



The  
University  
Of  
Sheffield.

Semantics 2018 - 12 September 2018

# Adapted TextRank for Term Extraction: A Generic Method of Improving Automatic Term Extraction Algorithms

Ziqi Zhang<sup>1</sup>, **Johann Petrak**<sup>2</sup>, Diana Maynard<sup>2</sup>  
[ziqi.zhang@sheffield.ac.uk](mailto:ziqi.zhang@sheffield.ac.uk), [johann.petrak@sheffield.ac.uk](mailto:johann.petrak@sheffield.ac.uk),  
[d.maynard@sheffield.ac.uk](mailto:d.maynard@sheffield.ac.uk)

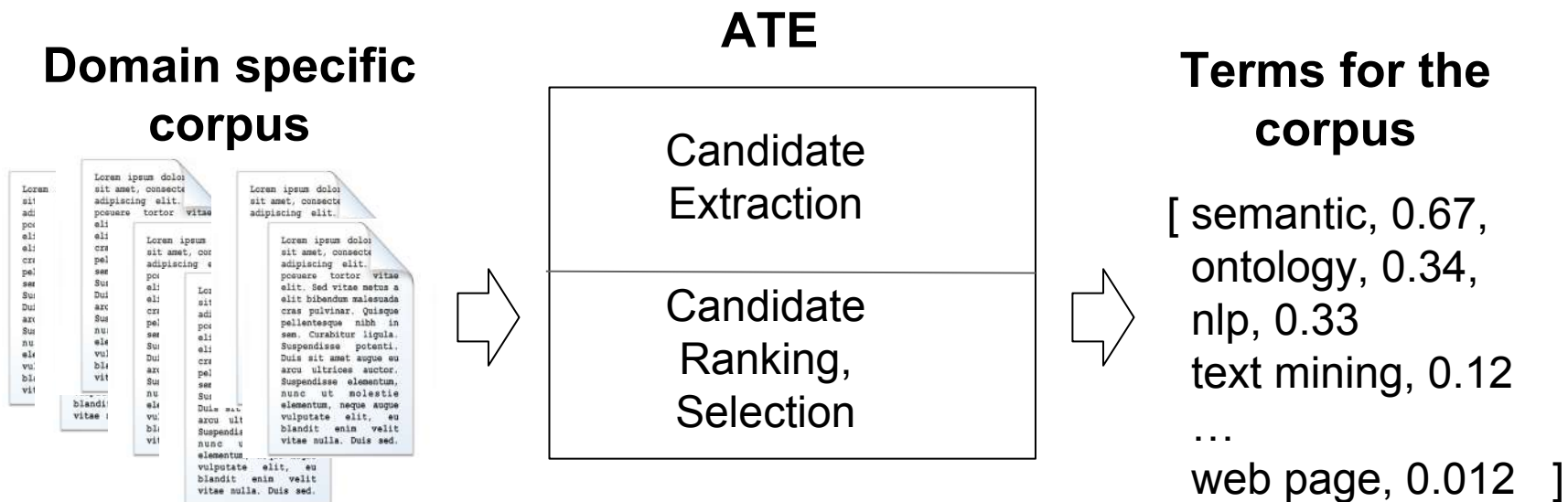
1. Information School, The University of Sheffield, UK

2. **Department of Computer Science, The University of Sheffield, UK**



# The Task of ATE

- **Input:** (reasonably large) domain specific, focused corpus
- **Output:** list terms from the corpus, representing the domain
- **Approach**
  - Candidate extraction: domain-dependent, usually noun phrases, n-grams, or sequence matched by PoS patterns
  - Candidate ranking & selection: scoring candidates based on corpus statistics, selection by threshold, or machine learning





# The Task of ATE

- **A classic text mining problem**
  - Dating back to 1990s (Bourigault 1992)
  - To date still an active area of research
- **A fundamental step to many complex tasks**
  - Ontology engineering
  - Dictionary, terminology construction
  - Information Retrieval
  - Translation
  - ...
- **Context of this work: KNOWMAK (<https://www.knowmak.eu/>)**



# The Task of ATE

## Differentiation from related tasks

		ATE
<b>Keyword Extraction</b>	<ul style="list-style-type: none"> <li>- document specific</li> <li>- only a handful</li> <li>- mainly for indexing</li> </ul>	<ul style="list-style-type: none"> <li>- domain specific</li> <li>- # depends on corpus</li> <li>- mainly knowledge acquisition</li> </ul>
<b>NER</b>	<ul style="list-style-type: none"> <li>- usually real world named entities</li> <li>- sentence context is more important</li> <li>- semantic typing</li> </ul>	<ul style="list-style-type: none"> <li>- domain specific terms</li> <li>- corpus level statistics are more important</li> <li>- no typing</li> </ul>

descriptions of electronic products from several shops that offer Microdata markup. We present for each step of the data integration process information extraction, product classification, product extraction, identity resolution, and data fusion. We evaluate our processing pipeline using 1.9 million products from 9240 e-shops which we extracted from the Crawl 2012, a large public Web corpus.

