Stefanie Dipper / Julia Krasselt / Simone Schultz-Balluff

# Creating synopses of 'parallel' historical manuscripts and early prints
## Alignment guidelines, evaluation, and applications

**Abstract**

In this paper we introduce the task of aligning parallel historical texts, to create synopses for comparing similarities and deviations between them. We present guidelines for manually annotating corresponding words and phrases. A test annotation reveals that there is considerable high inter-annotator agreement, ranging from kappa = 0.76 to 0.98, depending on the specific text. In an application scenario we show a typical use case for which token and phrase alignments are of value.

## 1. Introduction

In this paper we introduce methods for aligning parallel texts. Alignment refers to the task of linking corresponding elements (words, phrases, paragraphs, etc.) between related documents. Documents can be related because, e.g., they go back to the same source or one document is a (close or loose) copy of the other. As a result, the documents are similar to each other to different degrees and deviations can provide interesting clues for linguistic or historico-cultural investigations.

The goal of this paper is to present and evaluate guidelines we developed for aligning different versions of the medieval passion treatise *Interrogatio Sancti Anselmi de Passione Domini* (henceforth referred to as *Anselm*). In this text, St. Anselm of Canterbury asks Virgin Mary to reveal the passion of her son Jesus Christ, beginning with the Last Supper and ending with his resurrection. There are 70 vernacular manuscripts and prints, which date from the 14th–16th century and represent dialects of Early New High German (ENHG), Middle Low German and Middle Dutch.[1,2] The degree of similarity

---

[1] The further textual evidence covers 162 Latin texts and one Middle English text. The text is transmitted in different versions: verse, short prose, and long prose. For the alignment task introduced in this paper, we only deal with prose versions.

[2] Currently, 64 texts are transcribed and available for analysis (51 manuscripts and 10 prints in German, two manuscripts and one print in Dutch. 50 texts are preserved completely, 14 texts are fragments).

between individual *Anselm* texts varies[3]: there are pairs which are almost identical as well as pairs that differ to a high degree, e.g. in vocabulary, word order, or content. In our guidelines, we model these differences with different annotation layers, which range from token level alignments to phrasal alignments.

There is a variety of conceivable applications where the marking of concordant elements is of value. Alignments can be used for analyzing certain kinds of linguistic variance (e.g. in word order patterns). Furthermore, they can be of interest for philological research, e.g. in preparing printed or digital editions with a critical apparatus.

The paper is organized as follows. In Sec. 2, we address related work. Sec. 3 and 4 present the annotation guidelines and the results from our test annotation. In Sec. 5 we present an exemplary application where the results of the annotation process can be used for answering philological questions.

## 2. Related Work

Word alignment has a strong tradition in translation studies which focus on the detection of translation correspondences in multilingual environments. Guidelines for the annotation of corresponding tokens have been described, e.g., in the Blinker Annotation Project (Melamed 1998) or in Macken (2010a). Relating to this, techniques for the automatic detection of correspondences are investigated in the field of Machine Translation (cf. Och & Ney 2000). Compared to translation tasks of contemporary documents, our project has to cope with a constantly varying degree of similarity between the single texts, making the alignment procedure more complex.

There are several philological projects which deal with parallel historical texts: For instance, the University of Vienna provides an online synopsis of the *Nibelungenlied*[4] and the University of Bern aims at digitizing the complete tradition of the epic *Parzival* by Wolfram von Eschenbach.[5] Another project that will publish the complete tradition of the 12[th] century *Kaiserchronik* started recently at the University of Cambridge.[6] All these projects share the commonality that they concentrate on bringing together

---

[3] With the term *text* we refer to different instances (manuscripts or prints) of one (virtual) base text, the medieval passion treatise *Anselm.*
[4] http://germanistik.univie.ac.at/index.php?id=14531
[5] http://www.parzival.unibe.ch/home.html
[6] http://www.mml.cam.ac.uk/german/staff/kaischron.html

parallel versions in printed or online editions; they do not deal with word alignment.

## 3. Alignment Guidelines

The goal of the alignment task described here is to cover all extant texts of *Anselm* and to align corresponding elements. Thus, in contrast to many edition projects, there is no a priori defined "central" text; instead all texts are of equal importance. Depending on the particular research question, any *Anselm* text can be chosen as the "central" version and basis of comparison. The main focus is on intertextuality and we are interested in similarities and differences between entire documents.

