

# Keyphrase Extraction in Citation Networks: How do Citation Contexts Help?

Cornelia Caragea

Computer Science and Engineering  
University of North Texas

November 21, 2014

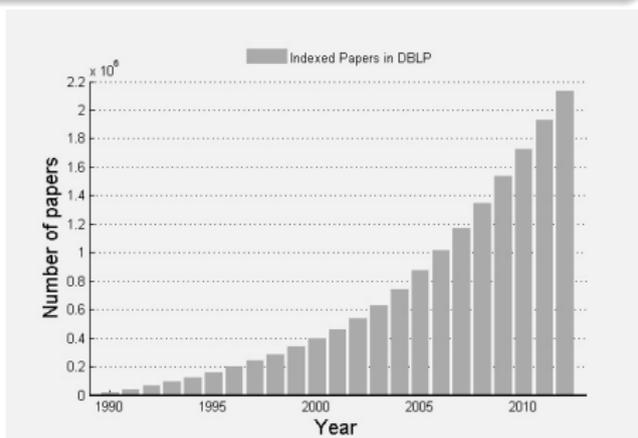
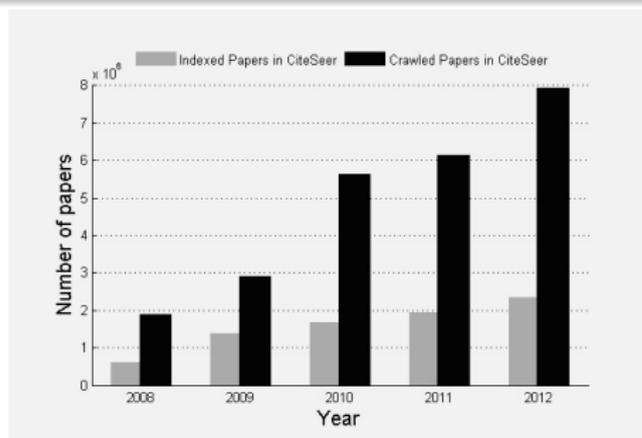
**MLg** Machine Learning Group  
*UNT Computer Science and Engineering*

# Scholarly Big Data

## Large number of scholarly documents on the Web

- PubMed currently has over 24 million documents
- Google Scholar is estimated to have 160 million documents

The growth in the number of papers indexed by CiteSeer and DBLP:



- Navigating in these digital libraries has become very challenging.

# Keyphrases

- **Keyphrases** provide a concise summary of a document

Example of keyphrases: A snippet from the 2010 best paper award winner in the WWW conference - abstract and author-input keyphrases

*Factorizing Personalized **Markov Chains** for **Next-Basket Recommendation***  
by Rendle, Freudenthaler, and Schmidt-Thieme

“**Recommender systems** are an important component of many websites. Two of the most popular approaches are based on **matrix factorization** (MF) and **Markov chains** (MC). MF methods learn the general taste of a user by factorizing the matrix over observed user-item preferences. [...] In this paper, we present a method bringing both approaches together. Our method is based on personalized transition graphs over underlying **Markov chains**. [...] We show that our factorized personalized MC (FPMC) model subsumes both a common **Markov chain** and the normal **matrix factorization** model. For learning the model parameters, we introduce an adaption of the Bayesian Personalized Ranking (BPR) framework for sequential basket data. [...]”

- Useful in applications such as **topic tracking, information filtering and search, query formulation, document clustering, classification, and summarization.**

# Keyphrase Extraction

- Manually annotated keyphrases are not always provided with the documents:
  - Need to be gleaned from the content of documents.
- Given the size of current scholarly digital libraries, manual annotation of keyphrases has become infeasible.
- Hence, accurate approaches are required for **keyphrase extraction** from research documents.
  - **Keyphrase extraction** is defined as the problem of automatically extracting **descriptive phrases** or **concepts** from documents.

