

A Method for Human-Interpretable Paraphrasticity Prediction

Maria Moritz¹, Johannes Hellrich^{2,3}, Sven Buechel²



LaTeCH@COLING, August 25 2018, Santa Fe, New Mexico, USA

Reuse in historical texts is difficult to identify if it is heavily modified, and algorithmic support for its identification & analysis is limited.

Existing reuse detection techniques can tell if and how frequently a text is modified, while our technique also determines the degree and characteristics of modification (i.e., the “features”) that constitute a reuse, which is an important prerequisite for the analysis of historical text reuse giving scholars hints for deeper investigation.

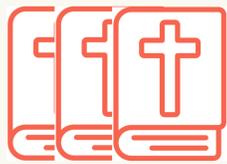
Our technique is both, human-interpretable and semantically informed setting it apart from recent developments based on distributional semantics which do not allow for easy manual inspection of individual models.

OVERVIEW

DATA



plagiarism corpus⁴



Bibles⁵ (1749-1859)
3 lit. translations
5 revisions



Bernard

ENG	ENG	Medieval LAT	language
5,000	315,000	990	POSITIVE
1,500	35,000	110	training
	(15 pairings)		test
5,000	270,000	990	NEGATIVE
1,500	30,000	110	training
rand. Sampled	(13 pairings)	randomly	test
(overlap of 4)		sampled	
PAN 2010 plagiarism detection challenge (Madnani et al. 2012)	www.biblestudytools.com www.mysword.info Parallel Text Project (Mayer & Cysouw 2014).	Bibindex (Mellerin 2014)	reference

⁴ Icons made by Freepik from Flaticon is licensed by CC 3.0 BY
⁵ Icons made by Smashicons from www.flaticon.com is licensed by CC 3.0 BY

1) mmoritz@etrap.eu



SPONSORED BY THE



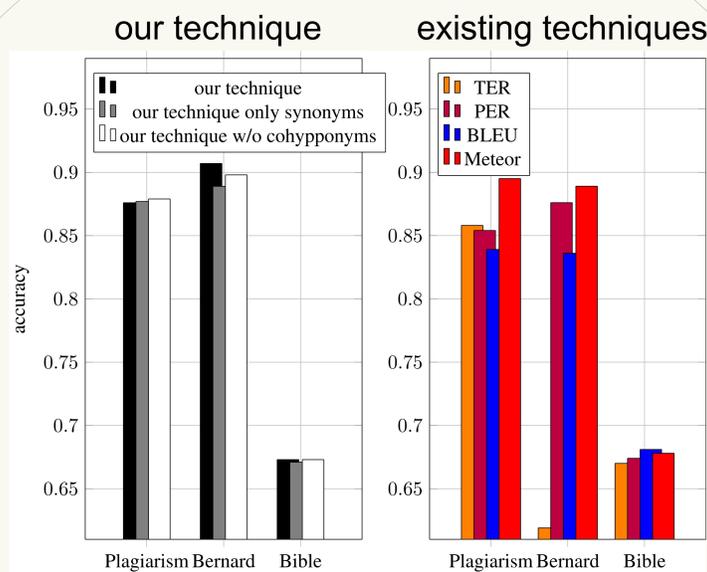
2) johannes.hellrich@uni-jena.de
sven-buechel@uni-jena.de



3) johannes.hellrich@uni-jena.de



RESULTS



Co-hyponym replacement is common in historical text reuse (Moritz et al. 2016). Our method is well suited to predict this reuse, because it considers

such substitutions. Meteor outperforms all other methods on the plagiarism data.

As such, our method is especially useful for applications in the humanities as operation frequencies and their feature weights are open to manual inspection.

Modeling reuse in historical text using semantic relations as we propose, achieves results comparable to using features derived from machine translation metrics.

In future work, we plan to tune parameters and to qualitatively analyze weaknesses of our method.

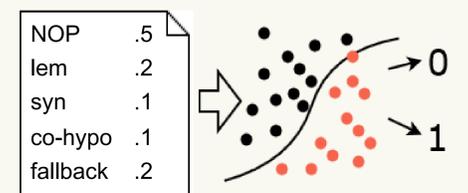
word alignment

they mount up with wings as eagles walk
lower syn NOP syn NOP lem
They raise up the pinion as eagles go

operation modeling

operation	example
No OPeration necessary	above, above
lower-casing match	LORD, Lord
normalizing match	desireth, desires
lemmatizing match	mine, my
derivation match	help, helper
short edit distance match	Phinehas, Phinees
words are synonyms	went, departed
word1 is hypernym of word2	coat, doublet
word1 is hyponym of word2	spears, arms
words are co-hyponyms	steps, feet
fallback	—

prediction of paraphrastic reuse using operation frequencies as features



OUR METHOD

DISCUSSION

We presented a method for paraphrase detection that describes reuse based on the frequency of specific modification operations, especially semantic relations beyond synonyms.