# Candidate Extraction Impact on Automatic Keyphrase Extraction

**Adrien Bougouin** and **Florian Boudin** and **Béatrice Daille**
Université de Nantes, LINA, France
{adrien.bougouin,florian.boudin,beatrice.daille}@univ-nantes.fr

## Abstract

8+2 pages maximum...

## 1 Introduction

Keyphrases are single or multi-word expressions that represent the main topics of a document. Keyphrases are useful in many tasks such as information retrieval (Medelyan and Witten, 2008), document summarization (Litvak and Last, 2008) or document clustering (Han et al., 2007). Although scientific articles usually provide them, most of the documents have no associated keyphrases. Therefore, the problem of automatically assigning keyphrases to documents is an active field of research.

Automatic keyphrase extraction methods are divided into two categories: supervised and unsupervised methods. Supervised methods typically recast keyphrase extraction as a binary classification task (Witten et al., 1999; Sujian et al., 2003; Eichler and Neumann, 2010). For unsupervised methods, keyphrase extraction is often considered as a ranking task and many approaches are used (Barker and Cornacchia, 2000; Tomokiyo and Hurst, 2003; Mihalcea and Tarau, 2004). As distinct as they are, both supervised and unsupervised methods rely on a preliminary candidate extraction step which identifies single and multi-word expressions that have the same syntactic properties than a keyphrase. These expressions are the only textual units that can be extracted as keyphrases.

In this paper, we focus on the candidate extraction step and show its impact on the performance of automatic keyphrase extraction. Various methods are commonly employed to extract keyphrase candidates[1]. Usually, a set of either

---

[1]In this work, we do not consider methods which use a manually defined controlled vocabulary.

single words, n-grams filtered by stop words, NP-chunks or sequences of words matching given patterns is extracted (Hulth, 2003). According to the chosen method, the extracted set contains more or less candidates, and the amount of these that match with the ground truth keyphrases may vary. Hence, a few questions arise. How the different sets influence the keyphrase extraction? Do large candidate sets introduce noise that affects the performance of some keyphrase extraction methods?

We seek to better understand the impact of candidate extraction methods on keyphrase extraction by studying the aforementioned questions. We first quantify the differences between the candidate sets obtained by the commonly used methods. Also, we propose to use another method developed to extract noun-phrases for document indexing (Evans and Zhai, 1996) and we argue that such term detection method (Castellví et al., 2001) provides solid keyphrase candidates. Then, we evaluate the impact of the candidate extraction methods on three dissimilar keyphrase extraction methods. We select KEA (Witten et al., 1999) to represent supervised methods, TF-IDF (Spärck Jones, 1972) to represent unsupervised methods that require a collection of documents and TopicRank (Bougouin et al., 2013) to represent unsupervised methods that only make use of the document to analyse.

Results show that...

## 2 Definition of Candidate Keyphrases

Candidate keyphrases are textual units which can be selected as keyphrases of the document they are extracted from. Hence, they must have the same syntactic and linguistic properties than ground truth keyphrases. This section aims to determine those properties by analysing three standard evaluation datasets, for keyphrase extraction, and by providing statistics about their reference

keyphrases (ground truth keyphrases).

### 2.1 Keyphrase Extraction Datasets

Keyphrase extraction datasets are used to train or evaluate keyphrase extraction methods. Hence, the datasets are collections of documents paired with reference keyphrases, given by authors, readers or both. Unlike the studied methods, human annotators do not only extract keyphrases which are contained into the document. This problem of missing keyphrases leads to a bias of the training or evaluation of keyphrase extraction methods. In this work, we use three standard datasets which differ in terms of document type and/or language. The problem of missing keyphrases is partially bypassed using their stemmed forms during comparisons, when training or evaluating methods.

The **DUC** dataset (Over, 2001) is a collection of 308 English news articles covering about 30 news topics. This is the part of the dataset made for the DUC 2001 summarization evaluation campaign that has been annotated by Wan and Xiao (2008) for keyphrase extraction evaluation purpose. We split this into two sets: a training set containing 208 documents and a test set containing 100 documents.