# Previous Approaches to Keyphrase Extraction

- Many approaches have been studied:
  - Supervised approaches [Frank et al., 1999; Turney, 2000; Hulth, 2003]
    - Formulated as binary classification, where candidate phrases are classified as either positive (i.e., keyphrases) or negative (i.e., non-keyphrases)
  - Unsupervised approaches
    - Formulated as a ranking problem, where keyphrases are ranked using various measures such as tf, tf-idf, PageRank scores and other centrality measures [Mihalcea and Tarau, 2004; Wan and Xiao, 2008; Liu et al., 2010; Lahiri et al., 2014]
- Generally, previous approaches
  - Use only the textual content of the target document [Mihalcea and Tarau, 2004; Liu et al., 2010].
  - Incorporate a local neighborhood of a document for extracting keyphrases [Wan and Xiao, 2008]
    - However, the neighborhood is limited to textually-similar documents.

# Our Questions

- In addition to a document's textual content and textually-similar neighbors, are there other informative neighborhoods that exist in research document collections?
- Can these neighborhoods improve keyphrase extraction?

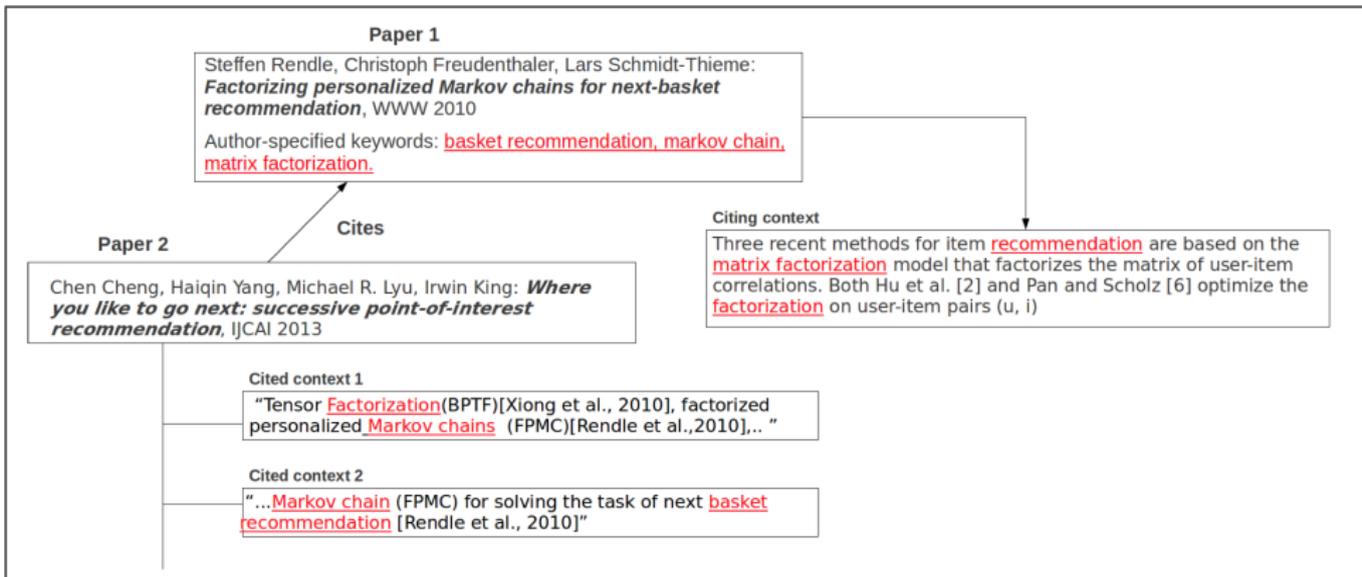
# From Data to Knowledge

- A typical scientific research paper:
  - Proposes new problems or extends the state-of-the-art for existing research problems.
  - Cites relevant, previously-published papers in appropriate *contexts*.
- The citations between research papers give rise to an interlinked document network, commonly referred to as the *citation network*.

# Citation Networks

- In a citation network, information flows from one paper to another via the citation relation [Shi et al, 2010]
- Citation contexts capture the influence of one paper on another as well as the flow of information
- Citation contexts or the short text segments surrounding a paper's mention serve as "micro summaries" of a cited paper!

# A Small Citation Network



- Citation contexts are very informative!

