# Topic extraction, expert finding and trend analysis from scientific literature

*John P. McCrae*

**Insight Centre for Data Analytics**

**National University of Ireland Galway**

# Knowledge Extraction from Text

## - with Saffron -

… act or process of **retrieving awareness** or understanding **of someone** or **something**, such as facts, **information**, **descriptions**, or **skills** out of text for further data processing … usually followed by data transformation and possibly the **addition of metadata** prior to export to another stage in the data workflow …

# Original Use Case: Expert Finding

ACL Anthology
A Digital Archive of Research Papers in Computational Linguistics

## NLP

*ACL HLT, COLING, EACL, ANLP, ACL Meetings*

Saffron provides insights in a research community or organization by analyzing its main topics of investigation and the experts associated with these topics.

Saffron analysis is fully automatic and is based on text mining and linked data principles.
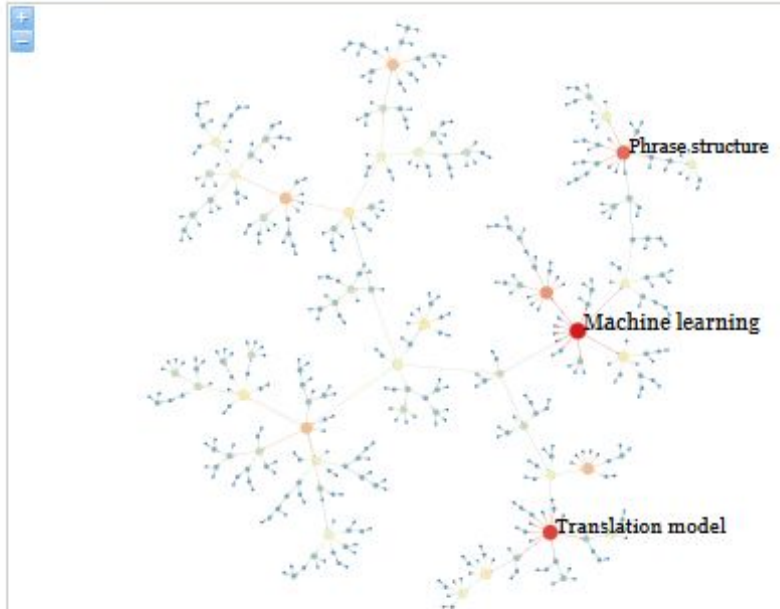
This instance of Saffron analyzes the research community in Natural Language Processing based on the proceedings of the conferences organized by the Association for Computational Linguistics (ACL).

### Hot Topics                                    more >>

1. Natural Language Processing
2. Natural Language
3. Language model
4. Statistical machine translation
5. Training data
6. Machine translation system
7. Machine translation
8. Hidden Markov Models
9. Support vector machines
10. Information retrieval
11. N-gram language model
12. Machine learning
13. Word sense disambiguation
14. Target language
15. Computer Science
16. Knowledge base
17. Human Language Technology
18. Translation model
19. Predicate-argument structure
20. Speech recognition
21. Natural language understanding
22. Feature structures
23. Spoken language systems
24. Language Processing
25. Syntactic structure
26. Natural language interface
27. Automatic evaluation of machine tr...
28. Natural language processing system...
29. Spoken dialogue systems
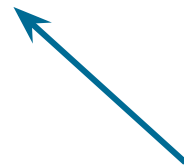30. Log-linear model

### Taxonomy

# Architecture
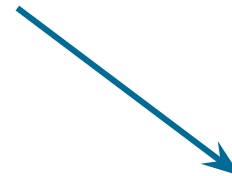
# Step 1 - Corpus Indexing

# Step 2 - Domain Modelling

…**concepts** such as *Machine Translation*…

Term

Trigger Words

…*Noun phrases* and other **elements**…

Term

# Step 3 - Topic (term) Extraction

NNS     JJ  IN  NNP    NNP

concepts such as Machine Translation

| Candidate | Weirdness | Relevance | Domain Pertinence | ... |
|---|---|---|---|---|
| *Concepts* | 0.1 | 0.6 | 0.8 | ... |
| *Machine Translation* | 0.8 | 0.7 | 0.7 | ... |

