

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	2
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	6
2.3.4 Today	7
3 Compiling the Corpus	8
3.1 Introduction	8
3.2 Collecting the Data	8
3.3 Web Scraping	8
3.4 Building the Corpus	10
3.4.1 HTML Parsing	10
3.4.2 Document Alignment	11
3.5 Manual alignment of unlinked documents	14
3.6 SQLite database	14
3.7 Summary	14
3.7.1 Statistics	15

4	Sentence Alignment	16
4.1	Introduction	16
4.1.1	Formal definition	16
4.2	Method Overview	17
4.2.1	Length Based	17
4.2.2	Partial Similarity Based	17
4.2.3	Translation based	18
4.2.4	Hybrid models	18
4.2.5	Summary	18
4.3	More Recent methods	19
4.3.1	Bleualign	19
4.3.2	Vecalign	20
4.4	Sentence alignment pipeline	20
4.4.1	Tool of choice	20
4.4.2	Pipeline	21
4.4.3	Sentence segmentation	21
4.4.4	Aligning language pairs	22
4.4.5	Filtering and tokenizing	23
4.5	Results	23
5	Word Alignment	25
5.1	Introduction	25
5.2	Overview of Methods	26
5.2.1	IBM Model 1	26
5.2.2	Higher IBM Models	27
5.3	Word Embeddings	28
5.3.1	Excursion: Words	28
5.3.2	Word Embeddings	29
5.3.3	Word Similarity	30
5.3.4	Multilingual Word Embeddings	30
5.3.5	Summary	31
5.4	Similarity Based Word Alignment	31
5.4.1	Method	31
5.4.2	Summary	33
6	Gold standard	34
6.1	Introduction	34
6.2	Sure and Possible Alignments	34
6.3	Evaluation Metrics	35
6.4	Gold standard for German-Romansh	35
6.4.1	Annotation tool	35
6.4.2	Guidelines	36
6.4.3	General principles	36

6.4.4	Examples	37
6.5	Flaws	39
7	Results	41
7.1	Evaluation Metrics	41
7.2	Baseline Systems	42
7.2.1	fast_align	42
7.2.2	eflomal	42
7.2.3	Performance	42
7.3	SimAlign	43
7.3.1	Performance	43
7.4	Discussion	43
7.4.1	General Problems with Evaluation	44
7.5	Summary	46
8	Concluding Words	48
8.1	Goals	48
8.2	Corpus Compilation	48
8.3	Gold Standard	48
8.4	Evaluation	49
8.5	Future	49
	List of Tables	51
	List of Figures	52
	Bibliography	53
A	Algnment Examples	58

Chapter 3

Compiling the Corpus

3.1 Introduction

The corpus at hand incorporates the press releases published by the Canton of Grisons/Graubünden. These press releases are a means of the cantonal government to publish news and information about topics such as politics, economy, health and culture. Graubünden, which is made up of German speaking, Italian speaking and Romansh speaking regions, is the only trilingual canton in Switzerland. As such, virtually all press releases are published in German, Italian and Romansh. This trilingual setting lends itself to be collected to a parallel trilingual corpus.

3.2 Collecting the Data

At first, I contacted the *Standeskanzlei* (“State Chancellery of Grisons”) which is the “the general administrative authority for questions of office, coordination and liaison with the cantonal parliament (‘Grosser Rat’), government and cantonal administration” (**staka**). The *Standeskanzlei*, with its *Übersetzungsdienst* (“Translation service”), is responsible for translating documents in service of the canton. I was hoping to receive the data directly from them – after all, we are not talking about private or commercial data, but about public translation work financed with taxpayers’ money.

I spoke to Mr. Mirco Frepp from the communication services (*Kommunikationsdienst*), which, although very friendly, had to inform me that it would be impossible for me to receive the data. The explanation was that the documents are not saved locally somewhere, but are saved in a database. The documents are extracted from the database and are generated as ad-hoc HTML documents whenever the website is accessed. It was also not possible to receive a dump of the database.