# Citation Contexts - Not a New Idea

- Using terms from citation contexts resembles the analysis of hyperlinks and the graph structure of the Web
  - Web search engines build on the intuition that the anchor text pointing to a page is a good descriptor of its content, and thus use anchor terms as additional index terms for a target webpage.
- Previously used for other tasks:
  - Indexing of cited papers [Ritchie et al. (2006)]
  - Author influence in document networks [Kataria et al., 2011]
  - Scientific paper summarization [Abu-Jbara and Radev, 2011; Qazvinian et al., 2010; Qazvinian and Radev, 2008; Mei and Zhai, 2008; Lehnert et al., 1990; Nakov et al., 2004]

# Citation Contexts to Keyphrase Extraction

- Citation contexts capture how one paper influences another along various aspects such as topicality, domain of study, and algorithms
- How can we use these contexts and how do they help in keyphrase extraction?
  - We propose **CiteTextRank (CTR)**: an unsupervised, graph-based algorithm that incorporates evidence from multiple sources (citation contexts as well as document content) in a flexible way to extract keyphrases [Das Gollapalli and Caragea, 2014].
  - We propose **Citation-enhanced Keyphrase Extraction (CeKE)**: a supervised binary classification model built on a combination of novel features that capture information from citation contexts and existing features from previous works [Caragea et al., 2014].

# Unsupervised Keyphrase Extraction

General steps for unsupervised keyphrase extraction algorithms:

- ① Extract candidate words or lexical units from the content of the target document by applying stopword and parts-of-speech filters.
- ② Score candidate words based on some criterion.
  - For example, in the TFIDF scoring scheme, a candidate word score is the product of its frequency in the document and its inverse document frequency in the collection.
- ③ Score consecutive words, phrases or  $n$ -grams using the sum of scores of individual words that comprise the phrase [Wan and Xiao, 2008].
- ④ Output the top-scoring phrases as the predicted keyphrases.

CiteTextRank incorporates information from *citation contexts* while scoring candidate words in step 2.

# CiteTextRank: Definitions and Notation

Let  $d$  be the target document and  $\mathcal{C}$  be a citation network such that  $d \in \mathcal{C}$ .

- Definitions:
  - A *cited context* for  $d$  is defined as a context in which  $d$  is cited by some paper  $d_i$  in the network.
  - A *citing context* for  $d$  is defined as a context in which  $d$  is citing some paper  $d_j$  in the network.
  - The content of  $d$  comprises its *global context*.
- Let  $T$  represent the types of available contexts for  $d$ 
  - The *global context* of  $d$
  - $\mathcal{N}_d^{Ctd}$  : the set of *cited* contexts for  $d$
  - $\mathcal{N}_d^{Ctg}$  : the set of *citing* contexts for  $d$
  - $\mathcal{N}_d^{Sim}$  : textually-similar global contexts

# Graph Construction in CiteTextRank

We construct an undirected graph,  $G = (V, E)$  for  $d$  as follows:

- ① For each unique candidate word from all available contexts of  $d$ , add a vertex in  $G$ .
- ② Add an undirected edge between two vertices  $v_i$  and  $v_j$  if the words corresponding to these vertices occur within a window of  $w$  contiguous tokens in any of the contexts.
- ③ The weight  $w_{ij}$  of an edge  $(v_i, v_j) \in E$  is given as

$$w_{ij} = w_{ji} = \sum_{t \in T} \sum_{c \in C_t} \lambda_t \cdot \text{cossim}(c, d) \cdot \#_c(v_i, v_j)$$

where  $\lambda_t$  is the weight for contexts of type  $t$  and  $C_t$  is the set of contexts of type  $t \in T$ .

# Parameterized Edge Weights in CiteTextRank

- Unlike simple graph edges with fixed weights, our equations correspond to *parameterized* edge weights.
- We incorporate the notion of “importance” of contexts of a certain type using the  $\lambda_t$  parameters.

