# NUBBI: An introduction and a non-examination of Sudanese Newspapers

## Peter Landwehr

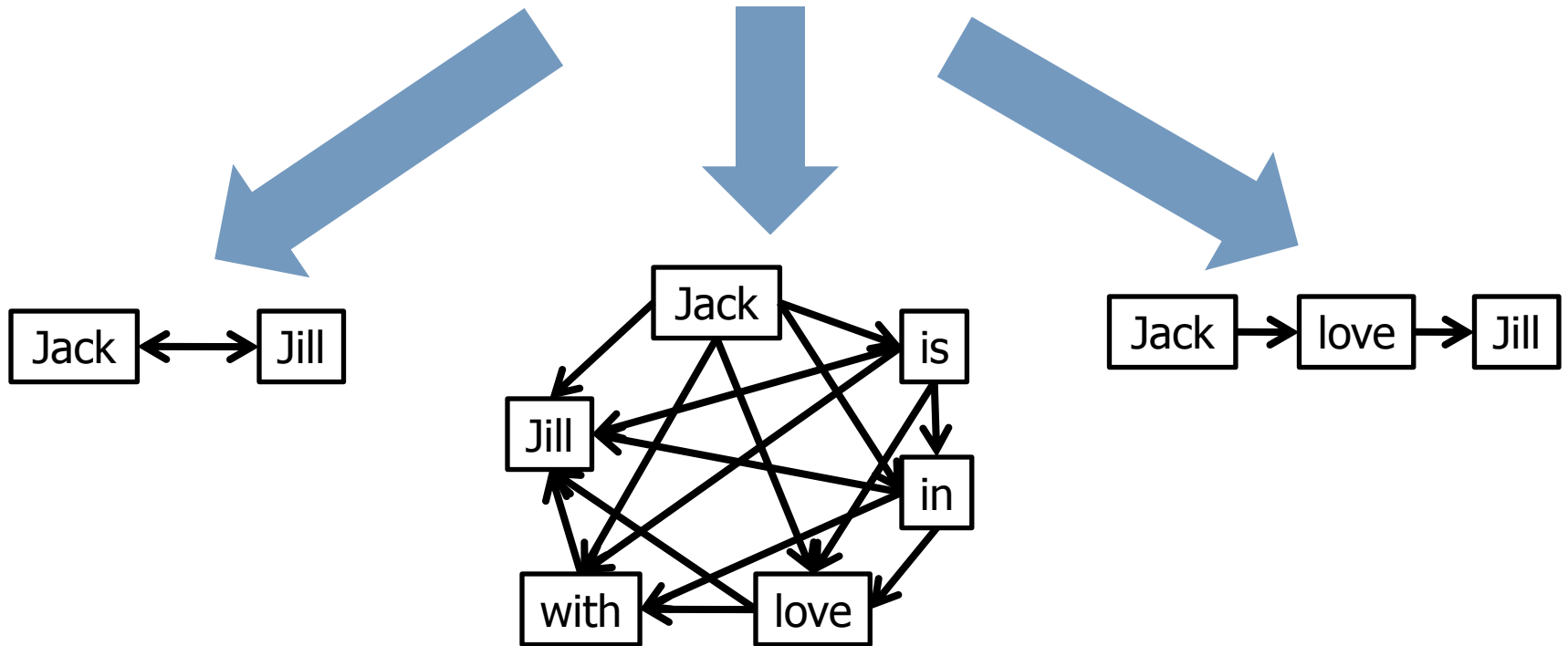plandweh@cs.cmu.edu

plandweh@cs.cmu.edu

# What I'm covering

1. A difficulty with semantic networks
2. One way to try and resolve it.
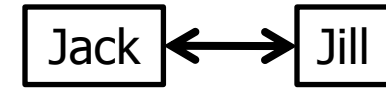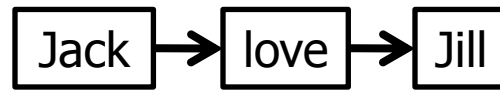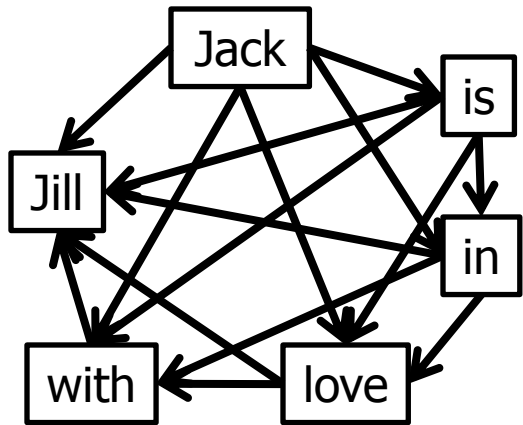3. My preliminary work with this method.

# Semantic Networks!

Jack is in love with Jill.
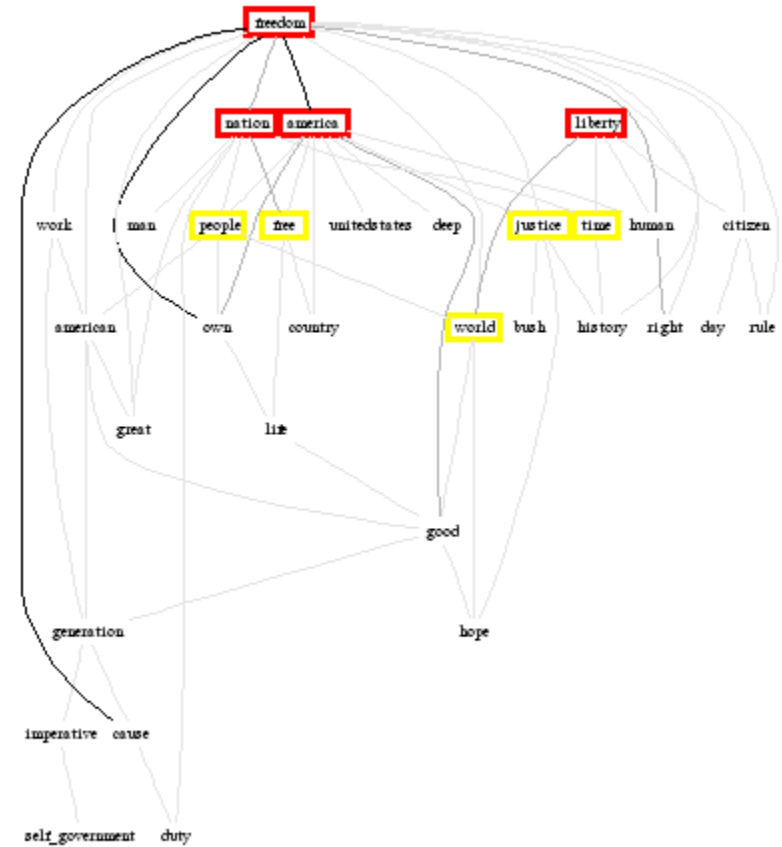
**Peter Landwehr**

# Map Theory



- Words in a document gain significance based on their juxtaposition with other words.

- Situate words in a network based on windowed proximity.

- Words can be classed into particular roles based on these network statistics.

[Carley 1997] [Carley & Kaufer 1993]

# Centering Resonance Analysis

- Drop all terms that aren't nouns

- Link nouns, noun phrases

- Compare betweenness centralities of words in networks to determine "resonance".



