

Contents

Abstract	i
Acknowledgements	i
1 Introduction	1
1.1 Motivation	1
1.2 Research Question and Goals	2
1.2.1 Research Questions	2
1.2.2 Goals	3
1.3 Structure	3
1.4 GitHub repository	3
2 Romansh	4
2.1 Rhaeto-Romance	4
2.2 Romansh	5
2.3 Rumantsch Grischun	6
2.3.1 Lia Rumantscha	6
2.3.2 Rumantsch Grischun	6
2.3.3 Properties	7
2.3.4 Today	7
3 Compiling the Corpus	9
3.1 Introduction	9
3.2 Collecting the Data	9
3.3 Web Scraping	10
3.4 Building the Corpus	11
3.4.1 HTML Parsing	11
3.4.2 Document Alignment	12
3.5 SQLite database	16
3.6 Summary	17
3.6.1 Statistics	17

4	Sentence Alignment	21
4.1	Introduction	21
4.1.1	Formal definition	21
4.2	Method Overview	22
4.2.1	Length Based	22
4.2.2	Partial Similarity Based	22
4.2.3	Translation based	23
4.2.4	Hybrid models	23
4.2.5	Summary	24
4.3	More Recent methods	24
4.3.1	Bleualign	24
4.3.2	Vecalign	25
4.4	Sentence alignment pipeline	26
4.4.1	Tool of choice	26
4.4.2	Pipeline	27
4.4.3	Sentence segmentation	27
4.4.4	Aligning Language Pairs	28
4.4.5	Filtering and Tokenizing	29
4.5	Results	29
5	Word Alignment	31
5.1	Introduction	31
5.2	Overview of Methods	32
5.2.1	IBM Model 1	32
5.2.2	Higher IBM Models	34
5.3	Word Embeddings	35
5.3.1	Excursion: Words	35
5.3.2	Word Embeddings	36
5.3.3	Word Similarity	37
5.3.4	Multilingual Word Embeddings	38
5.3.5	Summary	38
5.4	Similarity Based Word Alignment	38
5.4.1	Method	38
5.4.2	Summary	40
6	Gold standard	41
6.1	Introduction	41
6.2	Sure and Possible Alignments	42
6.3	Gold standard for German-Romansh	42

6.3.1	Annotation tool	43
6.3.2	Guidelines	43
6.3.3	General principles	43
6.3.4	Examples	44
6.4	Flaws	47
7	Results	49
7.1	Evaluation Metrics	49
7.2	Baseline Systems	50
7.2.1	fast_align	50
7.2.2	eflomal	51
7.2.3	Performance	51
7.3	SimAlign	51
7.3.1	Performance	52
7.4	Discussion	52
7.4.1	General Problems with Evaluation	53
7.5	Summary	55
8	Concluding Words	57
8.1	Goals	57
8.2	Corpus Compiliation	57
8.3	Gold Standard	58
8.4	Evaluation	58
8.5	Future	58
	List of Tables	60
	List of Figures	61
	List of Listings	62
	Bibliography	63
A	Algnment Examples	69

Chapter 6

Gold standard

6.1 Introduction

In the previous chapter, I discussed SimAlign, a method for computing word alignments based on measuring the similarity between multilingual word embeddings. The clear advantage of this method is that it does not rely on the existence of large amounts of parallel data. The multilingual word embeddings can be learned from monolingual data. Jalili Sabet et al. 2020 evaluated their method on language pairs which were all part of the training data for the language models in use (mBERT and XLM-R). In the course of this work, I was able to extract 79,548 sentence pairs for German-Romansh. I shall now proceed to test how well SimAlign performs on this language pair, considering the fact the Romansh is not part of the training data for these models, i.e., it is an unseen language.

In order to measure the quality of words alignments, a model's performance is measured on a test set, dubbed gold standard, which is created by human annotators. For the gold standard to be of good quality and consistent with itself, annotators have to follow strict guidelines. These guidelines address issues of ambiguity in word alignments. (Koehn 2009, p. 115).