In our guidelines for aligning parallel texts, we distinguish between four separate annotation layers: Cognates (Sec. 2.1.), Synonyms (Sec. 2.2.), Coreference (Sec. 2.3.), and Complex (phrasal) equivalents (Sec. 2.4.) The layers differ in the degree of similarity between the aligned tokens or phrases.[7] Each layer is annotated in a separate pass, i.e. all corresponding cognates are marked first, afterwards all synonyms are marked and so forth. Two tokens (or phrases) are taken to correspond if they share the same context. Ex. (1) shows such a comparable context from two manuscripts Ba1 and D4, where it is possible to align single corresponding tokens.

(1)

Bamberg (Ba1)
*Ain hoher lerer hiesz anshelmus, der pat vnser frauen lange weill vnd zeit wainent vasten vnd peten, Das sy im zu erkennen geb, wie vnser herre gemartert wer word*

'A high teacher was called Anselm, he asked our lady for a long time, crying, starving, praying, to show him how our lord was tortured'

Dessau (D4)
*Sant anszhelmüs / / bischoff hat gebetten lang zeit mit vasten / weinen vnnd betten / Maria die reinen Iuncfrowen vnd müter gots / das si Im wolt volkomenlich offenbaren / das leyden Ires lieben soenes cristi iesu*

'Saint Anselm, the bishop, has asked long time – while starving, crying and praying – Mary, the pure virgin and mother of god, to completely reveal the passion of her son Jesus Christ '

---

[7] Melamed (1998) and Macken (2010a) make similar distinctions, which Macken calls "regular links" (our first two levels) and "fuzzy links" (our levels three and four).

### 3.1. Cognates

On the first level of alignment, all corresponding tokens in two texts that are cognates are aligned. Cognates are words with a common etymological origin.[8] Cognates can belong to different word classes; e.g., nouns can be aligned with verbs as long as both have the same stem or root, see Ex. (2), where the noun *marter* 'martyrdom' is aligned with the participle *gemartert* 'martyred'. Only pairs of the form *token:token* are aligned at this level.

(2) Alignment of cognates[9]

Bamberg (Ba1)

| *[Das]*$_{L1.1}$ | *[sy]*$_{L1.2}$ | *[im]*$_{L1.3}$ | | *zu* | *erkennen* | *geb* | *wie* |
|---|---|---|---|---|---|---|---|
| that | she | him | | to | recognize | gives | how |
| *vnser* | *herre* | *[gemartert]*$_{L1.4}$ | | *wer* | *word* | | |
| our | lord | *tortured* | | has | been | | |

'That she show him how our lord was tortured'

Stuttgart (Stu1)

| *[daz]*$_{L1.1}$ | *[sy]*$_{L1.2}$ | | *[im]*$_{L1.3}$ | *kunt* | *taetty* | *irs* | *aingebornes* |
|---|---|---|---|---|---|---|---|
| that | she | | him | tell | does | her | only |
| *kindes* | *[marter]*$_{E1.4}$ | | | | | | |
| child's | martyrdom | | | | | | |

'That she tell him about her only childs martyry'

### 3.2. Synonyms

In the second annotation pass, *real* synonyms of the type *token:token* are aligned. Two tokens are taken to be synonyms if they are interchangeable without a resulting change in meaning, as in Ex. (3) where *getwagen* is aligned with *gewaeschen*, both meaning 'washed'.[10]

---

[8] This includes suppletive forms (e.g. as in the paradigm of German *sein* or the English equivalent *be*).

[9] Examples are organized as follows: L = level; L1 = level 1; L1.1 = Level 1, alignment pair 1. All tokens belonging to the same aligment pair have the same ID.

[10] In German, certain verbs, called prefix verbs, can occur discontinuously. The verb and its prefix count as one token in these cases, since they belong to the same lemma, listed in standard lexicons.