### Categories and Subject Descriptors

H.3.1 [Information Storage and Retrieval]: Analysis and Indexing; H.3.3 [Information Storage and Retrieval]: Information Search and Retrieval; H.3.4 [Information Storage and Retrieval]: Information Filtering

### Keywords

Microdata; Information Extraction; Data Integration; Identity Resolution

### 1. INTRODUCTION



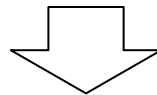


- **ATE still an unsolved problem**
  - No 'all-rounder' method
  - Performance always depends on data and domain
  - 'one-size-fits-all' solution feasible?
- **ATE methods are predominantly unsupervised**
  - For many domains there are already domain specific resources potentially useful, e.g., unlabelled corpus, pre-compiled named entity lists, partial ontologies, etc
  - Can we benefit from those?



# Motivation and Contribution

- **ATE still an unsolved problem**
  - No 'all-rounder' method
  - Performance always depends on data and domain
  - 'one-size-fits-all' solution feasible?
- **ATE methods are predominantly unsupervised**
  - For many domains there are already domain specific resources potentially useful, e.g., unlabelled corpus, pre-compiled named entity lists, partial ontologies, etc
  - Can we benefit from those?



A generic method that employs semantic relatedness to a set of **domain specific seed words** to potentially **improve any ATE** algorithms (by up to 25 percentage points in average precision in experiments).



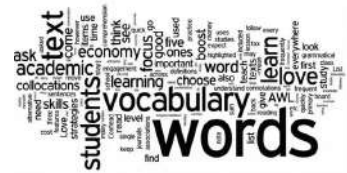
# AdaText - Overview

## Adapted TextRank for Automatic Term Extraction

Domain specific seed words/phrases

oolong  
cream  
herbal  
pu-erh  
ginger  
green  
vanilla  
dessert  
earl grey  
espresso  
chamomile  
americano  
cappuccino  
macchiato  
moccacino  
rooibus

Semantic relatedness



Filter by threshold



TextRank

Extract words

ATE (any algorithm)

[  $t_1=1.99$ ,  
 $t_2=1.21$ ,  
 $t_3=1.10$ ,  
... ]

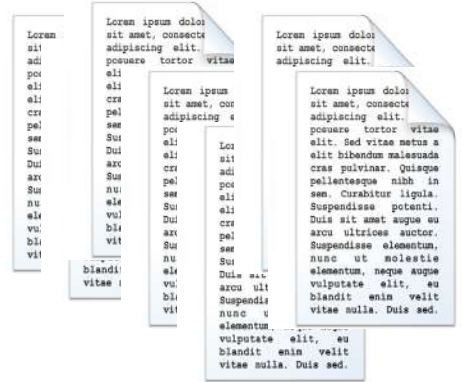
+

[  $w_1=0.67$ ,  
 $w_2=0.34$ ,  
 $w_3=0.22$ ,  
... ]

Re-rank

[  $t_1=2.19$ ,  
 $t_3=1.41$ ,  
 $t_2=1.29$ ,  
... ]

Domain specific corpus









# AdaText - Seeding

- **Input**

- $C$  - the target corpus from which terms are extracted
- $S$  - a set of 'seed' word/phrases representing the domain
  - taken from existing domain lexicons, or generated in an unsupervised way from available corpora
  - May not contain real terms from  $C$

- **Process**

- Extract words from  $C$ , as  $W$
- Compute pairwise semantic relatedness for  $S \times W$ 
  - Cosine similarity using GloVe embedding vectors
  - OOV ignored, phrase based on compositional averaging (Iyyer et al. 2015)

- **Output**

- $W_{sub}$  a subset of  $W$ , satisfying relatedness  $> \mathit{min}$   
Intuitively, they are more 'relevant' to the domain