Some problematic cases that might occur are function words¹ that have no clear equivalent in the other language. Koehn 2009 gives as an example the German-English sentence pair: *John wohnt hier nicht* and *John does not live here*. What German word should the English word *does* be aligned to? Three different choices can be made:

1. The word should remain unaligned since it has no clear equivalent in German.
2. The word *does* is connected with *live*; it holds information about number (singular) and tense (present tense), which, in German, is contained in one word *wohnt*, so it should be aligned to *wohnt*, together with *live*.

¹Function words form a closed class of words (a fixed set of words with virtually no new additions), they occur frequently and often have structuring uses in grammar. Pronouns, prepositions and conjunctions like *of*, *it*, *and*, or *you* are function words (Jurafsky and Martin 2019, p. 144).

3. *does* is part of the negation; without it, the sentence would not contain this word. Therefore, *does* should be aligned with *nicht* (the German negation).

There are several possibilities, all of them arguable, none of them plain wrong, which illustrates the need for clear guidelines.

6.2 Sure and Possible Alignments

An approach for solving problematic cases is the distinction between “Sure” and “Possible” alignments (Och and Ney 2000), which are also sometimes referred as “fuzzy alignments” (Clematide et al. 2018). Generally, these labels allow to distinguish between ambiguous and unambiguous links. Ambiguous links are labeled Possible and unambiguous links are labeled Sure (Lambert et al. 2005). The Possible label was conceived to be used especially for aligning words within idiomatic expressions, free translations and missing function words (Och and Ney 2000). This distinction also has an impact on the way the evaluation metrics are computed (see Section 7.1).

There seems to be no clear global definition about which alignments should be considered as unambiguous, thus marked as Sure, and which should be considered ambiguous, thus marked as Possible. For some created gold standards, no distinction between Sure and Possible alignments was made at all (Clematide et al. 2018). In another case, annotators were asked to first label all alignments as Sure and then refine their alignments with confidence labels (Holmqvist and Ahrenberg 2011). In the creation of the English-Icelandic gold standard in Steingrímsson, Loftsson, and Way 2021, annotators used only Sure links. Their annotations were then combined; all 1-to-1 alignments both annotators agreed upon (i.e., the intersection of their annotations) were marked as Sure and all other alignments were marked as Possible (Steingrímsson, Loftsson, and Way 2021). Different annotation schemes use Sure and Possible alignment in a different way.

6.3 Gold standard for German-Romansh

As explained before, in order to measure the performance of both models, the similarity and embedding based model (SimAlign) and the statistical models (fast_align and eflo-mal), on the language pair German-Romansh a gold standard is needed. Since no such gold standard exists, I took upon myself to create one. Although I am not a speaker of Romansh, my experience as a trained linguist, as well as my knowledge in related languages (Latin, Italian, French), allows me to confidently tackle this task. Additionally, whenever I was in doubt, I referred to the online dictionary Pledari Grond², which also offers a grammar overview.

²<https://www.pledarigrond.ch/rumantschgrischun>

6.3.1 Annotation tool

I used the tool *AlignMan* which was originally programmed for creating the gold standard for English-Icelandic (Steingrímsson, Loftsson, and Way 2021). It is quite easy to use and its code is readable. I also had to make some small changes to the code. For instance, the sentences to be aligned, while loaded into the database, were read in opposite order, such that the source language became the target language and vice versa. I fixed this issue, so that source (German) and target (Romansh) languages stay the same accross all applications.

As mentioned in Section 6.2, the annotation scheme used by Steingrímsson, Loftsson, and Way 2021 does not allow labeling of links with Sure and Possible. Instead, *AlignMan* treats the union of 1-to-1 alignments made by two annotators as Sure alignments and all other alignments as Possible. This means, each annotator is expected to only annotate Sure alignments. This also applied to me while annotating the German-Romansh gold standard.

6.3.2 Guidelines

As mentioned before, clear guidelines need to be defined for creating the gold standard in order to ensure quality and consistency. I shall now proceed to describe the guidelines I used for my annotation of the word alignments for the gold standard.

A motto cited often for annotating word alignments is “Align as small segments as possible, and as long segments as necessary” (Vronis and Langlais 2000, cited in Ahrenberg 2007). A variation of this is found in Clematide et al. 2018: “As few words as possible and as many words as necessary that carry the same meaning should be aligned,” referring to Lambert et al. 2005.