The **SemEval** dataset (Kim et al., 2010) contains 284 English papers collected from the ACM Digital Libraries (conference and workshop papers). The 284 scientific papers are divided into three sets: a trial set containing 40 documents (unused in this work), a training set containing 144 documents and a test set containing 100 documents. As for the associated keyphrases, these are provided by both authors and readers.

The **DEFT** dataset (Paroubek et al., 2012) is a collection of 244 French scientific papers that belongs to the Humanities and Social Sciences domain. As SemEval, DEFT is divided into threee sets: a trial set containing 50 documents (not used in this work), a training set containing 141 documents and a test set containing 93 documents. Unlike DUC and SemEval, the only available reference keyphrases are the ones given by authors.

Table 2 gives statistics about the datasets. As we aim to use these statistics to lead this work, we restrain the discussion to observations made with the training sets.

### 2.2 Keyphrase Analysis

This section focuses on the reference keyphrase statistics presented in Table 2. The aim is to deter-
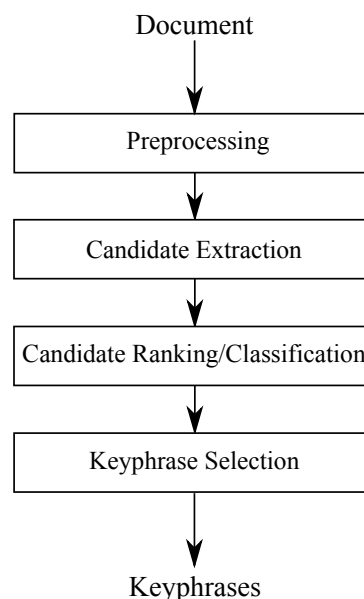


Figure 1: Processing steps of automatic keyphrase extraction methods.

mine the syntactic properties of most keyphrases, for English (combining information from DUC and SemEval) and for French (using DEFT information).

- Les Keyphrases sont principalement des unigrammes et bi-grammes (mais une tendance inversée entre anglais et français). - Presque toutes les keyphrase contiennent un nom. - L'usage d'adjectifs est fréquent (40 et 30%). - Très peu de verbes. - Usage fréquent de prépositions et determinants pour le français (uniquement).

DONNER DES EXEMPLES

Donner les séquences de POS les plus fréquentes dans le gold standard.

## 3 Candidate Extraction

Objectif + pré-requis.

### 3.1 N-Gram Extraction

### 3.2 NP-Chunk Extraction

### 3.3 Pattern Matching

### 3.4 Term Extraction

## 4 Keyphrase Extraction

Fonctionnement général.

|  | Statistics | Corpora | | |
|---|---|---|---|---|
|  |  | DUC | SemEval | DEFT |
| **Documents** | Language | English | English | French |
|  | Type | News | Papers | Papers |
|  | Documents | 208 | 144 | 141 |
|  | Tokens/document |  | 5134.6 | 7276.7 |
|  | Keyphrases/document | 8.1 | 15.4 | 5.4 |
|  | Missings keyphrases |  | 13.5% | 18.2% |
| **Keyphrases** | Unigrams | 26.2% | 20.2% | 66.4% |
|  | Bigrams | 54.1% | 53.4% | 20.7% |
|  | Trigrams and more | 19.7% | 26.4% | 12.9% |
|  | Containing nouns | 90.8% | 95.9% | 79.3% |
|  | Containing proper nouns | 18.7% | 5.8% | 16.8% |
|  | Containing adjectives | 41.6% | 40.5% | 28.8% |
|  | Containing verbs | 0.9% | 3.4% | 0.5% |
|  | Containing adverbs | 1.3% | 0.6% | 0.5% |
|  | Containing prepositions | 0.2% | 1.2% | 12.7% |
|  | Containing determiners | 0.0% | 0.0% | 8.1% |
|  | Containing others | 1.3% | 2.1% | 5.8% |

Table 1: Training dataset statistics. As a matter of consistency regarding the training and the evaluation of keyphrase extraction methods, the percentage of missing keyphrase is determined based on the stemmed form of the reference keyphrases.

### 4.1 TF-IDF

### 4.2 TopicRank

### 4.3 KEA

## 5 Evaluation

Expliquer les deux évaluations: intrinsèque et extrinsèque.

