# Contents

# 1 ACKNOWLEDGEMENTS

# 2 THE SPARV PIPELINE

## 2.1 SETTING UP THE SPARV PIPELINE

This section describes how to get the Sparv corpus pipeline up and running on your own machine. The source code is available from GitHub[1] under the MIT license[2].

Please note that different license terms may apply to any of the additional components that can be plugged into the Sparv Pipeline!

### 2.1.1 REQUIREMENTS

- A Unix-like environment (e.g. Linux, macOS)
- Python 3.4[3] or newer
- GNU Make[4]
- Java[5] (if you want to use the MaltParser, Sparv-wsd or hfst-SweNER)

Additional components (optional):

- Git[6] and Git Large File Storage[7] (for cloning the repository, strongly recommended!)
- Hunpos[8], with its path included in your `PATH` environment variable (for part-of-speech tagging)
- MaltParser[9] v. 1.7.2 (for dependency parsing)
- Sparv-wsd[10] (for word-sense disambiguation)
- hfst-SweNER[11] (for named entity recognition)
- FreeLing[12] 4.1 (if you want to annotate corpora in Catalan, English, French, Galician, German, Italian, Norwegian, Portuguese, Russian, Slovenian or Spanish)
- TreeTagger[13] (if you want to annotate corpora in Bulgarian, Dutch, Estonian, Finnish, Latin, Polish, Romanian or Slovak)
- fast_align[14] (if you want to run word-linking on parallel corpora)
- Corpus Workbench[15] (CWB) 3.2 or newer (if you are going to use the Korp backend for searching in your corpora)

### 2.1.2 INSTALLING THE REQUIRED SOFTWARE

The following information assumes that you are running Ubuntu, but will most likely work for any Linux-based OS.

1. Install the Sparv Pipeline by cloning the Git repository[16] into a directory of your choice. Please note that you must have Git Large File Storage[17] installed on your machine before cloning the repository. Some files will not be downloaded correctly otherwise. If you do not want to use Git you can download the latest release package from the release page[18].
2. Setup a new environment variable `SPARV_MAKEFILES` and point it to the directory `sparv-pipeline/makefiles`.

---

[1] https://github.com/spraakbanken/sparv-pipeline
[2] https://opensource.org/licenses/MIT
[3] http://python.org/
[4] https://www.gnu.org/software/make/
[5] http://www.oracle.com/technetwork/java/javase/downloads/jdk8-downloads-2133151.html
[6] https://git-scm.com/
[7] https://git-lfs.github.com/
[8] http://code.google.com/p/hunpos/
[9] http://www.maltparser.org/download.html
[10] https://github.com/spraakbanken/sparv-wsd
[11] http://www.ling.helsinki.fi/users/janiemi/finclarin/ner/hfst-swener-0.9.3.tgz
[12] https://github.com/TALP-UPC/FreeLing/releases/tag/4.1
[13] http://www.cis.uni-muenchen.de/~schmid/tools/TreeTagger/
[14] https://github.com/clab/fast_align
[15] http://cwb.sourceforge.net/beta.php
[16] https://github.com/spraakbanken/sparv-pipeline
[17] https://git-lfs.github.com/
[18] https://github.com/spraakbanken/sparv-pipeline/releases

#### 2.1.2.1 PYTHON VIRTUAL ENVIRONMENT

Set up a Python virtual environment as a subdirectory to the `sparv-pipeline` directory:

```
python3 -m venv venv
```

Activate the virtual environment and install the required Python packages (in some cases you are required to upgrade pip first):

```
source venv/bin/activate
pip install --upgrade pip
pip install -r requirements.txt
```

You can then deactivate the virtual environment:

```
deactivate
```

#### 2.1.2.2 CONFIGURING SOME VARIABLES

In `makefiles/Makefile.config` you will find a section called *Configuration*. Here you will need to specify the path to the pipeline directory by setting the variable `SPARV_PIPELINE_PATH`. If you are planning on using the pipeline to install corpora on a remote computer, you must (in addition to installing CWB) edit the `remote_host`, `remote_cwb_datadir` and `remote_cwb_registry` variables. If you are going to use anything database related (generating lemgram index or Word Picture data for Korp) you also have to edit the value of `rel_mysql_dbname`.

