

# Non-Literal Text Reuse in Historical Texts: An Approach to Identify Reuse Transformations and its Application to Bible Reuse

Maria Moritz<sup>1</sup>, Andreas Wiederhold<sup>1</sup>, Barbara Pavlek<sup>2</sup>, Yuri Bizzoni<sup>3</sup>, Marco Büchler<sup>1</sup>

## Research Questions

### Motivation

Text reuse is the spoken and written repetition of text across time and space. It can be a quotation, an allusion or translation. Detection methods of historical text reuse are needed in different scholarly fields, e.g. to detect redundancies in digital libraries, to trace transmissions of historical thought or to identify fragmentary authors.

However, text is often modified during the reuse process, which makes the detection challenging. Therefore, we analyze the non-literal share in historical text reuse to obtain an understanding of the requirements for contemporary detection methods.

**1 What** is the extent of non-literal reuse in our datasets?

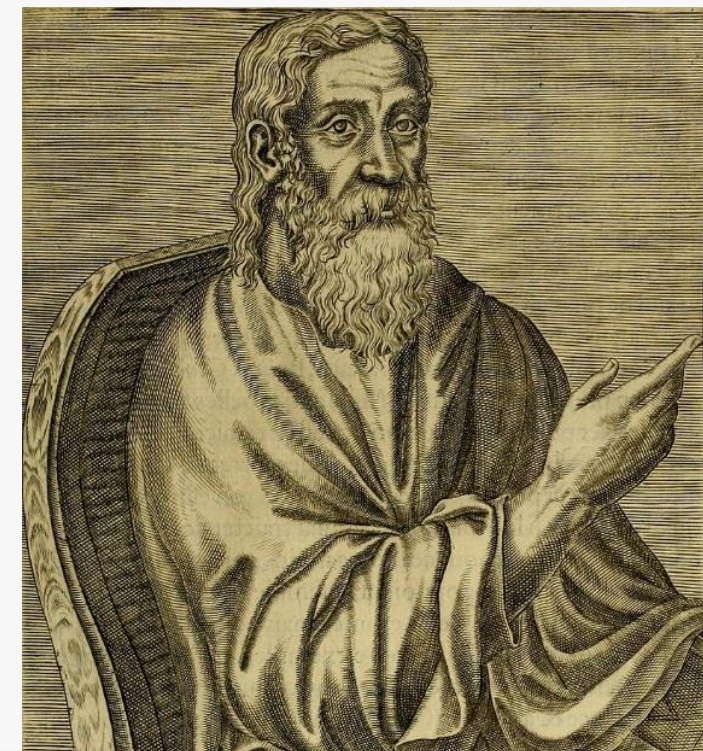
**2 How** is the non-literally reused text modified in our datasets?

**2.1 How** can linguistic resources support the discovery of non-literal reuse?

**2.2 What** are the limitations of an automated classification approach relying on linguistic resources?

## Data

Clement of Alexandria

Christian theologian  
from the 2<sup>nd</sup> century.

Bernard of Clairvaux

French abbot from  
the 12<sup>th</sup> century.We obtain Bible verse reuse pairs:  
199 & 162

literal	Bible verse	Bernard reuse
Prov 18 3	impus cum in profundum venerit peccatorum contemnit sed sequitur eum ignominia et obprobrium ( <i>When the wicked man is come into the depth of sins, also contempt comes but ignominy and reproach follow him</i> )	Impius, cum venerit in profundum malorum, contemnit ( <i>When the wicked man is come into the depth of evil</i> )
less literal	Bible verse	Clement reuse
1Cor 13 13	νυνὶ δὲ μένει πίστις, ἐλπίς, ἀγάπη, τὰ τρία ταῦτα μέζων δὲ τούτων ἡ ἀγάπη ( <i>And now remain faith, hope, love, these three; but the greatest of those is love.</i> )	ἀγάπην, πίστιν, ἐλπίδα ( <i>love, faith, hope – in accusative case</i> ) μένει δὲ τὰ τρία ταῦτα, πίστις, ἐλπίς, ἀγάπη· μέζων δὲ ἐν τούτοις ἡ ἀγάπη ( <i>and remain these three, faith, hope, love; but the greatest among them is love</i> )
non-literal	Bible verse	Clement reuse
Mt 12 35	ὁ ἀγαθὸς ἄνθρωπος ἐκ τοῦ ἀγαθοῦ θησαυροῦ ἐκβάλλει ἀγαθὰ, καὶ ὁ πονηρὸς ἄνθρωπος ἐκ τοῦ πονηροῦ θησαυροῦ ἐκβάλλει πονηρά. ( <i>A good man out of good storage brings out good things, and an evil man out of the evil storage brings evil things.</i> )	Ψυχῆς, τὰ δὲ ἐκτός, κἂν μὲν ἡ ψυχὴ χρήται καλῶς, καλὰ καὶ ταῦτα δοκεῖ, ἐὰν δὲ πονηρῶς, πονηρά, ὁ κελεύων ἀπαλλοτριῶν τὰ ὑπάρχοντα ( <i>[are in the] soul, and some are out, and if the soul uses them good, those things are also thought of as good, but if [used as] bad, [they are thought of] bad; he who commands the renouncement of possessions</i> )