In the following sections I will list some general principles as well as more specific principles involving German and Romansh.

6.3.3 General principles

Principle I. Use only Sure alignments: Since the annotating tool I was using does not provide the use of confidence labels (cf. Section 6.3.1: Annotation tool), I only aligned words which would be considered Sure alignments, i.e., they are unambiguous (cf. Section 6.2: Sure and Possible Alignments).

Principle II. Prefer 1-to-1 alignments over 1-to-n alignments or n-to-n alignments: Since all alignments are seen as Sure alignments, 1-to-n alignments should be avoided, unless a single word in the source sentence lexically corresponds to several words in the target sentence. This means alignments of phrases should be avoided. This is also due to the fact that we are testing models for automatic word alignment, and not phrase alignment.

Words that are repeated in one language, but not in the other, should only be linked once, leaving the repetition unaligned.

Principle III. Lexical alignments should always be preferred over all other alignments (part of speech (POS) alignments or morphosyntactical alignments). This means alignments should describe first and foremost lexical correspondences, i.e., both words have the same lexical meaning (but not necessarily share the same grammatical function or the same POS). Only words that are translations of each other also outside of the specific context of the sentence pair at hand should be aligned. This is in line with Clematide et al. 2018. In cases of paraphrasing during translations, words should remain unaligned.

6.3.4 Examples

I will now give some examples to illustrate the above principles.

Compound words Compounding is the formation of new lexemes by adjoining two or more lexemes (Bauer 1988). In German, compounds are productive and prominent means of word formation in German (Clematide et al. 2018). In a sample of 4,500 types examined by Clematide et al. 2018, 80% of German nouns were compounds. Romansh, in comparison, uses prepositions (usually *da*) for linking nouns, with one noun modifying the other (Tschärner and Denoth 2022). Other prepositions that can be found for linking words are *cunter* and *per*.³ In other cases, German compounds might be translated to Romansh using an adjective + noun, e.g., German *Gastkanton* was translated to *chantun ospitant* “hosting canton”. See table 6.1 for examples.

German compounds will be aligned to their equivalent lexical words, but not to function words, resulting in a 1-n alignment: *Webseite ~pagina [d'] internet, Gebäudeversicherung ~Assicuranza [d'] edifizis*. This is also inline with principles I, II and III in Clematide et al. 2018.

German preterite vs. Romansh perfect

In the corpus at hand, two tenses are used in German for referring to past events: the preterite and the perfect. The German preterite is a synthetic verb form, i.e., it is made up of a single conjugated form. Some examples are *nahm* (infinitive *nehmen* “take”) or *wurde* (infinitive *werden* “become”). The German perfect is an analytic construction made up of an auxiliary verb (*haben* “have” or *sein* “be”) and the past participle, e.g., *Die Präsidentenkonferenz hat nun entschieden* “The conference has decided”.

³Typologically, this is inline with other Romance languages such as French, which uses prepositions (*de*, *en* and *à*) for linking two nouns, e.g., *une robe de soie* “a silk dress” (Price 2008)[510].

German	Romansh	
<i>Beratungsstelle</i>	<i>post da cussegliaziun</i>	“consultation point”
<i>Gebäudeversicherung</i>	<i>Assicuranza d’edifizis</i>	“building insurance”
<i>Webseite</i>	<i>pagina d’internet</i>	“web site”
<i>Kindermasken</i>	<i>mascrinas per uffants</i>	“children masks”
<i>Brandversicherung</i>	<i>assicuranza cunter feu</i>	“fire insurance”
<i>Gastkanton</i>	<i>chantun ospitant</i>	“hosting canton”

Table 6.1: Translation examples of German compounds into Romansh

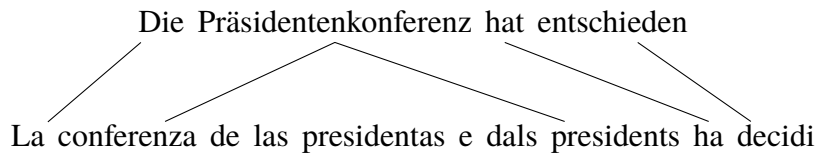


Figure 6.1: Aligning German perfect to Romansh perfect

In contrast to German, Romansh only has one tense referring to past events: the perfect. It is an analytic construction made, in a similar fashion as in German, of an auxiliary *habere* “have” for transitive verbs or *esse* “be” for intransitive verbs and the past participle (Bossong 1998, p. 189). The German sentence given above (*Die Präsidentenkonferenz hat nun entschieden*) was translated as *La conferenza da las presidentas e dals presidents ha usse decidi*. *ha* is the auxiliary and *decidi* is the past participle. This poses no real problem since we can link the German auxiliary to the Romansh auxiliary and the German participle to the Romansh participle.

