebra. In ICDM 2005: Proceedings of the 5th IEEE International Conference on Data Mining, pages 242–249, November 2005. (doi:10.1109/ICDM.2005.77)



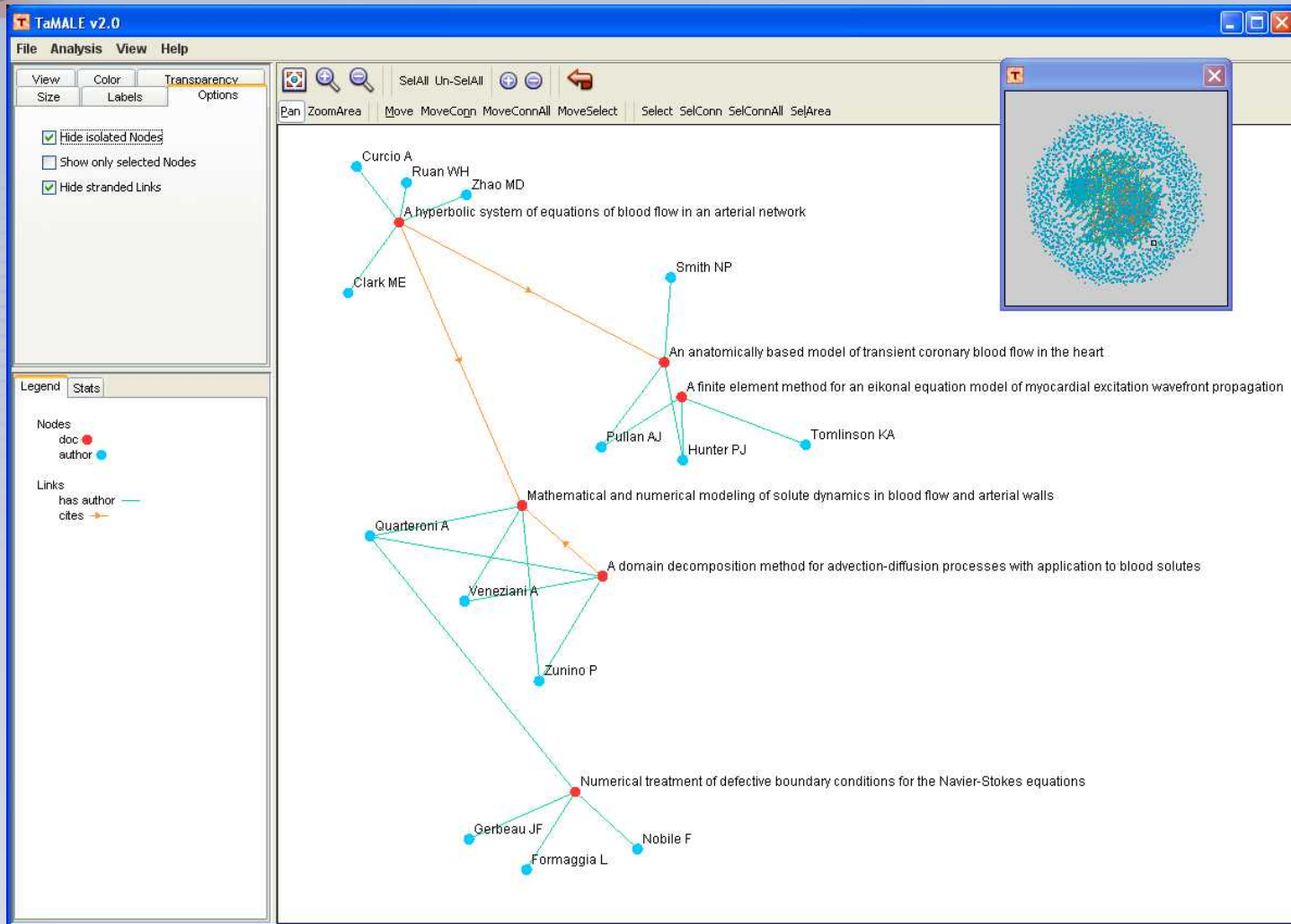
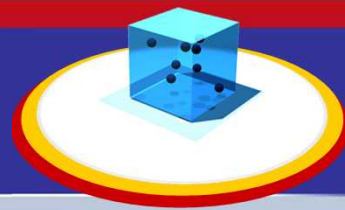
Detail: Tensor Toolbox Downloads by Continent