(3) Alignment of synonyms

Bamberg (Ba1)

| $do_{E1.1}$ | $mein_{L1.2}$ | $kindt_{L1.3}$ | $ir$ | $fuez_{L1.4}$ | **[getwagen]**$_{L2.1}$ | $het_{L1.5}$ |
|------|------|------|------|------|------|------|
| when | my | child | their | feet | washed | has |

*'when my child has washed their feet'*

Stuttgart (Stu1)

| $do_{L1.1}$ | $min_{L1.2}$ | | $kint_{L1.3}$ | $inen$ | $die$ | $fuezz_{L1.4}$ |
|------|------|------|------|------|------|------|
| when | my | | child | them | the | feet |

| $het_{L1.5}$ | **[gewaeschen]**$_{L2.1}$ |
|------|------|
| has | washed |

*'when my child has washed them their feet'*

## 3.3. Coreference

This level aligns phrases that are coreferent, i.e., a pro-form (pronoun or adverb) is used in one text and a corresponding full phrase (NP, PP, VP) in the other. In contrast to the two previous two layers, in this layer entire phrases can be aligned; see Ex. (4) where the full noun phrase *den knecht* 'the servant' is aligned with the pronoun *in* 'him'.

(4) Alignment of coreferential phrases

Stuttgart (Stu1)

| $vnd_{L1.1}$ | $macht_{L1.2}$ | **[den** | **knecht]**$_{L3.1}$ | $wider$ | $gesunt_{L1.3}$ |
|------|------|------|------|------|------|
| and | makes | the | servant | again | healthy |

*'and restored the servant's health again'*

Bamberg (Ba1)

| $vnd_{L1.1}$ | $machet_{L1.2}$ | **[in]**$_{L3.1}$ | $zehant$ | $gesundt_{L1.3}$ |
|------|------|------|------|------|
| and | makes | him | immediately | healthy |

*'and restored the servant's health'*

---

For our purpose, we use the *Duden* as a standard reference (www.duden.de). For extinct word forms (e.g *getwagen* 'washed' in text Ba1 in Ex. (2a)) we use a dictionary for historical German, in our case a dictionary for Middle High German (www.woerterbuchnetz.de/lexer).

### 3.4. Complex (phrasal) equivalents

In a last annotation pass, all remaining non-aligned tokens need to be checked for correspondences. Only entire phrases (NPs, PPs, VPs) can be aligned on this level. In Ex. (5) the individual elements of the verbal phrase *zu erkennen geb* ,to show' have no cognates or synonyms in the parallel text. But the phrase as a whole does have an equivalent in the other text: *kunt taetty* 'tell'. Hence, both phrases are aligned.

(5) Alignment of complex (phrasal) equivalents

Bamberg (Ba1)

| $Das_{L1.1}$ | $sy_{L1.2}$ | $im_{L1.3}$ | $[zu$ | $erkennen$ | $geb]_{L.4.1.}$ | $wie$ |
|---|---|---|---|---|---|---|
| that | she | him | to | recognize | gives | how |
| *vnser* | *herre* | *gemartert*$_{E1.4}$ | *wer* | *word* | | |
| our | lord | tortured | has | been | | |

*'That she signifies him how our lord was tortured'*

Stuttgart (Stu1)

| $daz_{L1.1}$ | $sy_{L1.2}$ | $im_{L1.3}$ | $[kunt$ | $taetty]_{L.4.1.}$ | $irs$ | $aingebornes$ |
|---|---|---|---|---|---|---|
| that | she | him | known | does | her | only |
| *kindes* | *marter*$_{E1.4}$ | | | | | |
| child's | martyrdom | | | | | |

*'That she tells him about her only childs martyrdom'*

## 4. Annotation

To evaluate the guidelines, we performed a test annotation with two annotators. In this section, we describe the annotation scenario, present results from the annotation, and compute the inter-annotator agreement.

### 4.1. Test Annotation

For evaluating the guidelines, we extracted three text fragments with comparable content, from two Anselm texts that are rather similar to each other (Ba1 and Ba2) and from two texts that are rather dissimilar (Ba1 and D4). The fragments were taken from the beginning of the texts and consist of roughly 500 tokens (Ba1: 570; Ba2: 561; D4: 529).

| | Ba1 : Ba2 | | | Ba1 : D4 | | |
|---|---|---|---|---|---|---|
| | 1:1 | 1:n, n:1 | n:m | 1:1 | 1:n, n:1 | n:m |
| 1 (cognates) | 543 | – | – | 170 | – | – |
| 2 (synonyms) | 4 | – | – | 60 | 1 | – |
| 3 (coref.) | – | – | – | – | 6 | – |
| 4 (phrases) | 0 | 4 | 3 | 3 | 20 | 26 |
| All | 547 | 4 | 3 | 233 | 27 | 26 |