#### 2.1.2.3 DOWNLOADING AND UPDATING THE MODELS

Certain annotations need models to work. To download and generate the models follow these simple steps:

1. Navigate to the `sparv-pipeline/models` directory in a terminal.
2. Run the command `make clean`, and then `make all`.
3. If no errors were reported while running the above commands you may run `make space` for saving disk space.

By default, the part-of-speech tagger also relies on the SALDO model. You can disable this dependency by commenting out the line beginning with `hunpos_morphtable` in Makefile.config, or changing its value to `$(SPARV_MODELS)/hunpos.suc.morphtable`.

### 2.1.3 INSTALLING ADDITIONAL COMPONENTS (OPTIONAL)

The following components are not part of the core Sparv Pipeline package. Whether or not you will need to install these components depends on how you want to use the Sparv Pipeline. Please note that different licenses may apply for the individual components.

#### 2.1.3.1 HUNPOS

Hunpos[19] is used for Swedish part-of-speech tagging and it is a prerequisite for all of the SALDO-annotations. Hunpos can be downloaded from here[20]. Installation is done by unpacking and then adding the path of the executables to your PATH environmental variable. If you are running a 64-bit OS, you might also have to install 32-bit compatibility libraries if Hunpos won't run:

---

[19]http://code.google.com/p/hunpos/
[20]http://code.google.com/p/hunpos/

```
sudo apt install ia32-libs
```

On Arch Linux, activate the `multilib` repository and install `lib32-gcc-libs`.

If that doesn't work, you might have to compile Hunpos from its source code.

#### 2.1.3.2 MALTPARSER

MaltParser[21] is used for Swedish dependency parsing. The version compatible with the Sparv pipeline is 1.7.2. Download and unpack the zip-file from the MaltParser home page and place its contents under `sparv-pipeline/bin/maltparser-1.7.2`.

Download the dependency model[22] and place the file `swemalt-1.7.2.mco` under `sparv-pipeline/models`.

#### 2.1.3.3 SPARV-WSD

The Sparv-wsd[23] is used for Swedish word-sense disambiguation. It is developed at Språkbanken and runs under the same license as the Sparv Pipeline core package. In order to use it within the Sparv Pipeline it is enough to download the saldowsd.jar from GitHub and place it inside the `sparv-pipeline/bin/wsd` directory:

```
wget https://github.com/spraakbanken/sparv-wsd/raw/master/bin/
    saldowsd.jar -P sparv-pipeline/bin/wsd/
```

Its models are added automatically when building the Sparv Pipeline models (see above).

#### 2.1.3.4 HFST-SWENER

The current version of hfst-SweNER[24] expects to be run in a Python 2 environment while the Sparv pipeline is written in Python 3. If you want to use hfst-SweNER's named entity recognition from within Sparv you need to make sure that your `python` command in your environment refers to Python 2. Alternatively, before installing hfst-SweNER, you can make sure that it will be run with the correct version of Python by replacing `python` with `python2` in all the Python scripts in the `hfst-swener-0.9.3/scripts` directory. The first line in every script will then look like this:

```
#! /usr/bin/env python2
```

#### 2.1.3.5 TREETAGGER AND FREELING

TreeTagger[25] and FreeLing[26] are used for POS-tagging and lemmatization of other languages than Swedish. Please install the software according to the instructions on the respective website or in the provided readme file.