## Methodology

Operation	Example
NOP(reuse_word, orig_word)	NOP(maledictus, maledictus)
lem(reuse_word, orig_word)	lem(penetrat, penetrabit)
repl_syn(reuse_word, orig_word)	repl_syn(magnificavit, glorificavit)
repl_hyper(reuse_word, orig_word)	hyper(cupit, habens)
repl_hypo(reuse_word, orig_word)	hypo(dederit, tollet)
repl_co-hypo(reuse_word, orig_word)	repl_co-hypo(magnificavit, fecit)
NOPmorph(reuse_tags, orig_tags)	NOPmorph(na, na)
repl_pos(reuse_tag, orig_tag)	repl_pos(n, a)
repl_case(reuse_tag, orig_tag)	repl_case(g, d)
lemma_missing(reuse_word, orig_word)	lemma_missing(tentari, inlectus)
no_rel_found(reuse_word, orig_word)	no_rel_found(gloria, arguitur)

**i We define operations (OPs)** reflecting literal reuse and semantic replacements (see above).

**ii Our algorithm** looks for identical & similar words and for morphological & semantic changes (see top right).

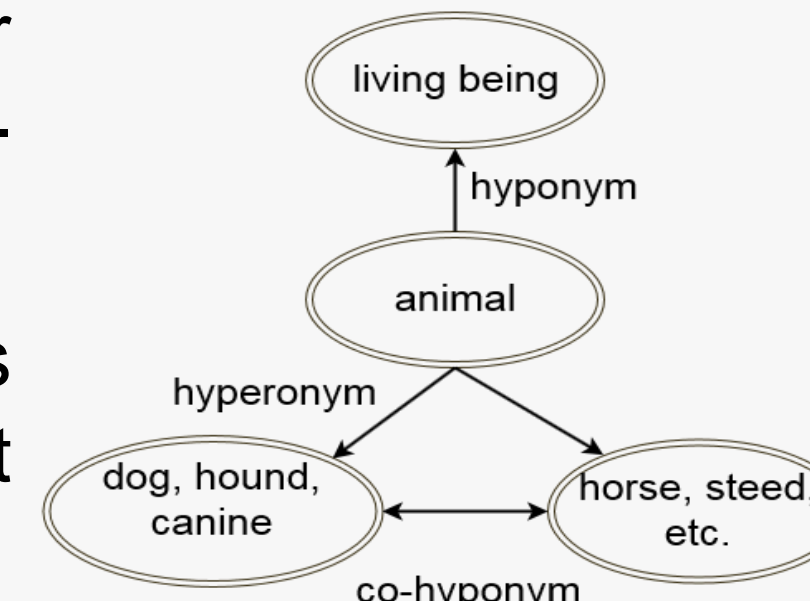
**iii We apply** both to our datasets using the Ancient Greek WordNet (see an example at the right).

**iv We complement the automated approach** by a manual analysis of a sample (60 & 100 instances resulting in 192 & 224 replacement operations) to find the limitations of our automated approach.