File: 2005 Bush.cra    Cutoff 0.025

[Corman et al. 2002]
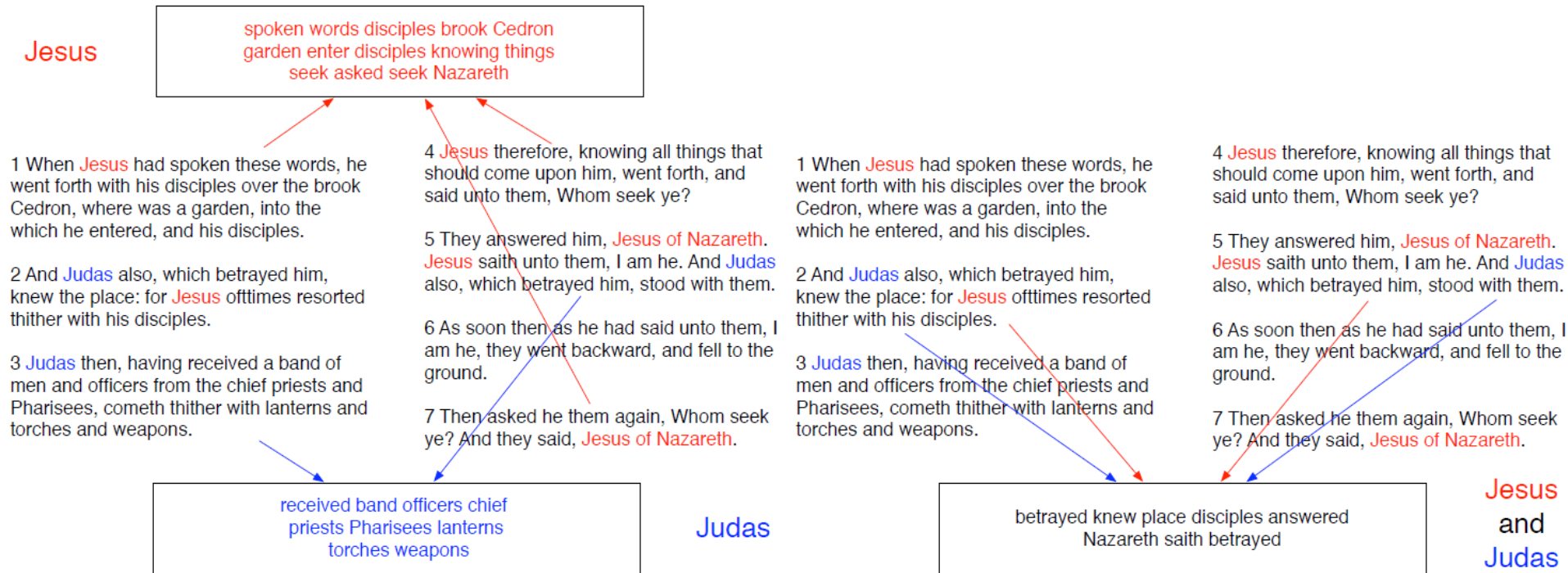
# An appealing improvement

# One way towards this: NUBBI

- Networks Uncovered By Bayesian Inference
- "Connections between the lines: augmenting social networks with text" by Chang, Boyd-Graber, & Blei. *KDD 2009*

# Stealing a diagram…

# Latent Dirichlet Allocation

- Decompose documents as mixtures of *N* topics
- Randomly assign topics to documents
- Based on random topic mixtures, randomly assign words to topics.
- Expectation Maximization or Gibbs sampling to converge.



[Blei et al. 2003][Blei 2011]

# Idealized description of the algorithm

- Assumes text is grouped into entity and pair contexts
- For each entity context, NUBBI assigns a mixture of entity topics.
- For each pair context, NUBBI splits the words into those describing entities and those describing the pair.
- For each pair-specific word in a pair-specific context, NUBBI assigns a pair topic.

# Idealized description of the analyst's role

- Define entity and pair counts
- Define three hyperparameters for topic proportions required by the model
  - Topic mixture
  - Pair mixture
  - Entity/pair selector mixture
- Select appropriate windowing parameters and clean text.