Tab. 1: Number of alignments at different layers with two similar (Ba1:Ba2) and two dissimilar (Ba1:D4) fragments

The fragments were annotated independently by two student annotators, who were well acquainted with the Anselm texts in general, but not with the alignment task. They had a short training phase: in a first meeting, they were introduced to the guidelines, next they aligned two short training fragments (of less than 500 tokens), followed by a discussion phase.
We used the annotation tool MMAX2 (Müller & Strube 2006), which we adapted to this task: The fragments to be aligned are placed next to each other in one MMAX window. Alignment links are then inserted using MMAX's facilities for coreference links. Words are displayed in different colors, depending on the type of alignment that they participate in.

### 4.2. Results from the Test Annotation

After the test annotation, the annotators produced an adjudicated gold standard. Table 1 shows the number and types of alignments in the gold corpus. The two similar texts show an extremely high number of correlation and most alignments link cognates, meaning that the fragments even share most of their vocabulary. The two dissimilar Anselm texts behave very differently, sharing fewer links in total and fewer cognates in particular. The vast majority of alignments in both texts are 1:1 alignments.
Fig. 1 confirms these findings. It plots the positions of the aligned tokens of both fragments. Aligning a text with itself would result in a diagonal. The plot on the left, displaying the links between Ba1 and Ba2, indicates that the correlation between both fragments is almost perfect. The plot of the dissimilar texts, Ba1 and D4, still clearly approximates the diagonal, which mirrors the fact that both fragments have the same topic. At the same time, it shows considerable deviations and alignment gaps.

|  | Ba1 : Ba2 | | Ba1 : D4 | |
| --- | --- | --- | --- | --- |
|  | WAA | kappa | WAA | kappa |
| 1 (cognates) | .99 | .83 | .95 | .84 |
| 2 (synonyms) | .99 | .39 | .95 | .63 |
| 3 (coref.) | – | – | .99 | .57 |
| 4 (phrases) | .98 | .58 | .85 | .67 |
| All | 1.00 | .98 | .91 | .76 |

Tab. 2: Inter-annotator agreement (Word Alignment Agreement score and kappa) for all layers and fragments
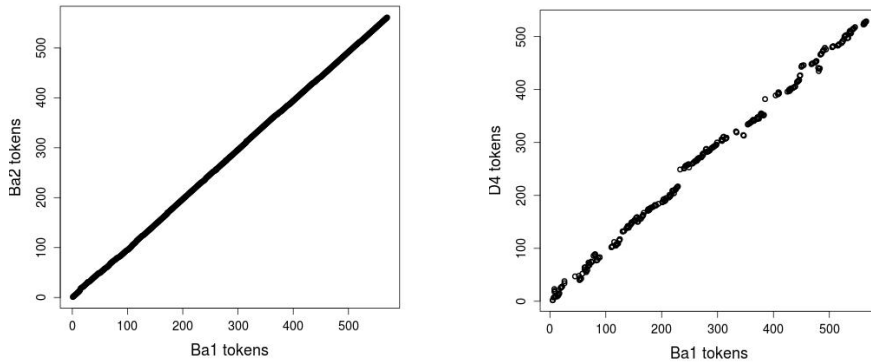


Fig. 1: Plot of the aligned token positions in the similar (Ba1:Ba2, left plot) and dissimilar (Ba1:D4, right plot) fragments.

## 4.3. Inter-annotator agreement

In similar alignment studies, Davis (2002), Daumé & Marcu (2005) and Macken (2010b) evaluated manual word and phrase alignment by computing Word Alignment Agreement (WAA, Davis 2002) and the chance-corrected kappa score (Cohen 1960). WAA/kappa = 1 means perfect agreement, WAA = 0 means no agreement, while kappa = 0 means no agreement *above chance*. The task is considered a classification task: given a pair of <source-word, target-word>, one must decide wether the two words should be aligned or not.

Tab. 2 shows the results for the different layers and fragments. Observed agreement (WAA) is very high for almost all layers. The kappa score, which

"punishes" highly skewed distributions, varies considerably. If we merge all layers into one, we see, however, that the annotator's agreement on the alignments *as such* is nearly perfect (.98) or substantial (.76), while the decision about the type of the alignment (= its layer) is more controversial.