Candidate selection **NEW** by voting

# Term Extraction – ACL Anthology

| | | | | |
|---|---|---|---|---|
| 1 | Natural Language Processing | | 51 | Word alignment |
| 2 | Natural Language | | 52 | Human Language Technology Conferen... |
| 3 | Language model | | 53 | Search engine |
| 4 | Statistical machine translation | | 54 | Natural language system |
| 5 | Training data | | 55 | Spoken language |
| 6 | Machine translation system | | 56 | Dependency structure |
| 7 | Machine translation | | 57 | Latent semantic analysis |
| 8 | Hidden Markov Models | | 58 | Natural language understanding sys... |
| 9 | Support vector machines | | 59 | Noun phrases |
| 10 | Information retrieval | | 60 | Dialogue systems |
| 11 | N-gram language model | | 61 | Parsing algorithm |
| 12 | Machine learning | | 62 | Content words |
| 13 | Word sense disambiguation | | 63 | Mutual information |
| 14 | Target language | | 64 | Discourse structure |
| 15 | Computer Science | | 65 | Machine learning techniques |
| 16 | Knowledge base | | 66 | Natural language text |
| 17 | Human Language Technology | | 67 | Natural Language Generation |
| 18 | Translation model | | 68 | Knowledge sources |
| 19 | Predicate-argument structure | | 69 | Vector space model |
| 20 | Speech recognition | | 70 | Semantic classes |
| 21 | Natural language understanding | | 71 | Dynamic programming |
| 22 | Feature structures | | 72 | Topic models |
| 23 | Spoken language systems | | 73 | Morphological analysis |
| 24 | Language Processing | | 74 | Data structure |
| 25 | Syntactic structure | | 75 | Learning algorithm |

# Step 4 - Author Consolidation
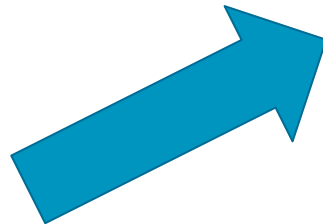
John McCrae

John P. McCrae

McCrae, J.P.

```
{
    "honorific": null          NEW
    "givenName" : "John",
    "middleInitial": "P",
    "familyName": "McCrae"
}
```

# Step 5 - DBpedia Lookup

"Machine Translation"



http://dbpedia.org/resource/Machine_translation
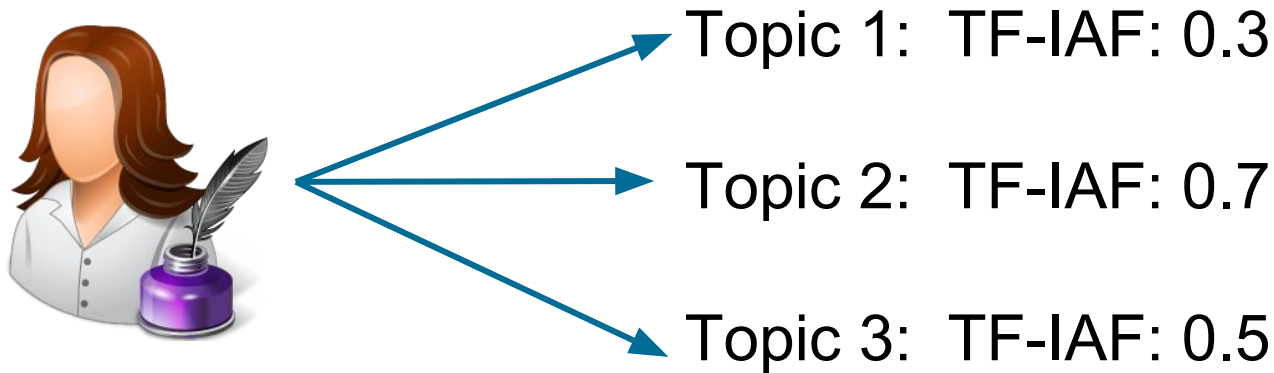
# Step 6 - Topic Statistics

Topic Generality

$$g(t) = \sum_{d \in \text{corpus}} \frac{PMI(t; d)}{p(t, d)}$$

Weaknesses:

- Favours common terms
- Denormalized PMI?