3.3 Web Scraping

Not being able to receive a dump of the database meant I had to scrape the canton’s website, extract the relevant content from the HTML files and construct my own database. In order to achieve this, I wrote a series of Python scripts that would take care of these tasks. All the scripts can be found on my GitHub/Gitlab repository. The scripts relevant for the database building are saved under the

```

corpus_builder
├── corpus_builder
│   ├── access_db.py
│   └── create_corpus.py
└── web_scraper
    ├── test_web_scraper.py
    └── web_scraper.py

```

Figure 3.1: Directory tree of corpus_builder

```

html
├── 1997
│   ├── 1997_12924_DE.html
│   ├── 1997_12936_IT.html
│   └── ...
├── 1998
├── 1999
├── ...
└── 2022
    ├── 2022_2022010301_DE.html
    ├── 2022_2022010301_IT.html
    ├── 2022_2022010301_RM.html
    └── ...

```

Figure 3.2: Directory scheme for saving the HTML files

folder corpus_builder.

Web Scraper

The script web_scraper.py goes to the index web page for each year and language. This page contains a links pointing to all the press releases that were released that year. It collects all those links and then downloads the HTML from each link. The HTML pages are saved to separate folders for each year. The filenames have the following format: year_file-id_language, e.g., 1997_12924_DE.html. The file-id is taken from the URL and will be later used to align the documents.

Since the script makes many requests to the website, one has to expect that server might stop responding, which will result in a request time-out. To avoid downloading HTML pages that were already downloaded, the script will skip and press release that already exists locally, providing the file size is greater than 0 bytes. This way, the script can be run at a later stage, after new press releases were published, to complete the local repository. To make sure the local copy of the press releases is complete, the script can be simply run until a message is printed that no new press releases were downloaded.

By default, the script will download the press releases for the entire year range (1997 to the

current year) and in all three languages. This can be limited by using the following optional arguments:

- `--year` – limit the scraping to a year or to a range of years separated by a comma, e.g., `--year 2022` or `--year 2020,2022`
- `--lang` – limit the scraping to one or more languages (comma separated), e.g., `--lang de,it`

3.4 Building the Corpus

All the scripts responsible for building the corpus can be found under the folder `corpus_builder`.

3.4.1 HTML Parsing

After the creation of a local copy of the HTML files containing the press releases, the content needs to be extracted from the HTML files and saved in a format that would be suitable for later processing.

Using the Python package BeautifulSoup to parse the HTML files, I extracted from each HTML file the title and the text of the press release, as well as some meta data: date, language and the original file-id and the original file name (for debugging purposes). The data was then saved to a JSON file, one per each year. See listing 3.1 for an example.

```
1 {
2   "0": {
3     "id": "12924",
4     "orig_file": "html/1997/1997_12924_DE.html",
5     "lang": "DE",
6     "title": "25 Jahre Arge Alp: Graubünden feiert tüchtig mit",
7     "date": "31.12.1997",
8     "content": "Die Arge Alp feiert heuer ihr 25-Jahre-Jubiläum. Aus diesem Grund
9                 ↳ finden vom 27. September bis 12. Oktober 1997 die Festwochen des
10                ↳ Alpenraums in Telfs-Mösern, Tirol, statt. ..."
11   },
12   "1": {
13     "id": "12926",
14     "orig_file": "html/1997/1997_12926_DE.html",
15     "lang": "DE",
16     "title": "Kanton will auch personelles Engagement bei der Bündner Kraftwerke AG
17             ↳ verstärken",
18     "date": "31.12.1997",
19     "content": "Nachdem der Kanton Graubünden letzten Herbst die Aktienmehrheit der
20             ↳ Bündner Kraftwerke AG übernommen hat, will er nun auch seine Vertretung
21             ↳ im Verwaltungsrat stärken...."
```

```

21     "lang": "DE",
22     "title": "Graubünden trifft präventive Massnahmen zur Bekämpfung der illegalen
        ↳ Einwanderung",
23     "date": "31.12.1997",
24     "content": "Die Fremdenpolizei des Kantons Graubünden trifft im Einvernehmen
        ↳ mit dem kantonalen Sozialamt, dem Amt für Zivilschutz sowie der
        ↳ Kantonspolizei Graubünden Massnahmen, um die illegale Einwanderung in
        ↳ den Südtälern des Kantons Graubünden zu bekämpfen...."
25 },

```

Listing 3.1: Example for a JSON file

3.4.2 Document Alignment

After extracting the relevant data from the HTML files and saving them in JSON files, the core task can begin: aligning the documents to get document-triples which are translations of each other.

Linked vs. unlinked

For all releases published after mid-2009 this is pretty simple. The file-id extracted from the URLs is common to all three releases in the three languages (see figure 3.2). This file-id can be used to link the press releases with each other. I shall refer to these press releases as “linked releases”.