## 5. Case study

The following case study shows how the alignments can be used to anwer philological questions. The treatise *Anselm* has been preserved both in verse and prose; the verse versions are a rather homogeneous corpus while the prose versions could be divided in two groups according to length, i.e. a short and a long form, depending on the word count. Our aim is to go beyond a quantitative classification of each Anselm text to a content-related grouping. For this grouping exercise, it is important to find a way of combining keywords so that they form significant clusters.

Keywords are terms (cognates and synonyms) or phrases (complex phrasal equivalents) selected for their frequency and significance. The aim is to mark up a number of keywords in each of the Anselm texts as a profiling exercise to arrive at the 'fingerprint' of individual texts made up by their specific use of the keywords. Our assumption is that groups and sub-groups can be visualized by using alignments and measuring the distance between the annotated keywords.

In the following we present a case study taken from the opening sequence of the text. We annotated three keywords, namely all references to the persons 'Anselm', 'Mary' and 'Jesus Christ' (for single terms this has been tried manually, cf. Dipper/Schultz-Balluff 2013, Wegera (2014).

The corpus comprises 51 text instances which have preserved the opening sequence i.e. 12 in verse, 20 in (long) prose, 15 in (short) prose, three Dutch copies, and one unclassified fragment. Seventeen texts start straight with the introduction, while the other 34 texts preface it by headings or preliminary remarks; most of them, however, added at a later stage of transmission or edition, e.g. by another writer. It can therefore be assumed that the basic text had no definite title. This makes the specific forms of references to persons within these pre-text sequences especially important as identifying features of the different versions.

Tab. 3 shows the opening clause of texts representing the long prose version (PL), the short prose version (PS) and the verse version (V). The versions differ considerably with regard to references to persons (printed in boldface).

| |
|---|
| PL : B2, fol. 48r,5-17 (Ms. germ. qu. 2025, Staatsbibliothek zu Berlin, Preußischer Kulturbesitz) |
| **Sante anſhelm** der bad *vnſer liebe frauwe von hymelriche* alczü lange zijt mit vil groſzer ynneger begerünge Mit faſten beden vnd mit wachen vnd mit andechtigem gebede vnd mit manichen heiſzen drenen daz ſie yme künd wolde dün <u>yres eyngebornes kindes</u> martele / wie ez von dem anbegynnen da erginge mit zü dem ende fynes lydens |
| PS: M10, fol. 58v,8-59r,2 (Clm 14945, Bayerische Staatsbibliothek München) |
| **Ein hocher lerer hiez anſhelmus** Der pat *vnſer frawn* lange wainent vnd vaſtent Daz ſi im zerkennen gebe wie <u>**vnſer h(er)re**</u> gemartert wer |
| V: O, fol. 1r,1-14 (Cim. I.74, Landesbibliothek Oldenburg) |
| **ANcelmus was ein heilich man** / De hadde langhe dar na ſtan / Dat he gherne hedde weten / Wat <u>**vnſe here**</u> hedde be ſeten / Nv moghe gi horen wu he dede / he was ſtede an ſinem bede / Beide nacht vn̄ dach / An ſiner venigen dat he lach / he ſprak *maria bloygende roſa* / *Lylia vn̄ ſittiloſa* Goddes dure balſ men ſchrin / Lat mir hute dir werden ſchin / Dattu mir moteſt rede ſaghen / van ſinen iām̄erliken plaghen |

Tab. 3: Opening clause in the long prose version, the short prose version, and the verse version

In the fairly homogeneous verse version, there is very little variation in lexis and syntax. The first mentioning of 'Anselm','Mary', and 'Jesus Christ' is consistent: all texts give the attribute ,holy man' to Anselm, compare Mary to flowers such as rose, lily, and perennial, and address Christ as the 'Lord' (Tab. 4). The prose versions use two main clusters in reduced or expanded form around two cores. Beyond this, there is a larger number of texts which show singular combinations which can however be grouped again in clusters (Fig. 2).

The possible relations of groups (core 1 and 2), sub-groups, and single combinations are shown in the following scheme: cluster 1 consists of the combination 'St Anselm' + 'Our Lady from Heaven' + 'only child'; cluster 2 of 'Honorable Teacher Anselm' + 'Our Lady' + 'Lord'.