However, a German preterite is always translated using the Romansh perfect. For example, in the sentence *Der Kanton Graubünden war letztmals 2003 Gastkanton* “The last time the Canton of Grisons was a host canton was in 2003” the verb *war* “was” is translated as *è stà*. This theoretically results in a 1-2 link. However, since the verb *è* here only carries grammatical information of tense and number, but no real lexical information, it should remain unaligned.

The German perfect should be aligned to the Romansh perfect using a 1-1 alignment; auxiliary to auxiliary and participle to participle. The German preterite should also be aligned using a 1-1 alignment to the Romansh participle, leaving the auxiliary unaligned and avoiding a 1-2 alignment.

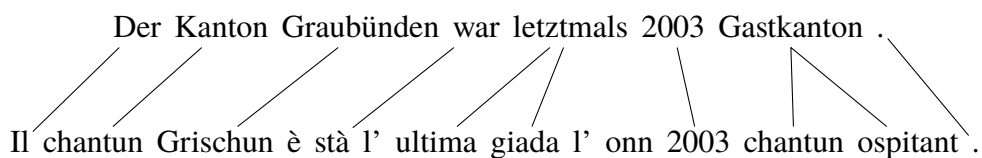


Figure 6.2: Alignment of German preterite to Romansh perfect

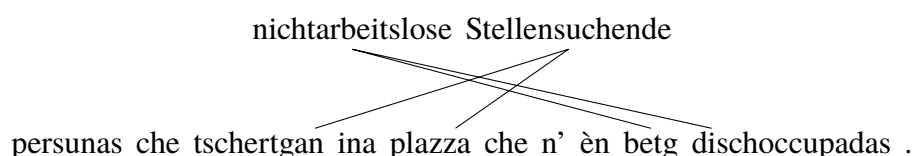


Figure 6.3: Aligning German present participles to Romansh relative clauses

German present participle

German present participles (known in German as *Partizip I*) are translated to Romansh using relative clauses. Moreover, adjectives (and participles in the function of adjectives), can be nominalized, meaning they become the head of a noun phrase and there is no need for an actual noun. A good example for that in the corpus is the German noun phrase *nichtarbeitslose Stellensuchende* (cf. ex. 1), which was translated as a noun phrase with a relative clause: *persunas che tschertgan ine plazza che n'èn betg dischoccupadas* “persons who look for a job who are not unemployed”.

- (1) nicht-arbeit-s-los-e Stellen-such-end-e
 not-work-GEN-less-PL job-search-PRES.PART-PL
 “People looking for jobs who are not unemployed”

In this case, these two phrases should not be aligned as phrases, but only the content words which lexically correspond to each other: *nichtarbeitslose* ~ *betg dischoccupadas*; *Stellensuchende* ~ *tschertgan [ina] plazza*.

Double negation

Negation in Romansh is built using two particles: *na* and *betg* to negate verbs or *nagin-* to negate nouns. Since we prefer 1-1 alignments, the German negations *nicht* (for verbs) and *kein-* for nouns should be aligned only to the second Romansh particle (*betg/nagin-*), leaving Romansh *na* unaligned. Granted, this is also in favor of the SimAlign output, but it is also linguistically motivated: when negating the imperative form, *na* can be omitted required TODO:cite Grammatica per l'instrucziun dal rumantsch grischun.

Articles and prepositions

German articles inflect in case, which expresses some syntactic relations between nouns. Romansh often uses prepositions for expressing the same relations. For instance *Zustimmung der Person* “the person’s agreement” is translated as *consentiment da la persuna*. I align the German article *der* with Romansh *da*, leaving *la* unaligned. Except for my preference for 1-1 alignments, the motivation for this is that it is the preposition *da* that expresses the genitival relations between the nouns.