Example:

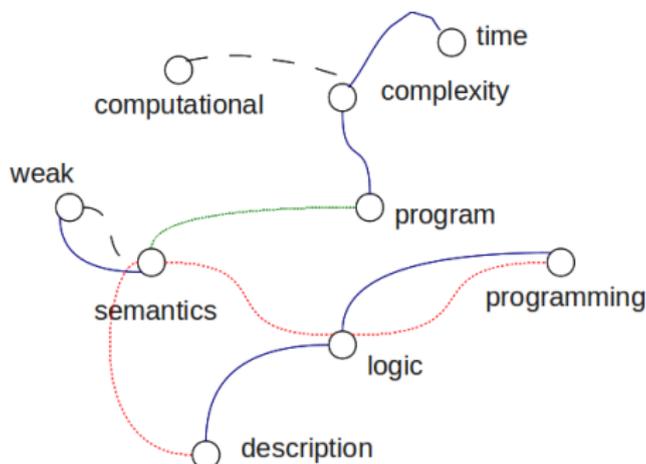


Figure: Visualization of our edges on a small word graph. Edges from different contexts are shown using different colors/line-styles.

# Vertex Scoring in CiteTextRank

We score vertices in  $G$  using their PageRank obtained by recursively computing:

$$s(v_i) = (1 - \alpha) + \alpha \sum_{v_j \in \text{Adj}(v_i)} \frac{w_{ji}}{\sum_{v_k \in \text{Adj}(v_j)} w_{jk}} s(v_j)$$

[Page et al., 1999]

- The PageRank score for a vertex provides a measure of its importance in the graph by taking into account global information computed recursively from the entire graph
- PageRank shown to be state-of-the-art in works involving word graphs for keyphrase extraction [Mihalcea and Tarau, 2004; Liu et al., 2010].

# Datasets

- We constructed three datasets of research papers and their associated citation networks using CiteSeerX.
- These datasets use:
  - The proceedings of the ACM Conference on Knowledge Discovery and Data Mining (KDD) and the World Wide Web Conference (WWW);
  - The UMD dataset from University of Maryland
- The author-input keywords were used as gold-standard for evaluation.

Conference	#Titles(Org)	#Titles(CiteSeer)	#Queries	AvgKeywords	AvgCitingContexts	AvgCitedContexts
UMD	490	439	163	3.93	20.15	34.65
WWW	2936	1350	406	4.81	15.91	17.39
KDD	1829	834	335	4.09	18.85	16.82

**Table 1:** Summary of datasets: #Queries represent the number of documents for which both citing and cited contexts were extracted from CiteSeerX and for which author-input keyphrases are available

All datasets are available upon request.

# Experiments and Results for CTR

Our experiments are organized around the following questions:

- How sensitive is CiteTextRank to its parameters?
- How well does citation network information aid in keyphrase extraction for research papers?
- How does CiteTextRank compare with state-of-the-art methods?
- We used citation contexts obtained from CiteSeerX, i.e., 50 words on each side of a citation mention

**Evaluation measures:** Precision, Recall, F1 and mean reciprocal rank, MRR

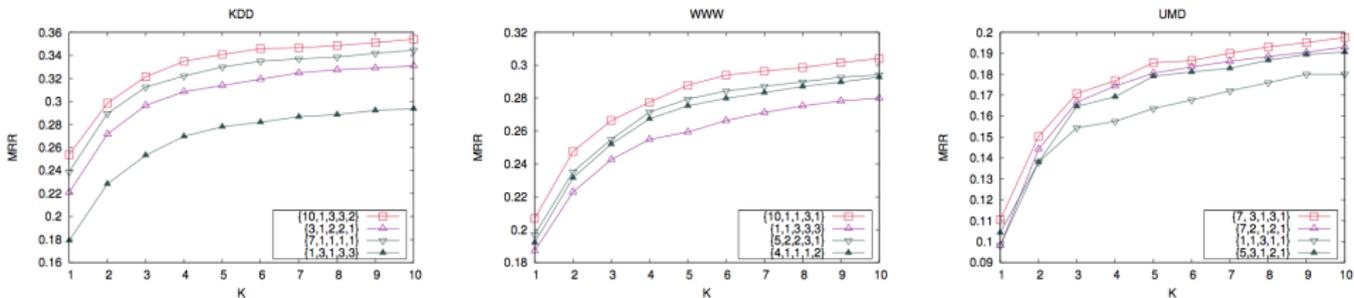
- We show results using MRR:

$$MRR = \frac{1}{|Q|} \sum_{q=1, \dots, |Q|} \frac{1}{r_q}$$

$r_q$  is the rank at which the first correct prediction was found for  $q \in Q$ .