The illustrations in Fig. 2 and Fig. 3 are meant as a spatial representation of the relation of different texts (referred to by their sigla, i.e. T, Hk etc.)[11] and editions.

---

[11] For a complete list of all German and Dutch texts see http://www.rub.de/schultz-balluff/sanktanselmus.

| Keywords (person) | Forms of reference |
|---|---|
| Anselm | *Heiliger Mann Anselm* ‚holy man Anselm' |
| Mary | *Maria, Rose, Lilie und Zeitlose* ‚Mary, rose, lily, and perennial |
| Jesus Christ | *Herr* ‚Lord' |

Tab. 4: Keywords and forms of references in the verse version

Eight to nine texts can be linked to these two fixed clusters. These can be extended, reduced, and combined, i.e. 'Our Lady from Heaven' can be extended to 'Our Dear Lady from Heaven' or the two clusters can be combined to form 'Teacher Anselm' + 'Our Dear Lady' + 'Dear Child' (Fig. 2). We can assume that these clusters developed first and that the other combinations are variations of them.



**Core 1**
(T, Hk, M4, Ka, Sb, M2, W, St2, M)

St. Anselm
+
unsere Frau vom Himmelreich
+
einziges Kind

(extended)
unsere **liebe** Frau vom Himmelreich
(sa, B2, Stu)
↓
(varied)
**liebes** Kind
(M3)

**combination of core 1 and 2**

(extended)
hoher Lehrer Anselm +
unsere Frau die Mutter Gottes
+ einziges Kind
(We)

(reduced and extended)
Lehrer Anselm +
unsere liebe Frau +
liebes Kind
(N3)

(or reduced core 1)
St. Anselm +
**unsere Frau** +
einziges Kind
(Be, N4, St)

(extended)
unsere **liebe** Frau
(H)
↓
(extended)
**liebes** einziges Kind
(SG)

**Core 2**
(Ba, Ba2, M5, M7, M9, M10, Me, N2)

hoher Lehrer Anselm
+
unsere Frau
+
Herr

(extended)
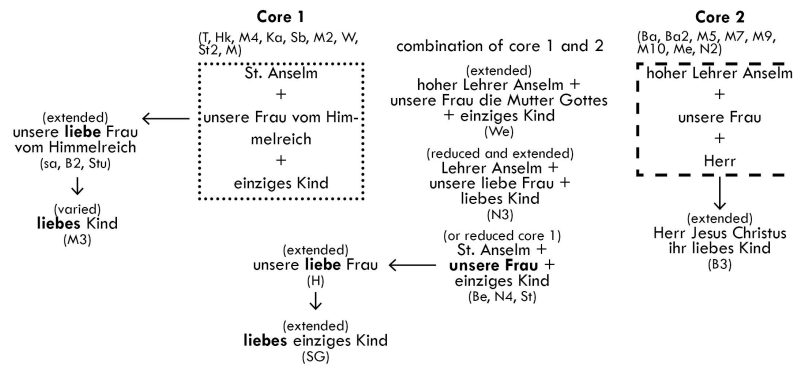Herr Jesus Christus
ihr liebes Kind
(B3)

Fig. 2: Main clusters in the prose version

The second scheme (Fig. 3) shows the correlation of disparate combinations and the identification of abstract cores. The ten Anselm texts show combinations of basic elements which suggest that they also are connected. It is possible to generate abstract clusters from the single elements which form virtual hubs, giving the ten text instances a clear position within the overall system.

**Core 3**

Anselm + Maria reine Magd + Herr (D3)

Anselm + Maria + Herr

St. Anselm + Maria + Herr Jesus Christus (B)

Bischof St. Anselm + Maria reine Jungfrau Mutter Gottes + lieber Sohn Jesus Christus (D4)

**Core 4**

Anselm + Mutter Gottes + Jesus Christus

heiliger Mann Anselm + Mutter Gottes + Herr Jesus Christus (Le)

heiliger Mann Anselm + Maria Mutter Gottes + Herr Jesus Christus (Am)

heiliger Anselm + selige Magd + Sohn (M6)

Anselm + unsere Frau + lieber süßer Herr Gott (Wo)