For releases published prior to that, each release has a unique file-id. This means it can’t be used for document alignment. I shall refer to these releases as “unlinked releases”. For unlinked releases I used a simple heuristic: if on one single date exactly three releases were published in three different languages, I assume they are translations of each other. The titles of press releases that weren’t aligned are saved to a CSV file which can be used for manual alignment.

Unfortunately, this means around 40% (TODO: what is the exact average?) of the releases each year prior to 2009 cannot be automatically added to the corpus, cf. figure 3.3.

Since the year 2009 contains both linked and unlinked releases, I wrote the script `split_2009.py` to split the data accordingly. It uses a very simple heuristic: if the file-id of a press release is longer than 5 digits, it is a linked press releases.

Aligned corpus

The aligned press releases are saved again to JSON files, with each row containing the three press releases in the three languages, along with metadata such as date and file-id. In the rare case that one language is missing (TODO: how rare?), i.e., the press releases wasn’t translated into that language for some reason, it is simply left blank. Press releases that are available only in one language are discarded from the corpus.

The script `create_corpus.py` deals with this task. Using the Python library Pandas, the JSON files are read into a Dataframe. For linked releases, all the unique ID’s are taken, and then for each ID the three languages are collected and saved into a new row. The dates are converted from their original format (DD.MM.YY) to an ISO 8601 format (YYYY-MM-DD) ([enwiki:1095673391](#)) for better compatibility and easier processing later.

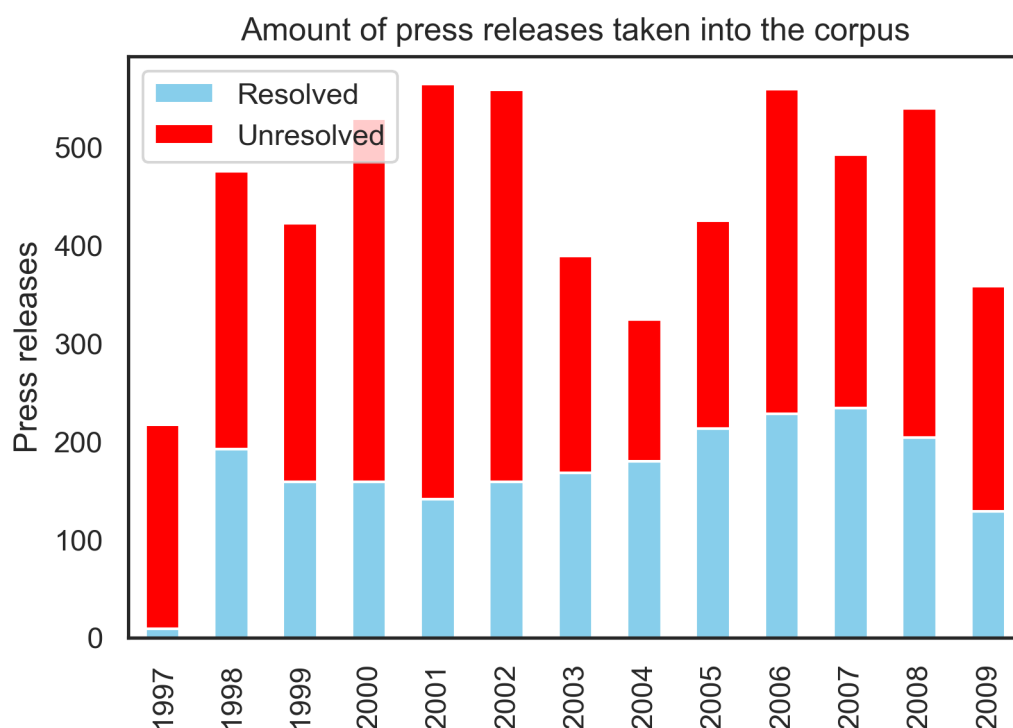


Figure 3.3: Portion of automatically aligned press releases up to 2009

For JSON files containing unlinked documents, the script `create_corpus` has to be run with the switch `--by-date`, which tells the program to use the date for aligning the documents instead of the file ID.