Unique downloads (by IP address) since Sept. 2006

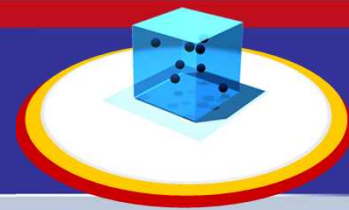


Detail: Bibliometric Graph in TaMALE





Detail: Refereed Papers



- **Journal Papers**

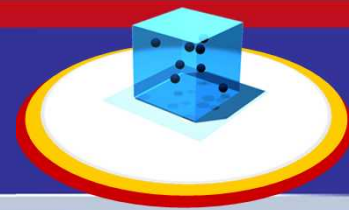
- Brett W. Bader and Tamara G. Kolda. **Efficient MATLAB computations with sparse and factored tensors**. *SIAM Journal on Scientific Computing*, Accepted for publication, July 2007.
- Brett W. Bader and Tamara G. Kolda. **Algorithm 862: MATLAB tensor classes for fast algorithm prototyping**. *ACM Transactions on Mathematical Software*, 32(4):635–653, December 2006. ([doi:10.1145/1186785.1186794](https://doi.org/10.1145/1186785.1186794))

- **Conference & Workshop Papers**

- Tamara G. Kolda, Brett W. Bader, and Joseph P. Kenny. **Higher-order web link analysis using multilinear algebra**. In *ICDM 2005: Proceedings of the 5th IEEE International Conference on Data Mining*, pages 242–249, November 2005. ([doi:10.1109/ICDM.2005.77](https://doi.org/10.1109/ICDM.2005.77))
- Brett W. Bader, Michael W. Berry, and Murray Browne. **Discussion tracking in Enron email using PARAFAC**. Text Mining Workshop, April 2007.
- Peter A. Chew, Brett W. Bader, Tamara G. Kolda, and Ahmed Abdelali. **Cross-language information retrieval using PARAFAC2**. In *KDD '07: Proceedings of the 13th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pages 143–152. ACM Press, 2007. ([doi:10.1145/1281192.1281211](https://doi.org/10.1145/1281192.1281211))
- Brett W. Bader, Richard Harshman, and Tamara G. Kolda. **Temporal analysis of semantic graphs using ASALSAN**. In *ICDM 2007: Proceedings of the 7th IEEE International Conference on Data Mining*, November 2007. To appear.



Detail: Other Papers and Presentations



- **Submitted Journal Papers**

- Tamara G. Kolda and Brett W. Bader. **Tensor decompositions and applications**. Submitted to *SIAM Review*, August 2007.

- **Other papers**

- Teresa M. Selee, Tamara G. Kolda, W. Philip Kegelmeyer, and Joshua D. Griffin. **Extracting clusters from large datasets with multiple similarity measures using IMSCAND**. To appear in CSRI2007 proceedings, August 2007.
- Daniel M. Dunlavy, Tamara G. Kolda, and W. Philip Kegelmeyer. **Multilinear algebra for analyzing data with multiple linkages**. Technical Report SAND2006-2079, Sandia National Laboratories, Albuquerque, NM and Livermore, CA, April 2006. Revised version to appear in *Array Based Graph Algorithms*
- Tamara G. Kolda. **Multilinear operators for higher-order decompositions**. Technical Report SAND2006-2081, Sandia National Laboratories, Albuquerque, NM and Livermore, CA, April 2006.

- **Refereed Tutorial**

- Christos Faloutsos, Tamara G. Kolda, and Jimeng Sun, **Mining Large Time-evolving Data Using Matrix and Tensor Tools**. Tutorial at SDM07, SIGMOD07, ICML07, KDD07.