# Issues with NUBBI

- Text in general
  - Deciding on context sizes. (e.g. Bible = verses)
  - Text cleaning.
  - Evaluating goodness of results.
- NUBBI-specific
  - # of entity topics = ?
  - # of pair topics = ?
  - Optimal mixing values
  - The best way of assigning multi-context text.
  - Relationships described in text can go beyond pairs to triples,

# A personal test

- Eight years of news articles from the Sudan Tribune, documenting key activities in the region
- Detailed thesauri containing ~17000 different actors.
- These texts have been used for a variety of CASOS-affiliated work
  - T. Van-Holt and J. C. Johnson, "A Text and Network Analysis of Natural Resource Conflict in Sudan," in *Proceedings Sunbelt XXXI*, St. Petersburg Beach, Florida, USA, 2011.
  - J. Diesner and K. M. Carley, "Mapping Socio-Cultural Networks of Sudan From Open-Source, Large-Scale Text Data," in *Proceedings of the 29th Annual Conference of the Sudan Studies Association*, West Lafayette, Indiana, 2010.

# Basic Notes

- Cleaned the text using AutoMap
  - Used standard AutoMap cleaning methods (clean white space, fix a predefined set of typos, convert British to American spellings, expand contractions & abbreviations, convert to lowercase, best-effort resolve pronouns and delete unresolved, remove punctuation and numbers.
  - "A stop list for general text" by Fox. *ACM SIGIR*, volume 24, issue 1-2. Fall 1989.
  - Specialized AutoMap thesaurus on Sudan. (Possible issue in order of application)

- Used the R implementation of NUBBI ("lda" package)

- 12 entity topics, 6 pair topics, windows of size 5.

# Corpora Notes

| Year | Articles in Corpora | Entities in Corpora | Meaningful Pairs | Vocabulary (Total) |
|------|---------------------|---------------------|------------------|--------------------|
| 2003 | 616 | 379 | 352 | 2007 (6539) |
| 2009 | 1056 | 536 | 460 | 3571 (13283) |
| 2010 | 1063 | 481 | 472 | 3641 (13644) |

# Topic Overlaps

- Pair topics
  - 2003 has no overlaps with 2008 or 2009
  - In an optimal matching, the 2008 and 2009 topics overlap by 24.6% (Top 25 words/topic.)
  - Negligible matches between entity pairs in each pair topic per year. (Top 25 entities/topic

- Entity topics
  - 20-24% optimal matching overlap between all entity topics. This says more about LDA than it does about our quality. (Top25 words/topic)
  - 25% and 22% overlaps of entities between 2003 and 2008 and 2009, but 35% overlap between 2008 and 2009

# Future Work

- Make this preliminary work less preliminary by…
  - improving data cleaning
  - Improve model parameter choices
  - Full exploration of all years of the corpus
- Improve the "NUBBI experience" by
  - Integrating it with a text cleaning tool (AutoMap) for easier use.
  - Merge with HDPs to avoid having to set topic counts.
  - Experiment with other corpora to develop deeper recommendations for tuning.

# References

1. D. M. Blei, "Introduction to Probabilistic Topic Models."
2. D. M. Blei, A. Y. Ng, and M. I. Jordan, "Latent Dirichlet Allocation," *Journal of Machine Learning Research*, vol. 3, pp. 993-1022, Jan. 2003.
3. K. M. Carley, "Extracting team mental models through textual analysis," *J. Organiz. Behav.*, vol. 18, pp. 533–558, 1997.
4. K. M. Carley and D. Kaufer, "Semantic Connectivity: An Approach for Analyzing Symbols in Semantic Networks," *Communication Theory*, vol. 3, no. 3, pp. 183–213, 1993.
5. J. Chang, J. Boyd-Graber, and D. M. Blei, "Connections between the lines: augmenting social networks with text," in *Proceedings of the 15th ACM SIGKDD international conference on Knowledge discovery and data mining*, Paris, France, 2009, pp. 169-178.
6. S. R. Corman, T. Kuhn, R. D. Mcphee, and K. J. Dooley, "Studying Complex Discursive Systems.," *Human Communication Research*, vol. 28, no. 2, pp. 157-206, 2002.

# The end

**Carnegie Mellon**