For an example of the resulting JSON files, with each row containing the aligned documents, see listing 3.2.

```

1 {
2   "0": {
3     "id": "2010010501",
4     "date": "2010-01-05",
5     "DE_title": "Neues Online-Angebot für das Bündner Rechtsbuch",
6     "DE_content": " Das im Internet verfügbare Bündner Rechtsbuch ist neu gestaltet
7       ↳ worden und enthält neue Funktionalitäten. ...",
8     "IT_title": "Nuova offerta online per la Collezione sistematica del diritto
9       ↳ cantonale grigionese",
10    "IT_content": " La Collezione sistematica del diritto cantonale grigionese
11      ↳ disponibile in internet è stata ristrutturata e contiene nuove funzioni.
12      ↳ ... ",
13    "RM_title": "Nova purschida d'internet per il cudesch da dretg grischun",
14    "RM_content": " Il cudesch da dretg grischun che stat a disposiziun en
15      ↳ l'internet ha survegnì in nov concept e novas funcziuns. ... "
16  },
17  "1": {
18    "id": "2010010502",
19    "date": "2010-01-05",

```



```

15     "DE_title": "Staupe bei Füchsen und Dachsen im Puschlav",
16     "DE_content": " Nachdem sich im Verlaufe des letzten Herbstes die
        ↳ Staupe-Krankheit bei Wildtieren in Nord- und Mittelbünden verbreitete,
        ↳ sind im Laufe der letzten Wochen nun auch im Puschlav bei Füchsen und
        ↳ Dachsen Infektionen mit dem Staupevirus nachgewiesen worden. ... ",
17     "IT_title": "Volpi e tassi affetti da cimurro in Valposchiavo",
18     "IT_content": " Dopo che nel corso dell'autunno il cimurro si è diffuso tra gli
        ↳ animali selvatici del Grigioni settentrionale e centrale, nelle ultime
        ↳ settimane la presenza del virus è stata rilevata anche tra volpi e tassi
        ↳ della Valposchiavo. ... ",
19     "RM_title": "Pesta dals chauns tar vulps e tar tass en il Puschlav",
20     "RM_content": " Suentar che la pesta da chauns è sa derasada tar la selvaschina
        ↳ dal Grischun dal nord e central en il decurs da l'atun passà, èn
        ↳ vegnidas cumprovadas en il decurs da las ultimas emnas ussa er
        ↳ infecziuns cun il virus da questa malsogna tar vulps e tar tass en il
        ↳ Puschlav. ... "
21 },
22 "2": {
23     "id": "2010010801",
24     "date": "2010-01-08",
25     "DE_title": "Projekt Sicherheitsfunknetz POLYCOM Graubünden mit
        ↳ Vertragsunterzeichnung offiziell gestartet",
26     "DE_content": " Die Vorsteherin des Departements für Justiz, Sicherheit und
        ↳ Gesundheit, Regierungsrätin Barbara Janom Steiner, und der Chef des
        ↳ Grenzwachtkorps, Jürg Noth, haben heute in Chur eine Vereinbarung zur
        ↳ Realisierung des Sicherheitsfunknetzes POLYCOM im Kanton unterzeichnet.
        ↳ ... ",
27     "IT_title": "Avviato ufficialmente con la sottoscrizione del contratto il
        ↳ progetto di rete radio di sicurezza POLYCOM Grigioni",
28     "IT_content": " La Consigliera di Stato Barbara Janom Steiner, direttrice del
        ↳ Dipartimento di giustizia, sicurezza e sanità, e il capo del Corpo delle
        ↳ guardie di confine, Jürg Noth, hanno sottoscritto oggi a Coira un
        ↳ accordo per la realizzazione nel Cantone della rete radio di sicurezza
        ↳ POLYCOM. ... ",
29     "RM_title": "Il project per la rait radiofonica da segirezza POLYCOM dal
        ↳ Grischun è vegnì lantschà uffizialmain cun suttascriber il contract",
30     "RM_content": " La scheffa dal departament da giustia, segirezza e sanadad,
        ↳ cussegliera governativa Barbara Janom Steiner, ed il schef dal corp da
        ↳ guardias da cunfin, Jürg Noth, han suttascrit oz a Cuiira ina cunvegna
        ↳ per realisar la rait radiofonica da segirezza POLYCOM en il chantun. ...
        ↳ "
31 },
32 }

```

Listing 3.2: Example for a JSON file containing aligned documents

3.5 Manual alignment of unlinked documents

As can be seen in figure 3.3, a big portion of the unlinked documents cannot be automatically aligned using the simple heuristic described in section 3.4.2. To deal with that, the script `create_corpus.py` will write the titles and file IDs of all the discarded documents to a CSV file, one for each year.