# AdaText - Corpus Level TextRank

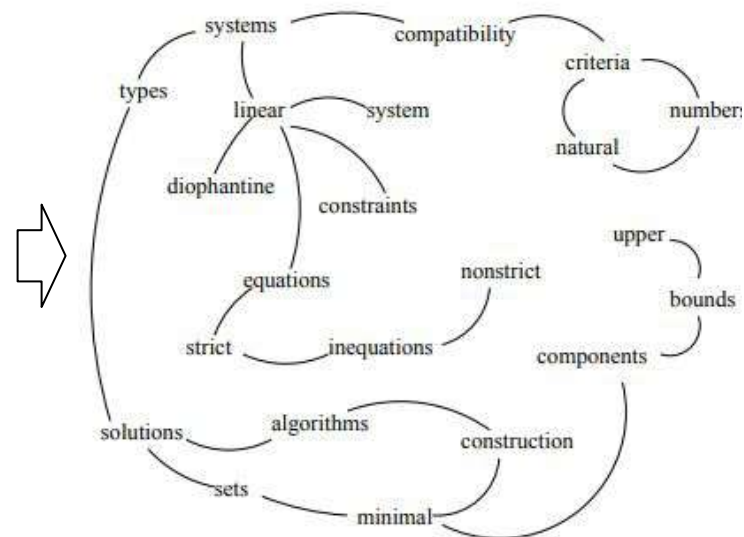
- **Input**

- $C$  - the target corpus from which terms are extracted
- $W_{sub}$  - the subset of words selected before

- **Process**

- Apply TextRank to the graph created for  $W_{sub}$  to compute a TextRank ( $tr$ ) score of every word  $w$  in  $W_{sub}$
- Traditional TextRank (Mihalcea et al., 2004) is a PageRank process to a graph of words from each document, where an edge is created if words co-occur in a context window of ***win***

Compatibility of systems of linear constraints over the set of natural numbers. Criteria of compatibility of a system of linear Diophantine equations, strict inequations, and nonstrict inequations are considered. Upper bounds for components of a minimal set of solutions and algorithms of construction of minimal generating sets of solutions for all types of systems are given. These criteria and the corresponding algorithms for constructing a minimal supporting set of solutions can be used in solving all the considered types systems and systems of mixed types





- **Input**

- $C$  - the target corpus from which terms are extracted
- $W_{sub}$  - the subset of words selected before

- **Process**

- Apply TextRank to the graph created for  $W_{sub}$  to compute a TextRank ( $tr$ ) score of every word  $w$  in  $W_{sub}$
- Here it is adapted in two ways
  - A graph of words from the entire corpus
  - An edge is created if two words appear within ***win*** *anywhere in the corpus* (in any document)

- **Output**

- $tr$  scores for every word  $w$  in  $W_{sub}$



# AdaText - Combining with ATE

- **Input**

- $C$  - the target corpus from which terms are extracted
- $ATE$  - some ATE algorithm
- $tr$  scores for every word  $w$  in  $W_{sub}$

- **Process**

- Apply  $ATE$  to  $C$  to extract and score candidate terms
- Revise each candidate term's score using  $tr$  scores for its composing words

$$score(t_i) = (1.0 + \frac{\sum_{w_i \in words(t_i)} tr(w_i)}{|words(t_i)|}) \times ate(t_i)$$

- Then re-rank candidate terms by the new score

- **Output**

- Re-ranked list of candidate terms



# Experiment and Findings

- **Base ATE methods** (as AdaText needs ATE scores of candidate terms)
  - Modified TFIDF (Zhang et al., 2016)
  - CValue (Ananiadou 1994)
  - Basic (Bordea et al., 2013)
  - RAKE (Rose et al., 2010)
  - Weirdness (Ahmad et al., 1999)
  - LinkProbability (LP, Astrakhantsev, 2016)
  - $X^2$  (Matsuo et al., 2003)
  - GlossEx (Park et al., 2002)
  - Positive Unlabelled (PU) learning (Astrakhantsev, 2016)
  - AvgRel - average relatedness score with seeds
- Use implementations:
  - JATE (<https://github.com/ziqizhang/jate>)
  - ATR4S (<https://github.com/ispras/atr4s>)



## Evaluation measures

- Precision for top  $K$  ranked candidate terms
- $K = \{50, 100, 500, 1000, 2000\}$
- **Average P@K for all five  $K$ 's**



## Datasets

- GENIA
  - 2,000 semantically annotated Medline abstracts
  - 434k words
  - 33k target terms
- ACLv2
  - 300 ACL paper abstracts
  - 32k words
  - 3k target terms