⇒ Multi-factor metric **NEW**

# Step 7 - Connect Authors

Topic 1:  TF-IAF: 0.3

Topic 2:  TF-IAF: 0.7

Topic 3:  TF-IAF: 0.5

$$\text{TF-IAF}(T; r) = \sum_{Doc \text{ if } T \in \text{Doc}, r \in \text{Authors}(\text{Doc})} \text{TF-IDF}(T)$$
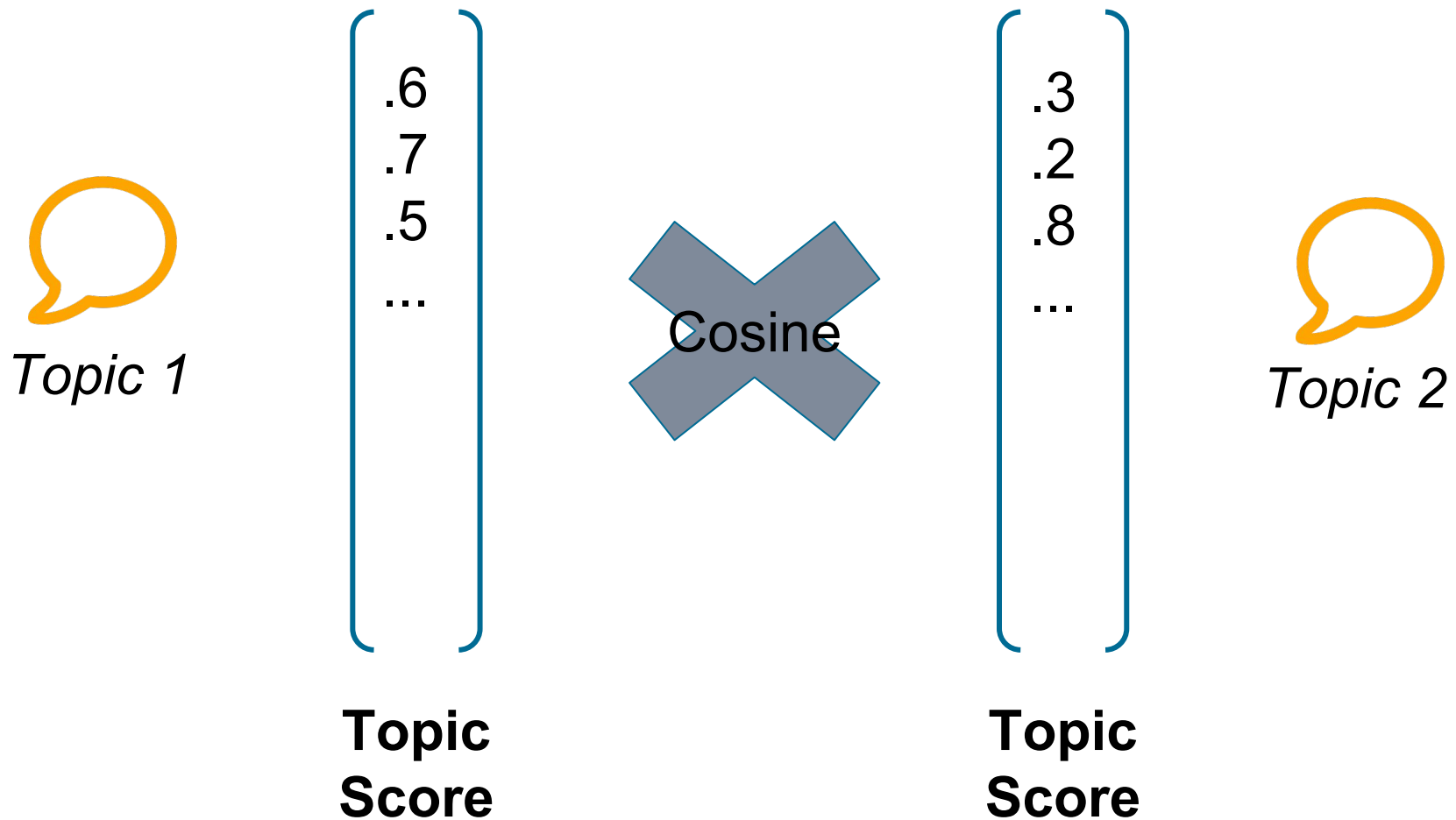
# Step 8 - Author Similarity