These CSV files can be used for manually aligning the corresponding documents. This is done by enumerating the documents while using the same digit for corresponding documents. The CSV files containing the now enumerated documents can be given to the script `create_corpus.py` using the argument `add-from-csv` to combine the enumerated documents as linked documents into the JSON file.

3.6 SQLite database

The query language SQL offers flexible and complex way to query databases. For this reason, I decided to save the resulting corpus in an SQLite database. I opted for SQLite because it doesn't require running a separate server and can be SQLite databases can be easily built, edited and accessed using `sqlite3`¹, a Python module delivered in the Python standard library².

The SQLite database contains two tables, `corpus` and `raw` with the exact same structure as the two JSON files described in listings 3.1 and 3.2.

3.7 Summary

For compiling the corpus, the following steps were taken:

1. Scrape website and save HTML documents locally
2. Extract relevant content from HTML files (date, language, title and content) and save to JSON files
3. Read the JSON files using Pandas Dataframe, align the documents and save to new JSON files
4. Feed both types of JSON files to an sqlite database

The final result is an sqlite database (`corpus.db`) containing two tables:

- `corpus`: all the aligned documents from 1997 until today. Each row contains following columns:
 - `id`: automatically incremented unique ID
 - `file_id`: original file ID
 - `date`: Release date
 - `DE_title`: Title of German document

¹<https://docs.python.org/3/library/sqlite3.html>

²<https://docs.python.org/3/tutorial/stdlib.html>

Year	Documents	Year	Documents
1997	3	2010	184
1998	64	2011	167
1999	53	2012	207
2000	53	2013	219
2001	47	2014	218
2002	53	2015	183
2003	56	2016	190
2004	60	2017	207
2005	71	2018	221
2006	76	2019	216
2007	78	2020	286
2008	68	2021	294
2009	109	2022	153

Table 3.1: Number of parallel documents per year, as of July 20, 2022.

- DE_content: Content of German document
 - IT_title: Title of Italian document
 - IT_content: Content of Italian document
 - RM_title: Title of Romansh document
 - RM_content: Content of Romansh document
- raw: all the documents contained in the HTML files scraped from the website. Each row contains the following columns:
 - id: Automatically incremented unique ID
 - file_id: Original file ID
 - orig_file: Original filename
 - lang: Document language (DE for German, IT for Italian, RM for Romansh)
 - title: Document title
 - date: Release date
 - content: Document content

3.7.1 Statistics

The corpus contains 3,536 parallel documents.

Glossary

Graubünden The Canton of Grisons. 1, 5

Acronyms

AER Average Error Rate. 33, 41, 43, 44, 46, 53

NER Named Entity Recognition. 2

POS Part of Speech. 2

List of Tables

2.1	Examples for choosing the forms for Rumanstch Grischun, based on liver1999	7
3.1	Number of parallel documents per year, as of July 20, 2022.	15
4.1	Parallel corpus in numbers, as of July 20, 2022. “Source” refers to the language on the left and “target” to the language on the right, and not necessarily to the actual source language of the translation.	23
6.1	Translation examples of German compounds into Romansh	37
7.1	Evaluation metrics for word alignments with the baseline models for different dataset sizes. “Dataset Size” refers to the number of sentence pairs.	43
7.2	Evaluation metrics for word alignments using SimAlign, with different embeddings and word/sub-word level. Best result per embedding type in bold.	44
7.3	Comparison of the best performance of each of the three methods. The best value in each column is in bold.	44

List of Figures

2.1	Distribution of Rhaeto-Romance, taken from haiman1992	5
3.1	Directory tree of corpus_builder	9
3.2	Directory scheme for saving the HTML files	9
3.3	Portion of automatically aligned press releases up to 2009	12
4.1	Sentence alignment pipeline	20
5.1	Word alignment example	25
5.2	Similarity matrix	32
5.3	Alignment matrix	32
5.4	The resulting word alignment	32
6.1	Aligning German perfect to Romansh perfect	38
6.2	Alignment of German preterite to Romansh perfect	38
6.3	Aligning German present participles to Romansh relative clauses	39
7.1	Comparing precision between the systems for different dataset sizes.	45
7.2	Comparing recall between the systems for different dataset sizes.	45
7.3	Comparing AER between the systems for different dataset sizes.	46
A.1	Word alignment example 1	58
A.2	Word alignment example 2	59
A.3	Word alignment example 3	59
A.4	Word alignment example 4	60
A.5	Word alignment example 5	60