### 5.1 Experimental Setting

### 5.2 Candidate Extraction

Donner le rappel max et comparer avec la taille des différents ensemble.

Quels sont les termes candidats communs aux ensembles, les propriétés ?

### 5.3 Keyphrase Extraction

Quelles sont les performances de chaque méthode avec chaque ensemble de termes candidats ?

## References

Ken Barker and Nadia Cornacchia. 2000. Using Noun Phrase Heads to Extract Document Keyphrases. In *Proceedings of the 13th Biennial Conference of the Canadian Society on Computational Studies of Intelligence: Advances in Artificial Intelligence*, pages 40–52, London, UK. Springer-Verlag.

Adrien Bougouin, Florian Boudin, and Béatrice Daille. 2013. Topicrank: Graph-Based Topic Ranking for Keyphrase Extraction. In *Proceedings of the 6th International Joint Conference on Natural Language Processing (IJCNLP)*, Nagoya, Japan, October.

M. Teresa Cabré Castellví, Rosa Estopà Bagot, and Jordi Vivaldi Palatresi. 2001. Automatic Term Detection: A Review of Current Systems. *Recent Advances in Computational Terminology*, pages 53–88.

Kathrin Eichler and Günter Neumann. 2010. DFKI KeyWE: Ranking Keyphrases Extracted from Scientific Articles. In *Proceedings of the 5th International Workshop on Semantic Evaluation*, pages 150–153, Stroudsburg, PA, USA. Association for Computational Linguistics.

David A. Evans and Chengxiang Zhai. 1996. Noun-Phrase Analysis in Unrestricted Text for Information Retrieval. In *Proceedings of the 34th annual meeting of the Association for Computational Linguistics*, pages 17–24, Stroudsburg, PA, USA. Association for Computational Linguistics.

Juhyun Han, Taehwan Kim, and Joongmin Choi. 2007. Web Document Clustering by Using Automatic Keyphrase Extraction. In *Proceedings of the 2007*

| Statistics | Corpora | | |
|---|---|---|---|
| | DUC | SemEval | DEFT |
| Language | English | English | French |
| Type | News | Papers | Papers |
| Documents | 100 | 100 | 93 |
| Tokens/document | | 5179.6 | 6844.0 |
| Keyphrases/document | | 14.7 | 5.2 |
| Missings keyphrases | | | |

Table 2: Test dataset statistics. As a matter of consistency regarding the training and the evaluation of keyphrase extraction methods, the percentage of missing keyphrase is determined based on the stemmed form of the reference keyphrases.

*IEEE/WIC/ACM International Conferences on Web Intelligence and Intelligent Agent Technology*, pages 56–59, Washington, DC, USA. IEEE Computer Society.

Anette Hulth. 2003. Improved Automatic Keyword Extraction Given More Linguistic Knowledge. In *Proceedings of the 2003 Conference on Empirical Methods in Natural Language Processing*, pages 216–223, Stroudsburg, PA, USA. Association for Computational Linguistics.

Su Nam Kim, Olena Medelyan, Min-Yen Kan, and Timothy Baldwin. 2010. SemEval-2010 task 5: Automatic Keyphrase Extraction from Scientific Articles. In *Proceedings of the 5th International Workshop on Semantic Evaluation*, pages 21–26, Stroudsburg, PA, USA. Association for Computational Linguistics.

Marina Litvak and Mark Last. 2008. Graph-Based Keyword Extraction for Single-Document Summarization. In *Proceedings of the Workshop on Multi-Source Multilingual Information Extraction and Summarization*, pages 17–24, Stroudsburg, PA, USA. Association for Computational Linguistics.

Olena Medelyan and Ian H. Witten. 2008. Domain-Independent Automatic Keyphrase Indexing with Small Training Sets. *Journal of the American Society for Information Science and Technology*, 59(7):1026–1040, may.

Rada Mihalcea and Paul Tarau. 2004. TextRank: Bringing Order Into Texts. In Dekang Lin and Dekai Wu, editors, *Proceedings of the 2004 Conference on Empirical Methods in Natural Language Processing*, pages 404–411, Barcelona, Spain, July. Association for Computational Linguistics.