.6
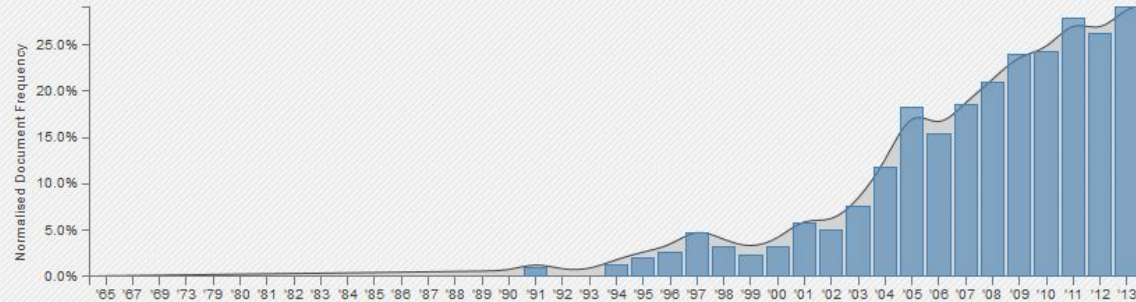.7
.5
…

Cosine

.3
.2
.8
…

**TF-IRF**

**TF-IRF**

# Step 9 - Topic Similarity

# Expertise Mining



**Statistical machine translation**

Source: http://dbpedia.org/resource/Statistical_machine_translation
See also: Statistical translation

**Experts**                                               more >>

1  Hermann Ney          +        6  Haizhou Li             +
2  Qun Liu              +        7  Kevin Knight           +
3  Kai-min K. Chang     +        8  Eiichiro Sumita        +
4  Ming Zhou            +        9  Tek Yong Lim           +
5  Stephan Vogel        +        10 Chris Callison-Burch   +

Georgeta Bordea (2013) Domain adaptive extraction of topical hierarchies for Expertise Mining. PhD Thesis, National University of Ireland, Galway

# Expertise Mining

## Qun Liu

### Topics

| | | | |
|---|---|---|---|
| 1 | Statistical machine translation + | 6 | Log-linear model + |
| 2 | Word alignment + | 7 | Translation quality + |
| 3 | BLEU score + | 8 | Translation rules + |
| 4 | Chinese word segmentation + | 9 | Dependency structure + |
| 5 | Bilingual phrases + | 10 | Phrase pairs + |

### Similar Researchers

| | | | |
|---|---|---|---|
| 1 | Shouxun Lin + | 6 | Jiajun Zhang + |
| 2 | Patrik Lambert + | 7 | Philipp Koehn + |
| 3 | Jinsong Su + | 8 | Masao Utiyama + |
| 4 | Ying Zhao + | 9 | Haitao Mi + |
| 5 | Dongdong Zhang + | 10 | Sankaranarayanan Ananthakrishnan + |

### Publications (39)

1. Improving Statistical Machine Translation using Lexicalized Rule Selection
   2008 - Zhongjun He, Qun Liu, Shouxun Lin