**Core 5**

Anselm + unsere Frau + Herr

heiliger Bischof St. Anselm + unsere liebe Frau + lieber Herr Jesus Christus (M8)

heiliger Mann St. Anselm + Maria Mutter Gottes + lieber Herr Jesus Christus (B1523)

heiliger Herr Anselm + unsere Frau + unser Herr ihr liebes Kind (sl)
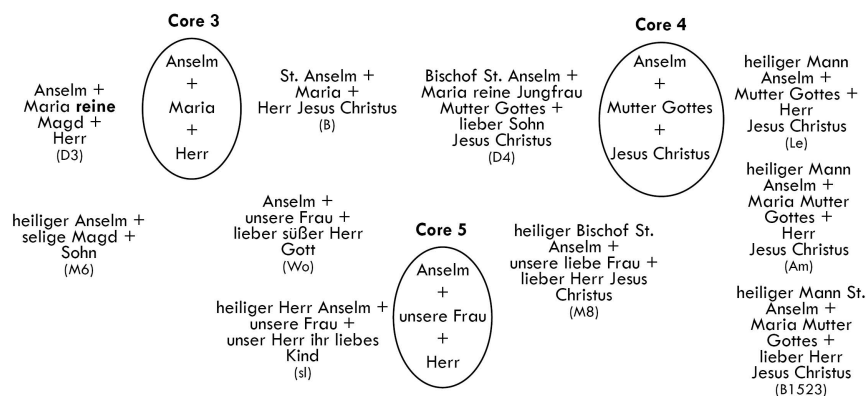
Fig. 3: Sub-groups and single combinations in the prose version

Against this background, the preliminary distinction between long and short prose texts has to be reconsidered. With the help of keywords which will be applied to each text instance in its entirety, new constellations will become visible. This will allow us to take a fresh look at text production and changes during the transmission process. It will hopefully form a solid textual basis to build on for further insights into textual criticism, literary trends, and cultural change – ultimately resulting in a stable framework within which to place the shape-shifting sets of questions with which St Anselm confronts Mary.

## 6. Conclusion

In this paper we introduced a method for aligning parallel historical texts. The annotation guidelines focus on the problems aligning texts which differ in their degree of similarity. An evaluation annotation revealed considerable inter-annotator agreement, even when the two aligned texts were very dissimilar. Furthermore, we showed how the alignment can be of use for clustering parallel texts according to their use of specific vocabulary terms.

**References**

Cohen, Jacob (1960): A coefficient of agreement for nominal scales. In: Educational and Psychological Measurement 20: 37–46.

Daumé III, Hal/Marcu, Daniel (2005): Induction of word and phrase alignments for automatic document summarization. In: Computational Linguistics 31(4): 505–530.

Davis, Paul C. (2002): Stone Soup Translation: The LinkedAutomata Model. Ph.d., Ohio State University.

Dipper, Stefanie/Schultz-Balluff, Simone (2013): The Anselm Corpus: Methods and Perspectives of a Parallel Aligned Corpus. In: Proceedings of the Workshop on computational historical linguistics at NODALIDA 2013. Oslo, S. 27–42 (NEALT Proceedings Series 18).

Macken, Lieve (2010a): Annotation Guidelines for Dutch-English Word Alignment. Version 1.0. Technical report, Language and Translation Technology Team, Faculty of Translation Studies. University College Ghent.

Macken, Lieve (2010b): An Annotation Scheme and Gold Standard for Dutch-English Word Alignment. In: Proceedings of the Seventh International Conference on Language Resources and Evaluation (LREC'10).

Melamed, I. Dan (1998): Annotation style guide for the Blinker project, Version 1.0.4. IRCS Technical Report #98-06. University of Pennsylvania, Philadelphia.

Müller, Christoph/Strube, Michael (2006): Multi-Level Annotation of Linguistic Data with MMAX2. In: Braun, Sabine/Kohn, Kurt/Mukherjee, Joybrato (eds.): Corpus Technology and Language Pedagogy. New Resources, New Tools, New Methods. Frankfurt: Peter Lang, 197–214.

Och, Franz Josef/Ney, Hermann (2000): A comparison of aligment models for statistical machine translation. In: Proceedings of the 18[th] International Conference on Computational Linguistics (COLING-ACL 2000).

Wegera, Klaus-Peter (2014): "Interrogatio St. Anselmi de Passione Domini, deutsch". Überlieferung – Edition – Perspektiven der Auswertung. Hrsg. von der Nordrhein-Westfälischen Akademie der Wissenschaften und Künste. Paderborn.