

Digging in the Dirt: Extracting Keyphrases from Texts with KD

Giovanni Moretti, Rachele Sprugnoli, Sara Tonelli
Digital Humanities Research Unit
Fondazione Bruno Kessler

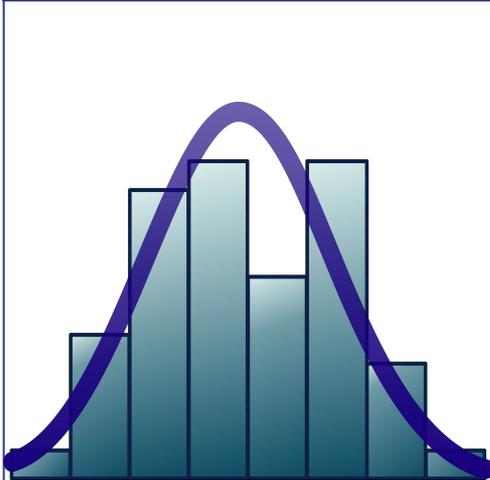
Website: <http://dh.fbk.eu>

Twitter: [@DH_FBK](https://twitter.com/DH_FBK)



KD = Keyphrase Digger

STATISTICAL MEASURES



LINGUISTIC INFORMATION



WEIGHTED LIST OF KEYPHRASES

rank	keyword synonyms	frequency	score
1	kd	33	43.07
2	types of texts	2	23.24
3	kx	17	22.36
4	multi-token expressions	4	14.04
5	level of customizability	2	12.64
6	keyphrase extraction	5	9.55
7	humanities scholars	2	8.31
8	document collections	4	7.08
9	configuration	4	6.87
10	set of key-concepts	2	6.63

NGRAMS capturing the main concepts of a document

Reimplementation of KX (Pianta & Tonelli, 2010)

Challenge 1

Better extraction performances



Challenge 2

Higher processing speed



Challenge 3

More customizability



Challenge 1

Extraction Performances

- SEMEVAL 2010 dataset 
 - scientific papers
 - best performing rule-based system

KX: 7th place  KD: 2nd place

- DE GASPERI corpus 
 - historical documents

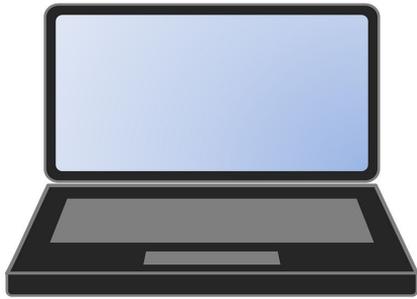
Precision over expert keyphrases: 42%



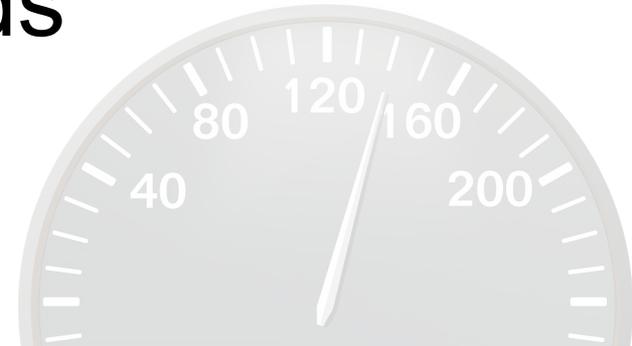
Challenge 2

Processing Speed

Same machine, same corpus (101K tokens), comparable settings



- **KX: 3.4 minutes**
- **KD: 7 seconds**



Challenge 3

Customizability

- Accepts different pre-processed texts
e.g. TEXTPRO - TREETAGGER - STANFORD
- Uses custom configuration files
e.g. POS-PATTERNS - STOPWORDS
- Many options available
e.g. NUMBER of KEYPHRASES - NGRAM LENGTH
- Easily extendible to other languages



Applications and Availability

1. ALCIDE PLATFORM

[Analysis of Language and Content In a Digital Environment](#)

2. ON-LINE DEMO

http://celct.fbk.eu:8080/KD_KeyDigger/

3. SOFTWARE PACKAGE

<http://dh.fbk.eu/technologies/kd>

Digging in the Dirt: Extracting Keyphrases from Texts with KD

CENTER FOR INFORMATION TECHNOLOGY
DIGITAL HUMANITIES

KD = KEYPHRASE DIGGER

Keyphrase Digger is a rule-based system for keyphrase extraction. It is a **Java** re-implementation of KX tool (Pianta and Tonelli, 2010) with a **new architecture** and **new features**.

It combines **statistical measures** with **linguistic information** given by PoS patterns to identify and extract weighted keyphrases from texts.

MAIN FEATURES:

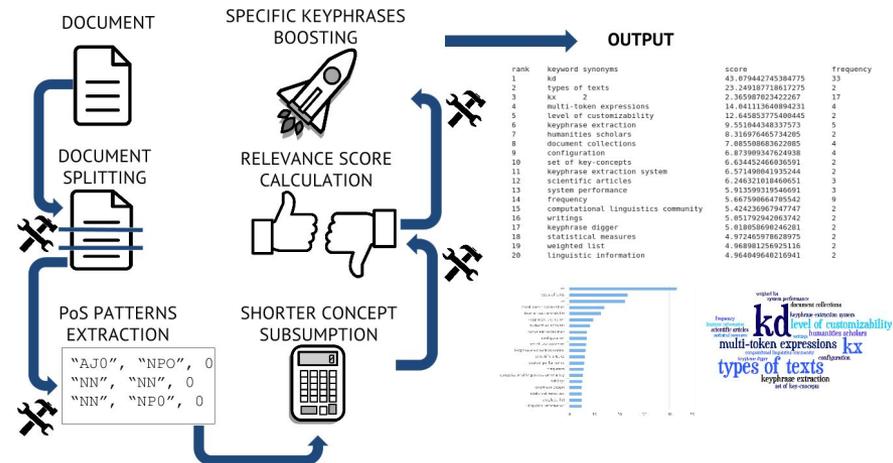
- Extraction of multi-words
- Multilinguality (EN and IT)
- Easily extendible to other languages
- Higher customizability than KX
- High processing speed
- Clustering of keyphrases under the same lemma
- Various accepted formats and POS tagsets
- Boost of specific PoS patterns
- Integration of Apache Lucene Library

NEW!



KD is about 25 times faster than KX:
very apt for web applications!

KD ARCHITECTURE



KD AVAILABILITY

ONLINE DEMO
http://celct.fbk.eu:8080/KD_KeyDigger/

JAVA RUNNABLE JAR and LIBRARY
<https://dh.fbk.eu/technologies/kd>

SEE YOU AT THE POSTER SESSION!