2. Word Lattice Reranking for Chinese Word Segmentation and Part-of-Speech Tagging
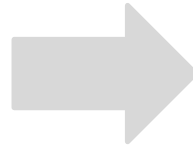   2008 - Wei-Bin Liang, Haitao Mi, Qun Liu

# Step 10 - Taxonomy Construction

- Reduce topic-topic graph to directed acyclic graph
  - Simpler hierarchical structure for corpus
- Minimum spanning tree
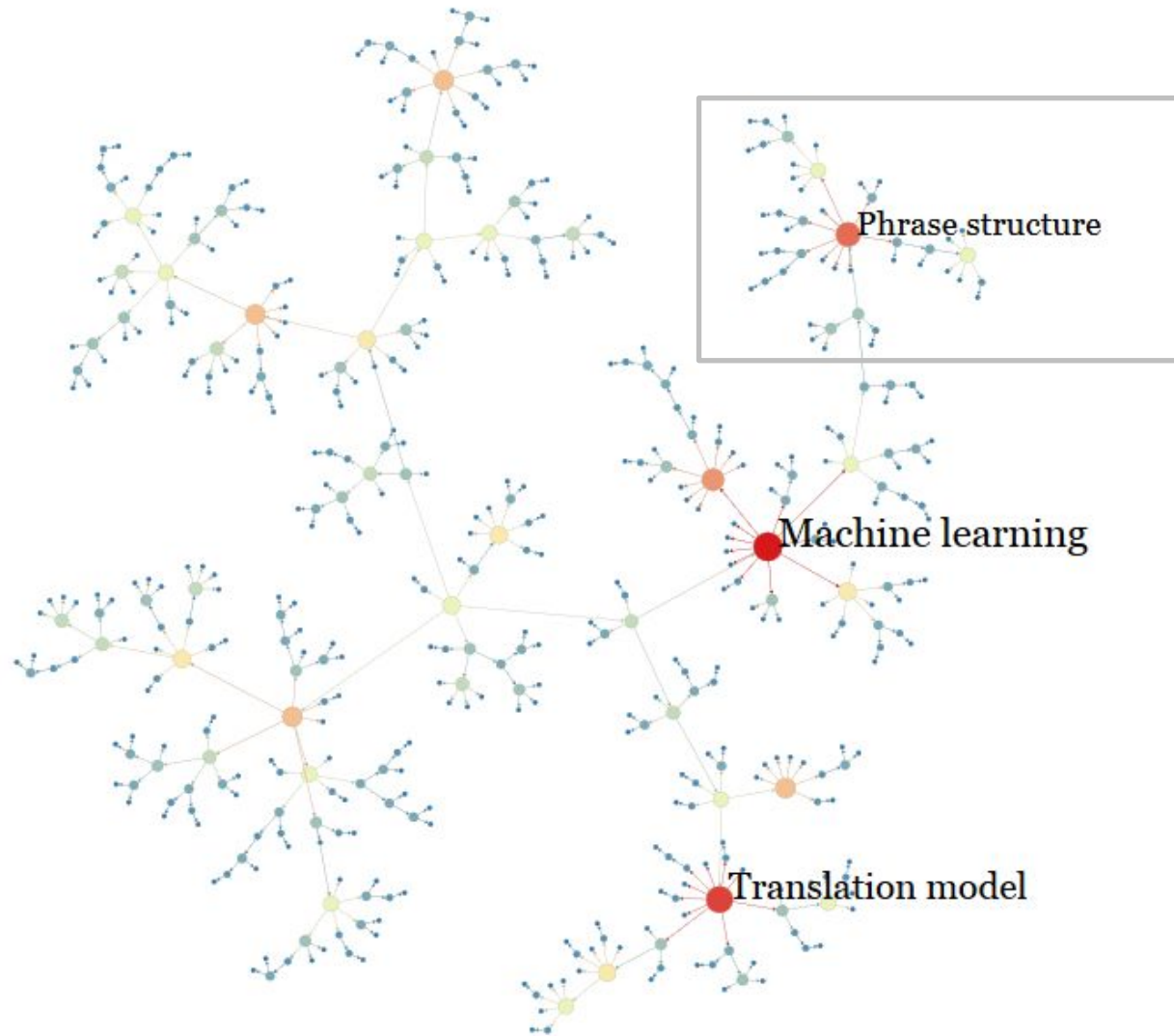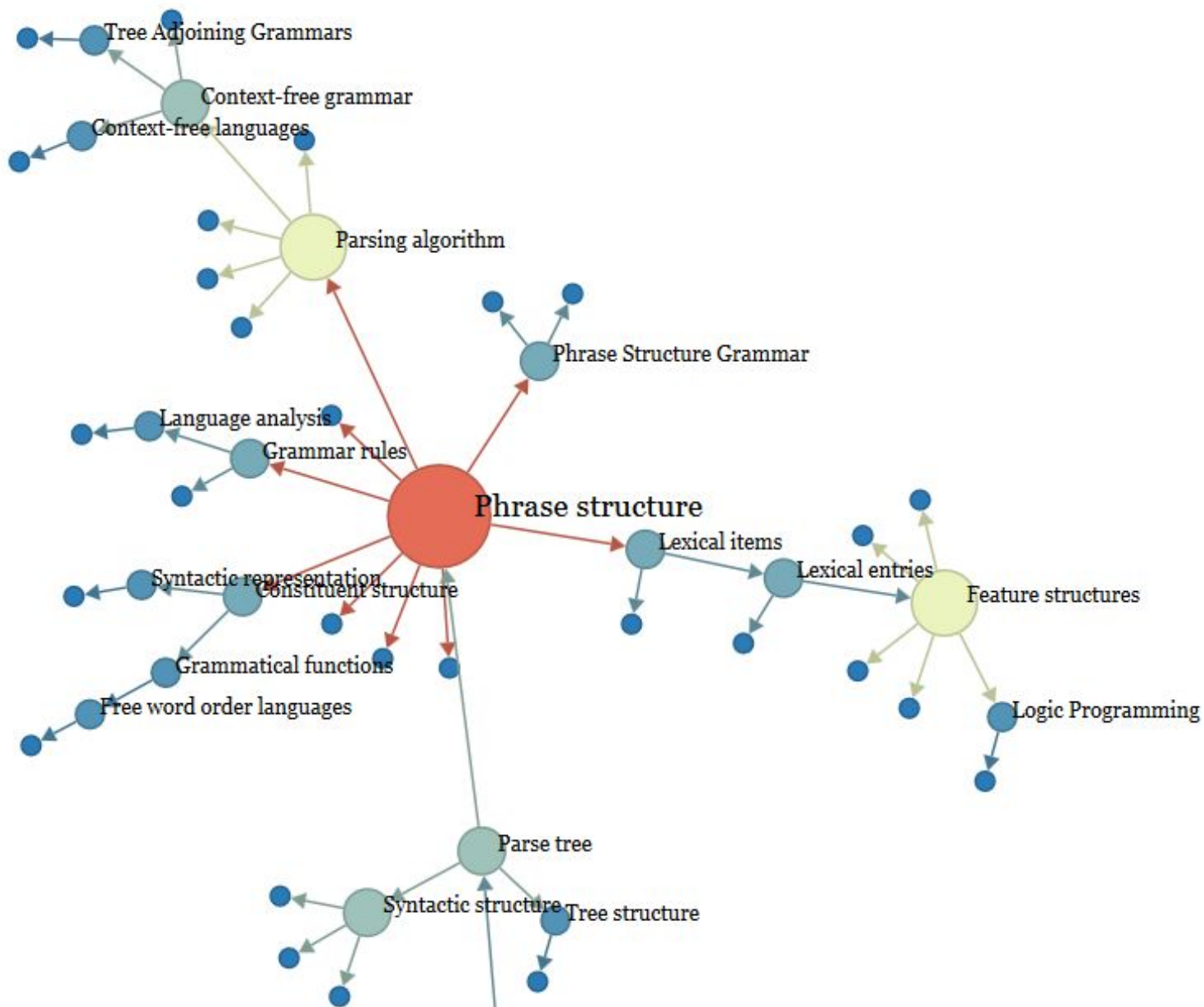- Directed to ensure most general nodes are at the top