## Seeds and parameters

- For GENIA:
  - 5,502 named entities from the BioNLP Shared Task 2011, only **25** match candidate terms
- For ACLv2:
  - 1,301 noun phrases from the titles of ACL, NAACL, and EACL papers (since 2000), **none** matches candidate terms
- Semantic relatedness threshold  $min=0.5$  to  $0.85$  with  $0.05$  increment (selects for GENIA/ACL ~ 50/70 % ... 10/5 %)
- TextRank context window  $win=5, 10$

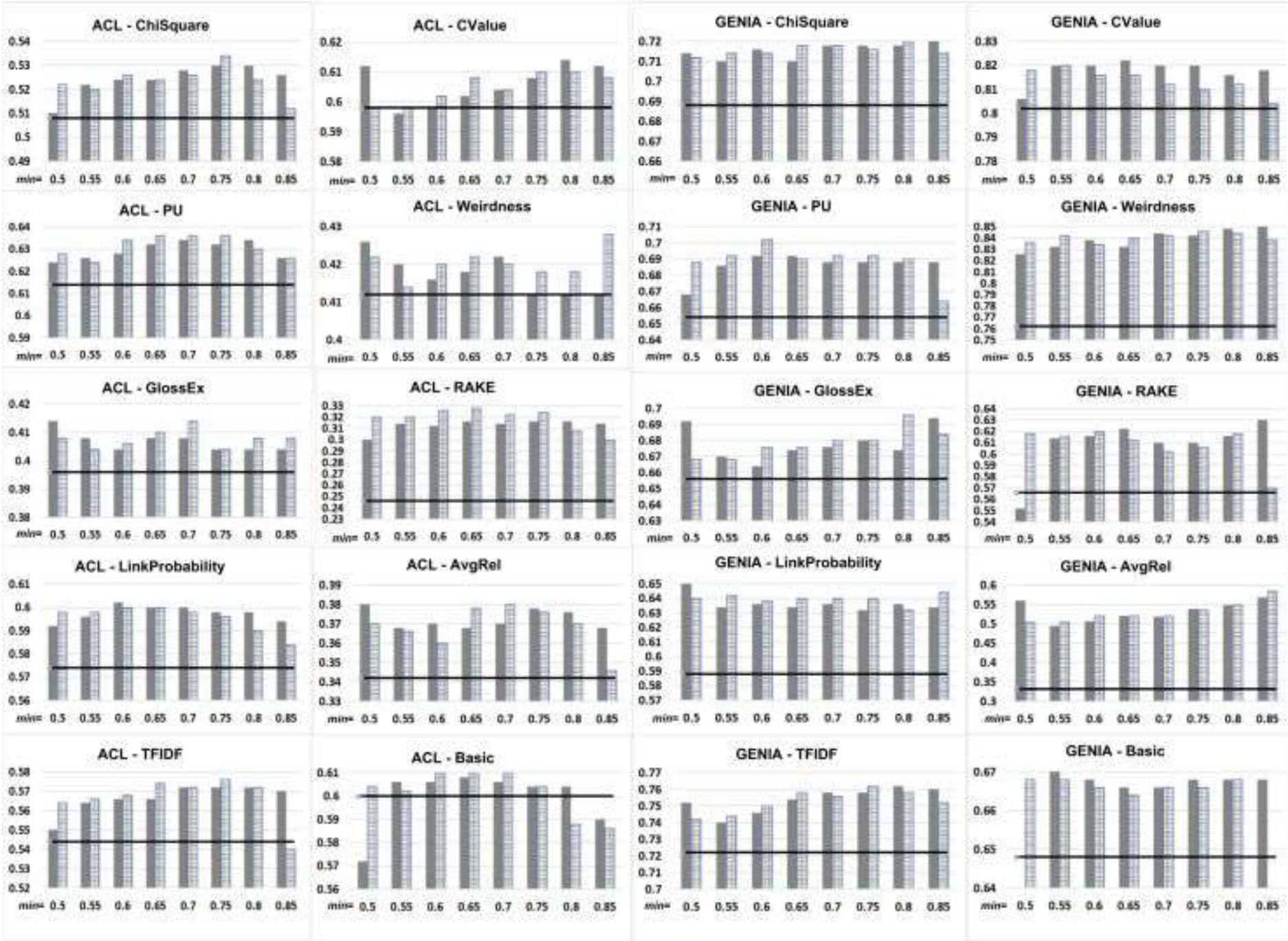




# Result - Base ATE

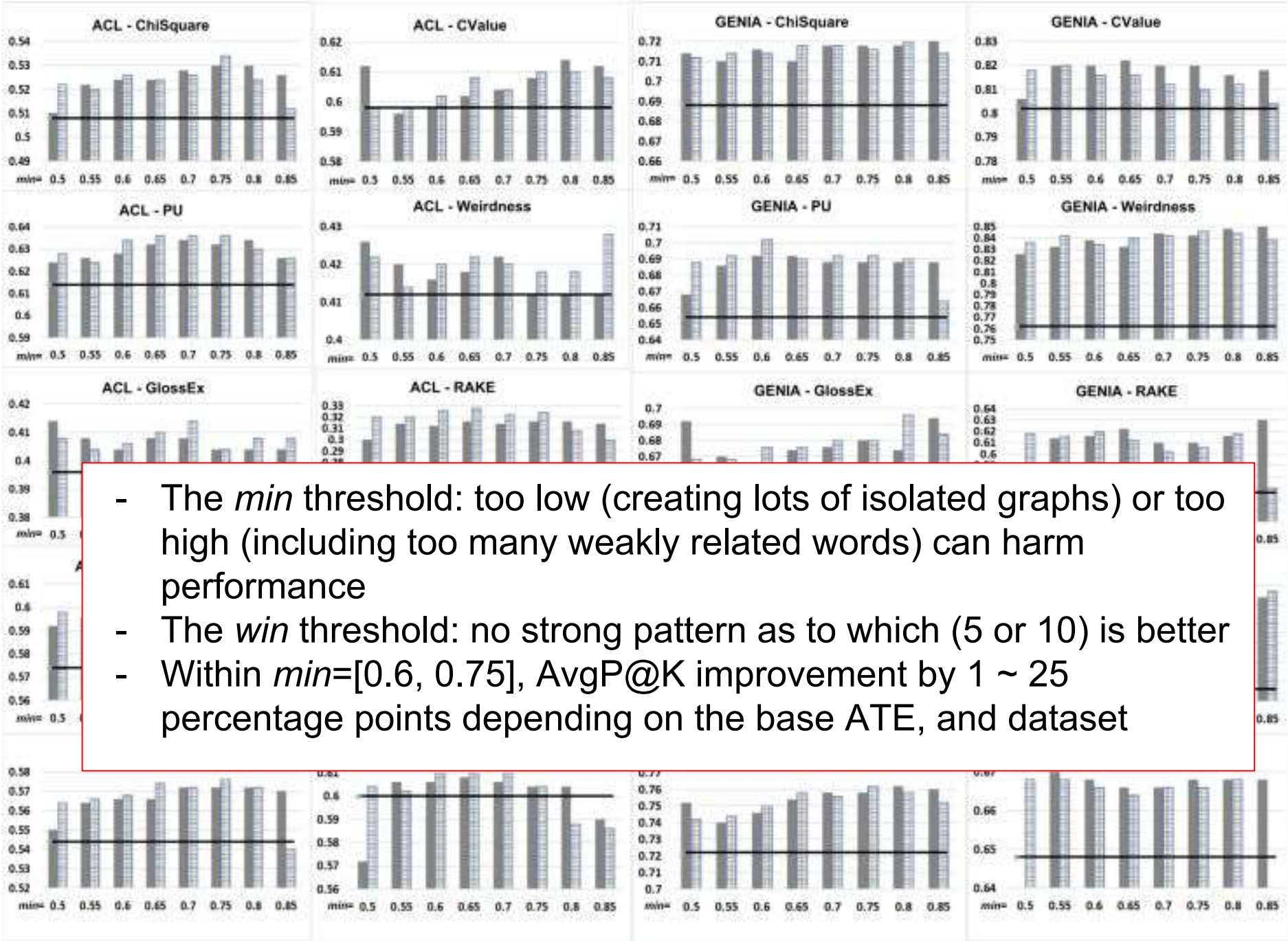
	Basic	LP	PU	CValue	GlossEx	RAKE	TFIDF	Weirdness	$\chi^2$	AvgRel
ACLv2										
P@50	<b>.84</b>	.72	.82	.62	.44	.18	.64	.40	.58	.30
P@100	.72	.69	<b>.82</b>	.69	.46	.15	.65	.50	.62	.34
P@500	.56	.56	.60	<b>.67</b>	.34	.29	.53	.36	.48	.39
P@1,000	.49	.51	.43	<b>.56</b>	.36	.29	.47	.40	.45	.35
P@2,000	.39	.39	.40	<b>.45</b>	.38	.32	.43	.40	.41	.33
AvgP@K	.60	.57	<b>.61</b>	.60	.40	.25	.54	.41	.52	.34
GENIA										
P@50	.80	.38	.74	.86	<b>.88</b>	.68	.68	.78	.66	.16
P@100	.74	.51	.69	<b>.83</b>	.82	.63	.65	.74	.69	.28
P@500	.64	.70	.65	<b>.80</b>	.58	.56	.74	.78	.71	.36
P@1,000	.57	.69	.61	<b>.78</b>	.53	.52	.77	.77	.71	.44
P@2,000	.49	.66	.58	.74	.47	.44	<b>.77</b>	.74	.67	.42
AvgP@K	.65	.59	.65	<b>.80</b>	.66	.57	.72	.76	.69	.33