Separable verbs

German uses many verbs to which an adverb or a preposition is affixed in order to delimit the verb's meaning (or sometimes completely change its meaning). In such cases, both the verb and its affix should be aligned to the corresponding Romansh verb, resulting in a 2-1 alignment.

6.4 Flaws

I shall now discuss the quality of my gold standard and some flaws it has.

The most obvious flaw is the fact that I created the gold standard alone. With more than one annotator, more intricate annotating schemes can be used in order to ensure higher quality, consistency and harmony. For instance the annotators' agreement can be measured using the so-called inner-annotator agreement (Holmqvist and Ahrenberg 2011). Further, the intersection of the annotators' *Sure* alignment can be used to build the final *Sure* alignments set and the reunion of the *Possible* alignments can be used to create the final *Possible* alignments set Mihalcea and Pedersen 2003. A third annotator can also revise and resolve conflicts between two annotators Mihalcea and Pedersen 2003. When several annotators work on the same task, they can also discuss conflicts and resolve them using a majority vote (Melamed 1998).

All of these possible schemes cannot be realized in my case.

Another flaw is the missing confidence labels (*Sure* and *Possible*), which may influence the evaluation scores. Doing without *Possible* links and using only *Sure* links is however precededented (Clematide et al. 2018; Mihalcea and Pedersen 2003) and hence defensible.

In order to test my own consistency, I have re-annotated the first 100 sentences in the sample. TODO: results

Despite of the flaws mentioned, I am certain that gold standard is of high quality and consistency, due to the fact that I was also the one to define the guidelines.

Glossary

Graubünden The Canton of Grisons. 1, 5, 9

Standeskanzlei State Chancellery of Grisons. 9

Acronyms

AER average error rate. 40, 49, 52, 53, 54, 62

EM expectation-maximization. 33

NER named entity recognition. 2

NMT neural machine translation. 26, 29

POS part of speech. 2, 44

SMT statistical machine translation. 25

List of Tables

2.1	Examples for choosing the forms for Rumanstch Grischun, based on liver1999	8
3.1	Description of the table corpus in corpus.db	16
3.2	Description of the table raw in corpus.db	16
3.3	Number of parallel documents per year, as of July 20, 2022.	18
3.4	Number of documents per language and year as of 20 July, 2022.	19
3.5	Twenty most frequent tokens in each language in the corpus.	20
4.1	Parallel corpus in numbers	30
6.1	Translation examples of German compounds into Romansh	45
7.1	Evaluation metrics for word alignments with the baseline models for different dataset sizes. “Dataset Size” refers to the number of sentence pairs.	51
7.2	Evaluation metrics for word alignments using SimAlign, with different embeddings and word/sub-word level. Best result per embedding type in bold.	52
7.3	Comparison of the best performance of each of the three methods. The best value in each column is in bold.	53

List of Figures

2.1	Distribution of Rhaeto-Romance, taken from haiman1992	5
3.1	Directory tree of corpus_builder	10
3.2	Directory scheme for saving the HTML files	11
3.3	Portion of automatically aligned press releases up to 2009	13
3.4	Corpus creation pipeline	17
4.1	Sentence alignment pipeline	27
5.1	Word alignment example	31
5.2	Similarity matrix	39
5.3	Alignment matrix	39
5.4	The resulting word alignment	39
6.1	Aligning German perfect to Romansh perfect	45
6.2	Alignment of German preterite to Romansh perfect	45
6.3	Aligning German present participles to Romansh relative clauses	46
7.1	Comparing precision between the systems for different dataset sizes.	53
7.2	Comparing recall between the systems for different dataset sizes.	54
7.3	Comparing AER between the systems for different dataset sizes.	54
A.1	Word alignment example 1	69
A.2	Word alignment example 2	70
A.3	Word alignment example 3	70
A.4	Word alignment example 4	71
A.5	Word alignment example 5	71

List of Listings

- 3.1 Example for a JSON file containing the press releases extracted from the HTML files. 11
- 3.2 Example for a JSON file containing aligned documents 14
- 4.1 An excerpt from a file containing sentences for alignment. 28
- 4.2 An excerpt from the file containing sentence pairs in German–Romansh . 30