# Terms to Taxonomy - ACL Anthology

1. Natural Language Processing
2. Natural Language
3. Language model
4. Statistical machine translation
5. Training data
6. Machine translation system
7. Machine translation
8. Hidden Markov Models
9. Support vector machines
10. Information retrieval
11. N-gram language model
12. Machine learning
13. Word sense disambiguation
14. Target language
15. Computer Science
16. Knowledge base
17. Human Language Technology
18. Translation model
19. Predicate-argument structure
20. Speech recognition
21. Natural language understanding
22. Feature structures
23. Spoken language systems
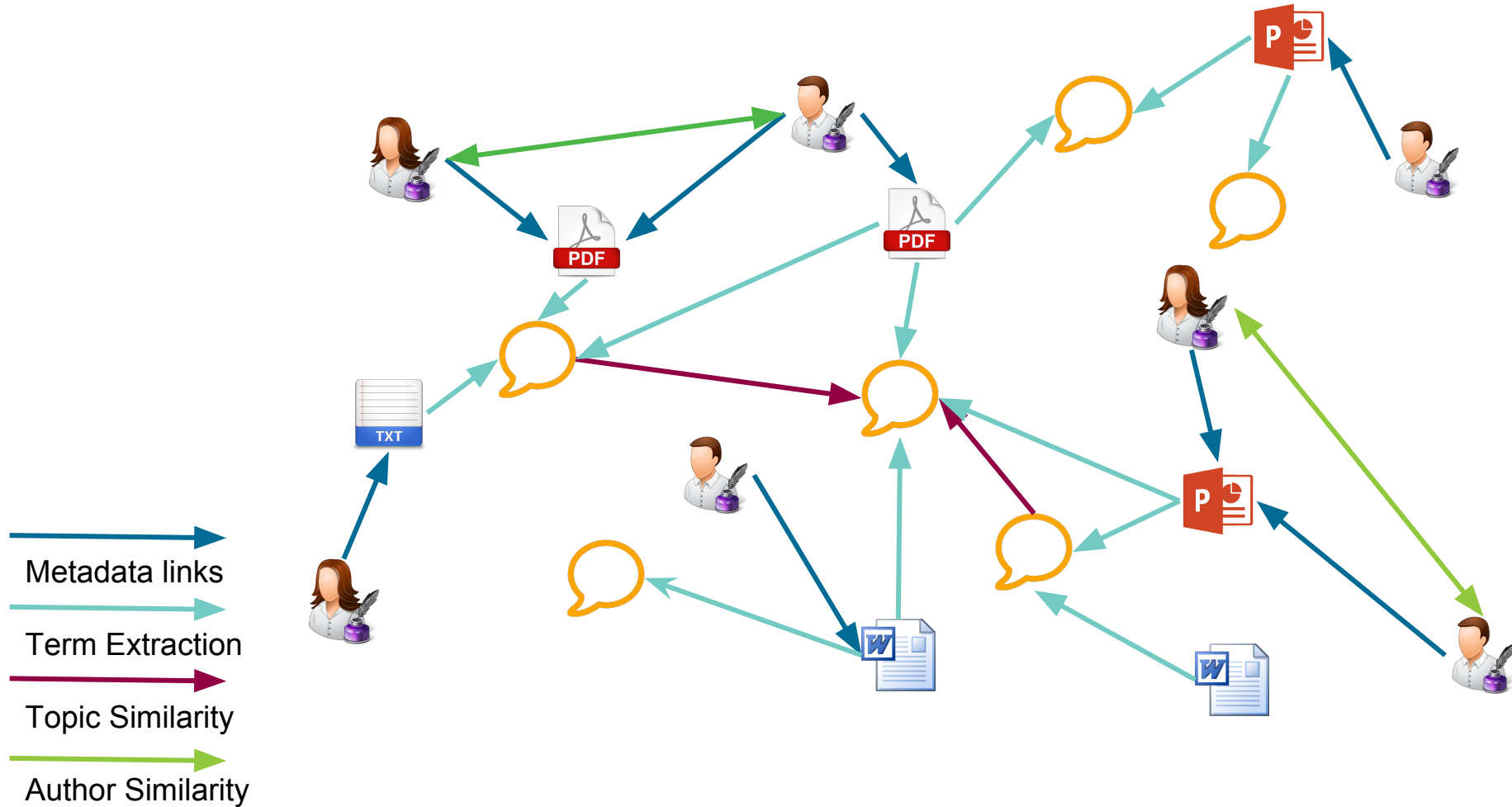24. Language Processing
25. Syntactic structure