The following is a list over the languages currently supported by the corpus pipeline, their language codes and which tool Sparv uses to analyze them:

| Language | Code | Analysis Tool |
|----------|------|---------------|
| Bulgarian | bg | TreeTagger |

---

[21] http://www.maltparser.org/download.html
[22] http://maltparser.org/mco/swedish_parser/swemalt.html
[23] https://github.com/spraakbanken/sparv-wsd
[24] http://www.ling.helsinki.fi/users/janiemi/finclarin/ner/hfst-swener-0.9.3.tgz
[25] http://www.cis.uni-muenchen.de/~schmid/tools/TreeTagger/
[26] https://github.com/TALP-UPC/FreeLing/releases/tag/4.1

| Language | Code | Analysis Tool |
| --- | --- | --- |
| Catalan | ca | FreeLing |
| Dutch | nl | TreeTagger |
| Estonian | et | TreeTagger |
| English | en | FreeLing |
| French | fr | FreeLing |
| Finnish | fi | TreeTagger |
| Galician | gl | FreeLing |
| German | de | FreeLing |
| Italian | it | FreeLing |
| Latin | la | TreeTagger |
| Norwegian | no | FreeLing |
| Polish | pl | TreeTagger |
| Portuguese | pt | FreeLing |
| Romanian | ro | TreeTagger |
| Russian | ru | FreeLing |
| Slovak | sk | TreeTagger |
| Slovenian | sl | FreeLing |
| Spanish | es | FreeLing |
| Swedish | sv | Sparv |
| Swedish 1800's | sv-1800 | Sparv |
| Swedish development mode | sv-dev | Sparv |

If you are using TreeTagger, please copy the `tree-tagger` binary file into the `sparv-pipeline/bin/treetagger` directory belonging to the pipeline. The TreeTagger models (parameter files) need to be downloaded separately and saved in the `sparv-pipeline/models/treetagger` directory. The parameter files need to be re-named to a two-letter language code followed by the file ending `.par`, e.g. the Dutch parameter file is called `nl.par`.

When using Freeling you will need the sparv-freeling extension which is available via GitHub[27] and runs under the AGPL license[28]. Follow the installation instructions for the sparv-freeling module on GitHub[29] in order to set up the models correctly.

### 2.1.3.6 FAST_ALIGN

fast_align[30] is used for word-linking on parallel corpora. Follow the installation instructions given in the fast_align repository and copy the binary files `atools` and `fast_align` into the `sparv-pipeline/bin/word_alignment` folder of the pipeline.

### 2.1.4 PREPARING A CORPUS

Create a new directory for your corpus (e.g. `mycorpus`). Under your new directory create another directory containing your input texts (e.g. `original`).

The input text must meet the following criteria:

- The files must be encoded in **UTF-8**.
- The format must be **valid XML**, *with the following two exceptions*:
    1. No root element is needed. It is however required that all text in the file must be contained within elements.
    The simplest valid file possible consists of a single element, for example `<text>`, containing nothing but raw text.

---

[27] https://github.com/spraakbanken/sparv-freeling
[28] http://www.gnu.org/licenses/agpl.html
[29] https://github.com/spraakbanken/sparv-freeling
[30] https://github.com/clab/fast_align

2. Overlapping elements are allowed (`<a><b></a></b>`).
- No file should be larger than **~10-20 MB**, or you might run into memory problems. Split larger files into smaller files.

Every corpus needs a *Makefile* in which you configure the format of your input material and what annotations you want. You can use the included `Makefile.example` as a base for a simple corpus, or `Makefile.template` if you want to dive into the more advanced stuff. Whichever file you choose, you should put it in the directory you created for your corpus, and name the file `Makefile`.

By default `Makefile.example` is configured to read the source files from a directory named `original`, and assumes that all the text is contained within a `<text>` element. This can easily be changed by editing the `original_dir` and `xml_elements` variables in the Makefile. For a description of every available variable and setting, see `Makefile.template`.

Once you have edited the Makefile to fit your source material and needs, you are ready to execute it by running one of the following commands in a terminal:

`make TEXT`
This will parse your source files, which is a good place to start just to make sure everything works as it should. If there are any problems with the format of your source files, the script will complain.