# How Sensitive is CiteTextRank to its Parameters?

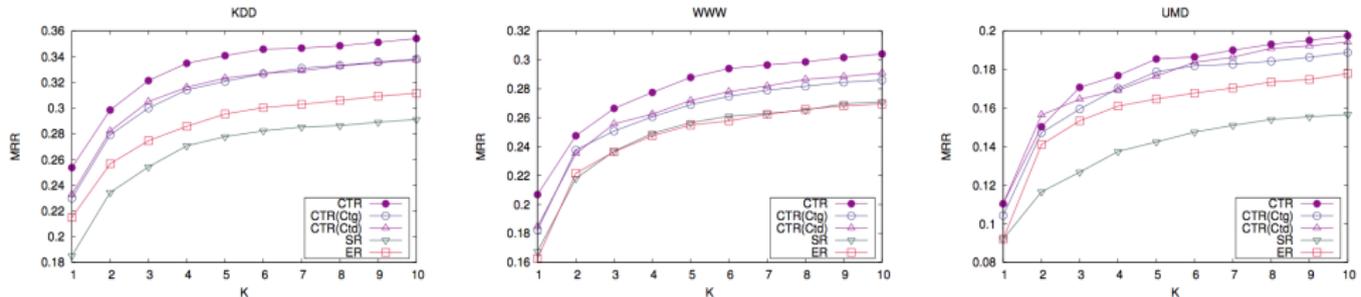
Values 1-10 were tested for each parameter in steps of 1.



**Figure:** Parameter tuning for CTR. Sample configurations are shown.  $\{a,b,c,d,e\}$  indicates that the window parameter is set to “a” with “b”, “c”, “d”, “e” as weights for textually-similar neighbors, cited, citing, and global contexts, respectively.

The varying performance of **CiteTextRank** with different  $\lambda_r$  parameters illustrates the flexibility that our model allows in treating each type of evidence differently.

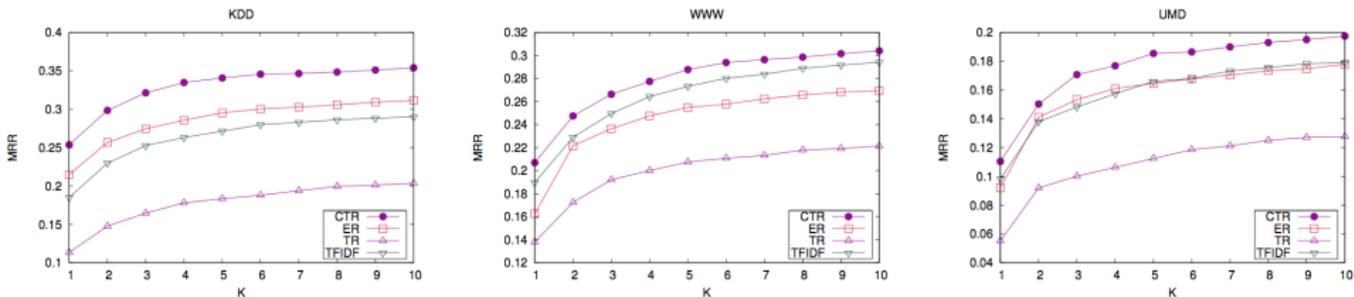
# How Well Does Citation Network Information Aid in Keyphrase Extraction for Research Papers?



**Figure:** Effect of citation network information on keyphrase extraction. CTR that uses citation network neighbors is compared with ExpandRank (ER) that uses textually-similar neighbors and SingleRank (SR) that only uses the target document content.

**CiteTextRank** substantially outperforms models that take into account only textually-similar documents. Cited and citing contexts contain significant hints that aid keyphrase extraction.

# How Does CiteTextRank Compare with Other Existing State-of-the-Art Methods?



**Figure:** MRR curves for different keyphrase extraction methods. CTR is compared with the baselines: TFIDF, TextRank (TR), and ExpandRank (ER).