Paul Over. 2001. Introduction to DUC-2001: an Intrinsic Evaluation of Generic News Text Summarization Systems. In *Proceedings of DUC 2001 Document Understanding Conference*.

Patrick Paroubek, Pierre Zweigenbaum, Dominic Forest, and Cyril Grouin. 2012. Indexation libre et contrôlée d'articles scientifiques. Présentation et résultats du défi fouille de textes DEFT2012 (Controlled and Free Indexing of Scientific Papers. Presentation and Results of the DEFT2012 Text-Mining Challenge) [in French]. In *Proceedings of the Joint Conference JEP-TALN-RECITAL 2012, Workshop DEFT 2012: DÉfi Fouille de Textes (DEFT 2012 Workshop: Text Mining Challenge)*, pages 1–13, Grenoble, France, June. ATALA/AFCP.

Karen Spärck Jones. 1972. A Statistical Interpretation of Term Specificity and its Application in Retrieval. *Journal of Documentation*, 28(1):11–21.

Li Sujian, Wang Houfeng, Yu Shiwen, and Xin Chengsheng. 2003. News-Oriented Keyword Indexing with Maximum Entropy Principle. In *Proceedings of the 17th Pacific Asia Conference*. COLIPS Publications.

Takashi Tomokiyo and Matthew Hurst. 2003. A Language Model Approach to Keyphrase Extraction. In *Proceedings of the ACL 2003 Workshop on Multiword Expressions: Analysis, Acquisition and Treatment - Volume 18*, pages 33–40, Stroudsburg, PA, USA. Association for Computational Linguistics.

Xiaojun Wan and Jianguo Xiao. 2008. Single Document Keyphrase Extraction Using Neighborhood Knowledge. In *Proceedings of the 23rd National Conference on Artificial Intelligence - Volume 2*, pages 855–860. AAAI Press.

Ian H. Witten, Gordon W. Paynter, Eibe Frank, Carl Gutwin, and Craig G. Nevill Manning. 1999. KEA: Practical Automatic Keyphrase Extraction. In *Proceedings of the 4th ACM Conference on Digital Libraries*, pages 254–255, New York, NY, USA. ACM.

| Methods | DUC | | SemEval | | DEFT | |
|---|---|---|---|---|---|---|
| | Candidates | Rmax | Candidates | Rmax | Candidates | Rmax |
| {1..2}-grams | | | | | | |
| {1..3}-grams | | | | | | |
| {1..4}-grams | | | | | | |
| {1..5}-grams | | | | | | |
| NP chunks | | | | | | |
| Longest NPs | | | | | | |
| Best patterns | | | | | | |
| TermSuite | | | | | | |
| CLARIT'96 | | | | | | |

Table 3: Candidate extraction statistics. Rmax stands for maximum recall, i.e. it is the percentage of candidates that match with reference keyphrases.

| Methods | DUC | | | SemEval | | | DEFT | | |
|---|---|---|---|---|---|---|---|---|---|
| | P | R | F | P | R | F | P | R | F |
| {1..2}-grams | 00.0 | 00.0 | 00.0 | 00.0 | 00.0 | 00.0 | 00.0 | 00.0 | 00.0 |
| {1..3}-grams | 00.0 | 00.0 | 00.0 | 00.0 | 00.0 | 00.0 | 00.0 | 00.0 | 00.0 |
| {1..4}-grams | 00.0 | 00.0 | 00.0 | 00.0 | 00.0 | 00.0 | 00.0 | 00.0 | 00.0 |
| {1..5}-grams | 00.0 | 00.0 | 00.0 | 00.0 | 00.0 | 00.0 | 00.0 | 00.0 | 00.0 |
| NP chunks | 00.0 | 00.0 | 00.0 | 00.0 | 00.0 | 00.0 | 00.0 | 00.0 | 00.0 |
| Longest NPs | 00.0 | 00.0 | 00.0 | 00.0 | 00.0 | 00.0 | 00.0 | 00.0 | 00.0 |
| Best patterns | 00.0 | 00.0 | 00.0 | 00.0 | 00.0 | 00.0 | 00.0 | 00.0 | 00.0 |
| TermSuite | 00.0 | 00.0 | 00.0 | 00.0 | 00.0 | 00.0 | 00.0 | 00.0 | 00.0 |
| CLARIT'96 | 00.0 | 00.0 | 00.0 | 00.0 | 00.0 | 00.0 | 00.0 | 00.0 | 00.0 |