`make vrt`
This will run every type of annotation that you have specified in the Makefile, and will output a VRT file for every input file. The VRT format is used as input for Corpus Workbench to create binary corpus files. The resulting files will be found in the `annotations` directory (created automatically).

`make export`
Same as `make vrt`, except the output format is XML and the files are saved to the `export.original` directory.

`make cwb`
This command will take the VRT files and convert them to a Corpus Workbench corpus.

`make install_corpus`
This will copy your Corpus Workbench corpus to a remote computer.

`make relations`
Provided that you are using syntactic parsing, this will generate the relations data for the Word Picture in Korp.

`make install_relations`
This will install the relations data to a remote MySQL server.

`make add`
If you remove or add source files to your corpus after having already run one of the commands above, you will have to run the above command.

For a complete list of available commands, run `make` without any arguments.

## 2.1.5 INSTALLING CORPUS WORKBENCH (OPTIONAL)

If you are not going to use the Korp backend, you can skip this step.

You will need the latest version of CWB for unicode support. Install by following these steps:

Check out the latest version of the source code from subversion by running the following command in a terminal:

```
svn co https://cwb.svn.sourceforge.net/svnroot/cwb/cwb/trunk cwb
```

Navigate to the new cwb directory and run the following command:

```
sudo ./install-scripts/cwb-install-ubuntu
```

CWB is now installed, and by default you will find it under `/usr/local/cwb-X.X.X/bin` (where `X.X.X` is the version number). Confirm that the installation was successful by typing:

```
/usr/local/cwb-X.X.X/bin/cqp -v
```

CWB needs two directories for storing the corpora. One for the data, and one for the corpus registry. You may create these directories wherever you want, but let's assume that you have created the following two:

```
/corpora/data
/corpora/registry
```

You also need to edit the following variables in `sparv-pipeline/makefiles/Makefile.config`, pointing to the directories created above:

```
export CWB_DATADIR ?= /corpora/data
export CORPUS_REGISTRY ?= /corpora/registry
```

If you're not running Ubuntu or if you run into any problems, refer to the INSTALL text file in the cwb dir for further instructions.

## 2.2 DEVELOPING THE SPARV PIPELINE

This section describes the internal structure of the Sparv pipeline and its usage and can be used as a developer's guide. If you are working at Språkbanken you should use the latest version from Subversion. Check the internal wiki page for more information.

### 2.2.1 DOWNLOADING AND INSTALLING THE SPARV PIPELINE

The latest distribution of the pipeline can be downloaded from the distribution page[31].

Annotating corpora with the Sparv pipeline can be done in the same manner as on Språkbanken's own servers which is described here[32] (in Swedish only).

The import format is described here[33].

### 2.2.2 THE PIPELINE'S FOLDER STRUCTURE

```
annotate/
    bin/       binary files
    makefiles/ Makefiles
    models/    linguistic models
    python/    Python modules
    tests/     tests
```

---

[31]https://spraakbanken.gu.se/eng/research/infrastructure/sparv/distribution/pipeline
[32]https://spraakbanken.gu.se/swe/forskning/infrastruktur/korp/korpusimport
[33]https://spraakbanken.gu.se/eng/research/infrastructure/sparv/importformat

### 2.2.3 THE PIPELINE STRUCTURE

The Sparv corpus pipeline consists of Python modules which are located in the `python` directory. Most of the modules contain code to solve a single task, e.g. part-of-speech tagging, tokenisation etc. Some of the modules call external tools, e.g. Hunpos or the MaltParser.