**OPs used manually:** ins(word), del(word) and NOP, lem, repl\_syn, repl\_hyper, repl\_hypo, repl\_co-hypo

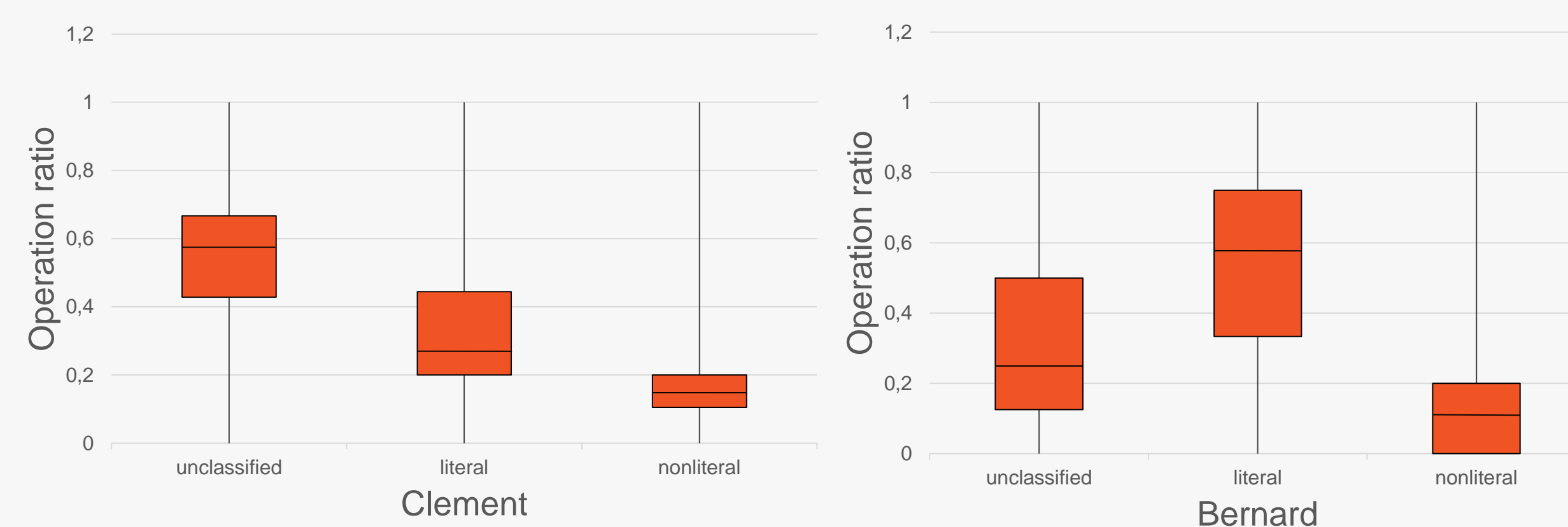
**Morphological information:** from Perseus' tag set (Bamman & Crane 2011), e.g. repl\_num\_s\_p

```
input : L ← set of word lemma pairs obtained from the lemma resources
input : S ← set of synsets from AGWN; each synset contains an id and a parent id
input : T ← list of words of reuse instance (containing part-of-speech information)
input : B ← list of words of Bible verse (containing part-of-speech information)
output: OP ← list of sets containing up to 3 parameterized operations
s1, s2 ← any two synsets ∈ S
tmp_op ← temporary variable which presents the absence of a relation but not of a lemma
for t in T do
  for b in B do
    if t=b then
      OP ← OP ∪ {NOP(t, b), check(morph(t), morph(b))}
    else if lowerCase(t) = b then
      OP ← OP ∪ {lower(t, b), check(morph(t), morph(b))}
    else if lowerCase(b) = t then
      OP ← OP ∪ {upper(t, b), check(morph(t), morph(b))}
    else if t ∈ L and b ∈ L then
      // lemma found for original (b) and reuse word (t)
      if connect(t) = connect(b) then
        OP ← OP ∪ {lem(t, b), check(morph(t), morph(b))}
      else if t ∈ s1 and b ∈ s2 and s1 ∈ S and s2 ∈ S then
        if s1 = s2 then
          // s1 is synonym of b
          OP ← OP ∪ {repl_syn(t, b)}
        else if s1 = parent(s2) then
          // s1 is hyperonym of b
          OP ← OP ∪ {repl_hyper(t, b)}
        else if parent(s1) = s2 then
          // s1 is hyponym of b
          OP ← OP ∪ {repl_hypo(t, b)}
        else if parent(s1) = parent(s2) then
          // s1 is co-hyponym of b and s2
          OP ← OP ∪ {repl_co-hypo(t, b)}
        // synset of t and synset of b both have the same
        // parent as present
        OP ← OP ∪ {repl_syn(t, b)}
    else
      tmp_op ← {no_rel_found(t, b)}
  end
end
return OP
```



## Results

**RQ1:** What is the extent of non-literal reuse in our datasets?



**Figure 1:** Ratios of operations in reuse instances. **literal:** NOP, lem, lower, etc.; **nonlit:** syn, hyper, etc.

The reuse is significantly non-literal and conceptualization might be preferred over stemming or semantic relations in the same POS category only.