Bibliography

- Ahrenberg, Lars (2007). *LinES 1.0 Annotation: Format, Contents and Guidelines*. Tech. rep.
- Bauer, Laurie (1988). *Introducing Linguistic Morphology*. Edinburgh University Press.
- Bosson, Georg (1998). *Die Romanischen Sprachen: Eine vergleichende Einführung*. Hamburg: Helmut Buske Verlag.
- Clematide, Simon et al. (2018). “A multilingual gold standard for translation spotting of German compounds and their corresponding multiword units in English, French, Italian and Spanish”. In: *Multiword Units in Machine Translation and Translation Technology*. Ed. by Ruslan Mitkov et al. John Benjamins, pp. 125–145. DOI: <https://doi.org/10.1075/cilt.341>.
- Holmqvist, Maria and Lars Ahrenberg (May 2011). “A Gold Standard for English-Swedish Word Alignment”. In: *Proceedings of the 18th Nordic Conference of Computational Linguistics (NODALIDA 2011)*. Riga, Latvia: Northern European Association for Language Technology (NEALT), pp. 106–113. URL: <https://aclanthology.org/W11-4615>.
- Jalili Sabet, Masoud et al. (Nov. 2020). “SimAlign: High Quality Word Alignments Without Parallel Training Data Using Static and Contextualized Embeddings”. In: *Findings of the Association for Computational Linguistics: EMNLP 2020*. Online: Association for Computational Linguistics, pp. 1627–1643. DOI: [10.18653/v1/2020.findings-emnlp.147](https://doi.org/10.18653/v1/2020.findings-emnlp.147). URL: <https://aclanthology.org/2020.findings-emnlp.147>.
- Jurafsky, Daniel and James H. Martin (2019). *Speech and Language Processing. An Introduction to Natural Language Processing, Computational Linguistics, and Speech Recognition*. Third Edition Draft. URL: <https://web.stanford.edu/~jurafsky/slp3/>.
- Koehn, Philipp (2009). *Statistical Machine Translation*. Cambridge University Press.
- Lambert, Patrik et al. (2005). “Guidelines for Word Alignment Evaluation and Manual Alignment”. In: *Language Resource and Evaluation* 39, pp. 267–285. DOI: [10.1007/s10579-005-4822-5](https://doi.org/10.1007/s10579-005-4822-5).
- Melamed, I. Dan (1998). “Annotation Style Guide for the Blinker Project”. In: *CoRR* cmp-lg/9805004. URL: <http://arxiv.org/abs/cmp-lg/9805004>.

- Mihalcea, Rada and Ted Pedersen (2003). “An Evaluation Exercise for Word Alignment”. In: *Proceedings of the HLT-NAACL 2003 Workshop on Building and Using Parallel Texts: Data Driven Machine Translation and Beyond*, pp. 1–10. URL: <https://aclanthology.org/W03-0301>.
- Och, Franz Josef and Hermann Ney (Oct. 2000). “Improved Statistical Alignment Models”. In: *Proceedings of the 38th Annual Meeting of the Association for Computational Linguistics*. Hong Kong: Association for Computational Linguistics, pp. 440–447. DOI: 10.3115/1075218.1075274. URL: <https://aclanthology.org/P00-1056>.
- Price, Glanville (2008). *A Comprehensive French Grammar*. Blackwell Publishing.
- Steingrímsson, Steinþór, Hrafn Loftsson, and Andy Way (2021). “CombAlign: a Tool for Obtaining High-Quality Word Alignments”. In: *Proceedings of the 23rd Nordic Conference on Computational Linguistics (NoDaLiDa)*. Reykjavik, Iceland (Online): Linköping University Electronic Press, Sweden, pp. 64–73. URL: <https://aclanthology.org/2021.nodalida-main.7>.
- Tscharner, Gion and Duri Denoth (2022). *Grammatikteil des Vallader / Grammatica valladar*. URL: http://www.udg.ch/dicziunari/files/grammatica_vallader.pdf (visited on 06/07/2022).
- Vronis, Jean and Philippe Langlais (2000). *Evaluation of parallel text alignment systems - The ARCADE project*.