Table 4: Comparison of candidate extraction methods, when extracting 10 keyphrases with the **TF-IDF** method. Results are expressed as a percentage of precision (P), recall (R) and f-score (F).

| Methods | DUC | | | SemEval | | | DEFT | | |
|---|---|---|---|---|---|---|---|---|---|
| | P | R | F | P | R | F | P | R | F |
| {1..2}-grams | 00.0 | 00.0 | 00.0 | 00.0 | 00.0 | 00.0 | 00.0 | 00.0 | 00.0 |
| {1..3}-grams | 00.0 | 00.0 | 00.0 | 00.0 | 00.0 | 00.0 | 00.0 | 00.0 | 00.0 |
| {1..4}-grams | 00.0 | 00.0 | 00.0 | 00.0 | 00.0 | 00.0 | 00.0 | 00.0 | 00.0 |
| {1..5}-grams | 00.0 | 00.0 | 00.0 | 00.0 | 00.0 | 00.0 | 00.0 | 00.0 | 00.0 |
| NP chunks | 00.0 | 00.0 | 00.0 | 00.0 | 00.0 | 00.0 | 00.0 | 00.0 | 00.0 |
| Longest NPs | 00.0 | 00.0 | 00.0 | 00.0 | 00.0 | 00.0 | 00.0 | 00.0 | 00.0 |
| Best patterns | 00.0 | 00.0 | 00.0 | 00.0 | 00.0 | 00.0 | 00.0 | 00.0 | 00.0 |
| TermSuite | 00.0 | 00.0 | 00.0 | 00.0 | 00.0 | 00.0 | 00.0 | 00.0 | 00.0 |
| CLARIT'96 | 00.0 | 00.0 | 00.0 | 00.0 | 00.0 | 00.0 | 00.0 | 00.0 | 00.0 |

Table 5: Comparison of candidate extraction methods, when extracting 10 keyphrases with **TopicRank**. Results are expressed as a percentage of precision (P), recall (R) and f-score (F).

| Methods | DUC | | | SemEval | | | DEFT | | |
|---|---|---|---|---|---|---|---|---|---|
| | P | R | F | P | R | F | P | R | F |
| {1..2}-grams | 00.0 | 00.0 | 00.0 | 00.0 | 00.0 | 00.0 | 00.0 | 00.0 | 00.0 |
| {1..3}-grams | 00.0 | 00.0 | 00.0 | 00.0 | 00.0 | 00.0 | 00.0 | 00.0 | 00.0 |
| {1..4}-grams | 00.0 | 00.0 | 00.0 | 00.0 | 00.0 | 00.0 | 00.0 | 00.0 | 00.0 |
| {1..5}-grams | 00.0 | 00.0 | 00.0 | 00.0 | 00.0 | 00.0 | 00.0 | 00.0 | 00.0 |
| NP chunks | 00.0 | 00.0 | 00.0 | 00.0 | 00.0 | 00.0 | 00.0 | 00.0 | 00.0 |
| Longest NPs | 00.0 | 00.0 | 00.0 | 00.0 | 00.0 | 00.0 | 00.0 | 00.0 | 00.0 |
| Best patterns | 00.0 | 00.0 | 00.0 | 00.0 | 00.0 | 00.0 | 00.0 | 00.0 | 00.0 |
| TermSuite | 00.0 | 00.0 | 00.0 | 00.0 | 00.0 | 00.0 | 00.0 | 00.0 | 00.0 |
| CLARIT'96 | 00.0 | 00.0 | 00.0 | 00.0 | 00.0 | 00.0 | 00.0 | 00.0 | 00.0 |

Table 6: Comparison of candidate extraction methods, when extracting 10 keyphrases with **KEA**. Results are expressed as a percentage of precision (P), recall (R) and f-score (F).