**CiteTextRank** effectively out-performs the state-of-the-art baseline models for keyphrase extraction.

# Supervised Keyphrase Extraction

- We propose Citation-enhanced Keyphrase Extraction (CeKE):
  - A supervised binary classification model built on a combination of novel features that capture information from citation contexts and existing features from previous works.

# Features for CeKE

Feature Name	Description
Existing features for keyphrase extraction	
<i>tf-idf</i>	term frequency * inverse document frequency, computed from a target paper; used in KEA
<i>relativePos</i>	the position of first occurrence of a phrase divided by the total number of tokens; used in KEA and Hulth's methods
POS	the part-of-speech tag of the phrase; used in Hulth's methods
Novel features - Citation Network Based	
<i>inCited</i>	if the phrase occurs in cited contexts
<i>inCiting</i>	if the phrase occurs in citing contexts
<i>citation tf-idf</i>	the <i>tf-idf</i> value of the phrase, computed from the aggregated citation contexts
Novel features - Extensions of Existing Features	
<i>first position</i>	the distance of the first occurrence of a phrase from the beginning of a paper
<i>tf-idf-Over</i>	<i>tf-idf</i> larger than a threshold $\theta$
<i>firstPosUnder</i>	the distance of the first occurrence of a phrase from the beginning of a paper is below some value $\beta$

# Experiments and Results for CeKE

The experiments for CeKE are organized around the following questions:

- How does CeKE compare with existing supervised models that use only information intrinsic to the data?
- How is our Citation-Enhanced algorithm comparing with recent unsupervised models?
- How well does our proposed model perform in the absence of either cited or citing contexts?

Evaluation measures:

- Precision, Recall, and F1-score.

# How Does CeKE Compare with Supervised Models?

Method	WWW			KDD		
	Precision	Recall	F1-score	Precision	Recall	F1-score
<b>Citation - Enhanced (CeKE)</b>	<b>0.227</b>	<b>0.386</b>	<b>0.284</b>	<b>0.213</b>	<b>0.413</b>	<b>0.280</b>
Hulth - $n$ -gram with tags	0.165	0.107	0.129	0.206	0.151	0.172
KEA	0.210	0.146	0.168	0.178	0.124	0.145

**Table:** Comparison of CeKE with Hulth's and KEA methods.

Features used in previous supervised methods:

- Hulth's features: *POS*, *relative position*, *document frequency* and *collection frequency*.
- KEA's features: *tf-idf* and *relative position*

# How Does CeKE Compare with Unsupervised Models?

Method	WWW			KDD		
	Precision	Recall	F1-score	Precision	Recall	F1-score
<b>Citation - Enhanced (CeKE)</b>	<b>0.227</b>	<b>0.386</b>	<b>0.284</b>	<b>0.213</b>	<b>0.413</b>	<b>0.280</b>
TF-IDF - Top 5	0.089	0.100	0.094	0.083	0.102	0.092
TF-IDF - Top 10	0.075	0.169	0.104	0.080	0.203	0.115
TextRank - Top 5	0.058	0.071	0.062	0.051	0.065	0.056
TextRank - Top 10	0.062	0.133	0.081	0.053	0.127	0.072
ExpandRank - 1 neigh. - Top 5	0.088	0.109	0.095	0.077	0.103	0.086
ExpandRank - 1 neigh. - Top 10	0.078	0.165	0.101	0.071	0.177	0.098
ExpandRank - 5 neigh. - Top 5	0.093	0.113	0.100	0.080	0.108	0.090
ExpandRank - 5 neigh. - Top 10	0.080	0.172	0.104	0.068	0.172	0.095
ExpandRank - 10 neigh. - Top 5	0.094	0.113	0.100	0.077	0.103	0.086
ExpandRank - 10 neigh. - Top 10	0.076	0.162	0.099	0.065	0.164	0.091

**Table:** Comparison of CeKE with state-of-the-art unsupervised systems.

- *TextRank*: window size is set to 2.
- *ExpandRank*: window size is set to 10.