**RQ2.1** How can linguistic resources support the discovery of non-literal reuse?

	literal				non-literal				unclassified		
	NOP	upper	lower	lem	syn	hyper	hypo	co-hypo	no_rel_found	lem_mssing	total
Clement	337	6	0	356	153	20	14	101	563	639	2189
Bernard	587	0	44	102	60	14	28	68	347	85	1335

Consider operations that successfully look up a lemma as:  
lem\_success={lem, syn, repl\_hyper, repl\_hypo, repl\_co-hypo, no\_rel\_found}, with lem\_missing holding tokens not found.

$$supp_{lem} = \frac{\sum_{Occ(o)} o \in \text{lem\_success}}{\sum_{Occ(o)} o \in \text{lem\_success} \cup \{\text{lem\_missing}\}}$$

$$supp_{lem}^{Clement} = 0.65 \text{ and } supp_{AGWN}^{Clement} = 0.34$$
$$supp_{lem}^{Bernard} = 0.88 \text{ and } supp_{AGWN}^{Bernard} = 0.33$$

**RQ2.2** What are the limitations of an automated classification approach relying on linguistic resources?

exception	quantity	
	Clement	Bernard
Word changed to antonym	1 <sup>1</sup>	0
Synonym and morphology changed	1	16
More than one morphological category changed	1	7
Synonym is multi-word expression	3	5
Many-to-many	0	12

Exceptions preventing applying our OPs.

1) "the God, the good (one)" (Clement) vs. "none is good but the God" (Bible).  
2) "judged calmly" (Bernard) vs. "fake friend" (Sal 12 18).

Language resources support the identification of reuse components. In our datasets, co-hyponyms are often used to rephrase an idea. Many-to-many relationships show that meanings can be hidden in structural or expert knowledge.

## Future plans

A more comprehensive study will strengthen the findings. For example using larger reuse datasets and additional languages, such as inflecting and non-inflecting languages.

A smarter automated approach for deriving an original text excerpt can be learning edit scripts, such as undetaken by Kehrer (2014) also considering the movement of reuse except within the reuse or the syntactical tree.

Deeper analyses of reuse statistics might be supported by the semantic relations that are presented in word nets.

### We aggregate

Bibindex' Lemmas  
(65.5K Biblical Greek entries; 315K Latin entries)

Classical Language Tool Kit (CLTK) (Johnson et al., 2014)  
954K Ancient Greek & 270K Latin entries

Greek New Testament of the Society of Biblical Literature<sup>3</sup> & Septuaginta<sup>4</sup>  
59.5K word-lemma-pairs

Ancient Greek WordNet (Bizzoni et al., 2014; Minozzi, 2009)  
99K synsets of which 33K contain Ancient Greek and 27K contain Latin words

### Greek Old Testament

Alfred Rahlfs, editor. 1935. Septuaginta, id est Vetus Testamentum Graece juxta LXX interpretes. Rahlfs. 2 vol., 1950.

### Greek New Testament

Kurt Aland and Barbara Aland, editors. 1966. The Greek New Testament. Deutsche Bibelgesellschaft-United Bible Societies, 27 edition.

### Latin Bible

Gribomont J. Weber R., Fischer B., editor. 1969, 1994, 2007.  
Biblia sacra juxta vulgatum versionem. Deutsche Bibelgesellschaft.

Clément d'Alexandrie, Quel riche sera sauvé ?, Quis dives salvetur, P. A. O à Sources Chrétiennes, col. 537, p. 100 ff., 2011.

### We acknowledge

Laurence Mellerin for providing the datasets and for valuable advice on their content. The German Federal Ministry of Education and Research for funding the work (grant 01UG1509).