# Taxonomy Extraction – ACL Anthology



Phrase structure

Machine learning

Translation model

# Heterogeneous graph



Metadata links

Term Extraction

Topic Similarity

Author Similarity

# Industry Applications

Content Analysis for Book Recommendation





Semantic Search on Digital News Archives

# THE IRISH TIMES

**Smart Insight Extractor**

ADVANCED SEARCH & E

http://smartie.ie/

RELATED ARTICLES

TIMELINE
Needs Briefcase

FEATURES
Needs Briefcase

AUTHORS

# THE IRISH TIMES
## Advanced Search & Browse

**BUSINESS**     CULTURE     LIFE & STYLE     NEWS     POLITICS     SPORT

search...     🔍

2014 | **2013** | 2012 | 2011 | 2010 | 2009 | 2005 | 2002_2007 |

Home

## Taxonomy



Consumer Price Index

Chief executive

Euro zone

Financial services

Federal Reserve

### Hot Tags                    More ›

National Asset Management Agency

Financial services

National Treasury Management Agency

Stoxx Europe 600 Index

Banking system

Property market

Technology companies

Consumer Price Index

Public service

Property prices

Asset Management

European Stability Mechanism

Property tax

Central Bank

FTSE 100 Index

Commercial property

Bailout programme

Operating system

International Financial Services Centre

Local property tax

Euro zone

Property developer

Business model

# Towards Saffron 3

- **Saffron** was developed primarily by Georgeta Bordea, Barry Coughlan (and many others)
- Technical improvements
  - One language (Java), one database (Lucene), one build system (Maven) etc.
  - Refactor code with existing libraries
    - V2.0: 14,500 Java LoC, 35,919 Python LoC
    - V3.0: 7,000 Java LoC

# Towards Saffron 3

- **Saffron** has attracted a lot of research and commercial attention
- But, **Saffron** is more importantly a research project.
- Next Step: Establish new baseline for
  - Term Extraction
    - Based on Astrakhanstev 2017
  - Taxonomy Learning
    - Use TExEval datasets (WordNet, E
    - New datasets that are taxonomic, ACM Computing Classification Sys
- Then: New **algorithms** :)

N. Astrakhantsev. ATR4S: Toolkit with State-of-the-art Automatic Terms Recognition Methods in Scala.
https://arxiv.org/abs/1611.07804
TExEval @ SemEval 2016: http://alt.qcri.org/semeval2016/task13/

# Conclusion

- Big document collections are **hard to understand**
  - In Academia
  - In Industry
- **Taxonomies** are the natural way to explore datasets
  - Evaluating the quality of a taxonomy is very hard
- Author metadata for documents lets us understand and **find experts**
- **Heterogeneous** graphs give new options for exploring document collections