# How Does CeKE Perform in the Absence of Either Cited or Citing Contexts?

Method	WWW			KDD		
	Precision	Recall	F1-score	Precision	Recall	F1-score
CeKE - Both contexts	<b>0.227</b>	<b>0.386</b>	<b>0.284</b>	<b>0.213</b>	<b>0.413</b>	<b>0.280</b>
CeKE - Only cited contexts	0.222	0.286	0.247	0.192	0.300	0.233
CeKE - Only citing contexts	0.203	0.342	0.253	0.195	0.351	0.250

**Table:** Results with both contexts and only cited/citing contexts.

# Anecdotal Evidence

- We considered an EMNLP paper by Poon and Domingos [2009].
  - Our classifier trained on both WWW and KDD
  - We gathered from the Web 49 cited contexts and 30 citing contexts
  - The classifier was tuned to return only high-confidence keyphrases

## Unsupervised Semantic Parsing<sup>0.997</sup>

We present the first unsupervised approach to the problem of learning a **semantic parser**<sup>1.000</sup>, using **Markov logic**<sup>0.991</sup>. Our **USP system**<sup>0.985</sup> transforms dependency trees into quasi-logical forms, recursively induces lambda forms from these, and clusters them to abstract away syntactic variations of the same meaning. The MAP **semantic parse**<sup>1.000</sup> of a sentence is obtained by recursively assigning its parts to lambda-form clusters and composing them. We evaluate our approach by using it to extract a knowledge base from biomedical abstracts and answer questions. **USP**<sup>1.000</sup> substantially outperforms TextRunner, DIRT and an informed baseline on both precision and recall on this task.

Human annotated keyphrases: *unsupervised semantic parsing, Markov logic, USP system*

*Grey* - filtered out words; *Black* - candidate phrases; **Bold red** - predicted keyphrases; *Numbers* - classifier's confidence.

# Conclusions and Future Directions

- We proposed supervised and unsupervised models for keyphrase extraction using multiple sources of evidence
  - The textual content of a document and its citing and cited contexts in the interlinked document network
- Our models give significant improvements over baseline models for multiple datasets of research papers in the Computer Science domain
- Future directions:
  - **Citation context lengths**: Interesting to incorporate more sophisticated approaches to identifying the text that is relevant to a target citation [Abu-Jbara and Radev, 2012; Teufel, 1999] and study the influence of context lengths on the quality of extracted keyphrase.
  - **Integrated terms not found in the target paper** to be predicted as keyphrases
  - Further evaluation of CTR on other domains, e.g., the ACL Anthology and PubMed.
  - Extend CTR for extracting document summaries similar to [Mihalcea and Tarau 2004; Qazvinian, Radev, and Özgür, 2010]

# References

- S. Das Gollapalli and C. Caragea (2014). Extracting Keyphrases from Research Papers using Citation Networks. In: *Proceedings of the 28th American Association for Artificial Intelligence (AAAI '14)*.
- C. Caragea, F. Bulgarov, A. Godea, and S. Das Gollapalli. (2014). Citation-Enhanced Keyphrase Extraction from Research Papers: A Supervised Approach. In: *Proceedings of the Conference on Empirical Methods in Natural Language Processing (EMNLP '14)*.
- R. Mihalcea and P. Tarau. (2004). TextRank: Bringing order into text. In: *Proceedings of the Conference on Empirical Methods in Natural Language Processing (EMNLP '04)*.
- X. Wan and J. Xiao (2008). Single document keyphrase extraction using neighborhood knowledge. In: *Proceedings of the Association for the Advancement of Artificial Intelligence (AAAI '08)*.
- Z. Liu, W. Huang, Y. Zheng, and M. Sun. (2010). Automatic keyphrase extraction via topic decomposition. In: *Proceedings of the Conference on Empirical Methods in Natural Language Processing (EMNLP '10)*.
- V. Qazvinian, D. Radev, and A. Özgür. 2010. Citation summarization through keyphrase extraction. In: *Proceedings of the 23rd International Conference on Computational Linguistics, COLING '10, pages 895 - 903*.

# Thank you!



Florin Bulgarov



Sujatha Das Gollapalli



Kishore Neppalli



Andrea Godea