- Base ATE performance varies significantly depending on datasets.
- No single, consistently winning method on all five  $K$ 's.
- E.g., PU is the best performing in **AvgP@K** on the ACL corpus, but the fourth worst performing on the GENIA corpus.



■ win=5    ▨ win=10





- The *min* threshold: too low (creating lots of isolated graphs) or too high (including too many weakly related words) can harm performance
- The *win* threshold: no strong pattern as to which (5 or 10) is better
- Within *min*=[0.6, 0.75], AvgP@K improvement by 1 ~ 25 percentage points depending on the base ATE, and dataset



- **The takeaway message**

- There is probably never a 'one-size-fit-all' ATE method, instead, think about improving existing ones
- AdaText makes use of existing domain resources and builds on the TextRank algorithm
- Generic method able to improve, potentially, any ATE method

- **Future work**

- Whether and how the size and source of the seed lexicon affects performance
- Adapt TextRank to a graph of both words and phrases, and see how this affects results



- **Data**

- Genia corpus, ACL corpus available
- Glove embeddings available

- **Software**

- JATE (<https://github.com/ziqizhang/jate>)
- ATR4S (<https://github.com/ispras/atr4s>)
- Code for this work: <https://github.com/ziqizhang/texpr>

- **Slides**

- <https://goo.gl/1sPuhg>



1. Bourigault, D. 1992. Surface grammatical analysis for the extraction of terminological noun phrases. In 14th International Conference on Computational Linguistics - COLING 92, 977–98
2. Iyyer, M., Manjunatha, V., Boyd-Graber, J., Daume, H., 2015. Deep unordered composition rivals syntactic methods for text classification, in: Association for Computational Linguistics. URL: docs/2015\_acl\_dan.pdf.
3. Mihalcea, R., Tarau, P., 2004. TextRank: Bringing order into texts, in: Proc. of EMNLP'04.
4. Zhang, Z., Gao, J., Ciravegna, F., 2016. Jate 2.0: Java automatic term extraction with apache solr, in: Proc. of LREC'16
5. Ananiadou, S., 1994. A methodology for automatic term recognition, in: Proc. of COLING1994, ACL, Stroudsburg, PA, USA. pp. 1034–1038.
6. Bordea, G., Buitelaar, P., Polajnar, T., 2013. Domain-independent term extraction through domain modelling, in: Proc. of the Conference on Terminology and Artificial Intelligence.
7. Astrakhantsev, N., 2015. Methods and software for terminology extraction from domainspecific text collection, in: Ph.D. thesis. Institute for System Programming of Russian Academy of Sciences.
8. Rose, S., Engel, D., Cramer, N., Cowley, W., 2010. Automatic keyword extraction from individual documents. John Wiley and Sons.
9. Ahmad, K., Gillam, L., Tostevin, L., 1999. University of surrey participation in trec 8: Weirdness indexing for logical document extrapolation and retrieval (wilder), in: Proc. of TREC1999.
10. Astrakhantsev, N., 2016. Atr4s: Toolkit with state-of-the-art automatic terms recognition methods in scala. arXiv preprint arXiv:1611.07804.
11. Matsuo, Y., Ishizuka, M., 2003. Keyword extraction from a single document using word co-occurrence statistical information. International Journal on Artificial Intelligence Tools 13, 157–169.
12. Park, Y., Byrd, R., Boguraev, B., 2002. Automatic glossary extraction: Beyond terminology identification, in: Proc. of COLING'02, Association for Computational Linguistics. pp. 1–7.



The  
University  
Of  
Sheffield.

# Acknowledgements

This work is supported by the European Union's Horizon 2020 research and innovation programme under grant agreement No. 726992 (KNOWMAK project)

<https://www.knowmak.eu/>





The  
University  
Of  
Sheffield.

Thank you!

Questions?