The following is a short summary of most of the existing modules and their basic functionalities. More information on every script can be found in the doc strings of the functions.

| Module | Description |
| --- | --- |
| *annotate.py* | Different types of small annotations. |
| *compound.py* | Compound analysis using SALDO. |
| *crf.py* | Sentence segmentation with Conditional Random Fields. |
| *cwb.py* | Creates the different output formats (VRT, XML, Corpus Workbench files). |
| *dateformat.py* | Date formatting. |
| *diapivot.py* | Creates the diachronic pivot, linking older Swedish texts to SALDO. |
| *fileid.py* | Creates IDs for every file in the corpus. |
| *freeling.py* | For annotation of other languages. Calls the Freeling software and processes its output. (Only available in the AGPL-version of Sparv) |
| *geo.py* | Annotates geographical features. |
| *hist.py* | Different functions used for annotation historical Swedish texts. |
| *hunpos_morphtable.py* | Creates a morphtable file for Hunpos, based on SALDO. |
| *hunpos.py* | Part-of-speech tagging with Hunpos. |
| *info.py* | Creates an info file which is used by CWB. |
| *install.py* | Installs corpora on other machines. |
| *lemgram_index.py* | Creates a lemgramindex which is used by Korp. |
| *lexical_classes.py* | Creates annotations for lexical classes with different resources (Blingbring, SweFN). |
| *lmflexicon.py* | Parses LMF lexicons. |
| *malt.py* | Syntactic parsing with the MaltParser. |
| *number.py* | Different types of numbering used e.g. in structural attributes. |
| *parent.py* | Adds annotations for parent links and/or children links (needed for pipeline internal purposes). |
| *readability.py* | Creates annotations for readability measures. |
| *relations.py* | Creates data used by the word picture. |
| *saldo.py* | Different types of SALDO-annotations. |
| *segment.py* | All kinds of segmentation/tokenisation (paragraph, sentence, word). |
| *sent_align.py* | Sentence linking of parallel corpora. |
| *sentiment.py* | Creates sentiment annotations. |
| *suc2hunpos.py* | Creates a test material based on SUC3 for Hunpos. |
| *swener.py* | Creates named-entity annotations with hfst-SweNER. |
| *timespan.py* | Handles time spans. |
| *train_nst_comp_model.py* | Trains a compound part-of-speech model used by the compound analysis. |
| *train_stats_model.py* | Trains a statistical model used by the compound analysis. |
| *treetagger.py* | For annotation of other languages. Calls the TreeTagger software and processes its output. |
| *vw.py* | Topic modelling with vowpal wabbit. |
| *word_align.py* | Word linking of parallel corpora with fast_align and cdec. |
| *wsd.py* | Crates annotations for word sense disambiguation. |
| *xmlanalyser.py* | An analyzer for pseudo-XML documents. |
| *xmlparser.py* | Parses XML and generates files in the pipeline's working format. |

The first thing that happens when annotating a corpus with the Sparv pipeline is that the XML input data is parsed by the module xmlparser.py. This module generates the following files for every input file (usually in the `annotations` directory):

- a `.@TEXT` file containing the source text with anchors
- one file for each structural attribute captured, where each line consists of a reference to two anchors indicating the span of the attribute, and the value of the attribute

The output of most modules are files in the above described format, with references to anchors and data for the indicated span. The part-of-speech tagger module for instance creates a file containing an anchor span per word and the part-of-speech as value.

### 2.2.4 ANNOTATIONS

Below is a list with some of the most common annotations and which module they are built with:

- *word*
  The string containing a token (word or punctuation). This annotation is created by *xmlparser.py*.
- *msd*
  A morphosyntactic descriptor. Created by *hunpos.py* with *word* as input data.
- *pos*
  A simplified MSD, preserving only its initial part. Created by *annotate.py* with *msd* as input data.
- *baseform, lemgram, sense*
  Citation form, lemgram, and SALDO-ID. Created by *saldo.py*, with *word*, *msd* and *ref* as input files. Uses a SALDO model which is created with *saldo.py*.
- *ref*
  En numrering av varje token inom varje mening. Skapas av *annotate.py*.
- *dephead, deprel*
  Dependency head and dependency relation. Created by *malt.py*, with *word*, *msd* and *pos* as input data. Uses the MALT parser with a model for Swedish.

### 2.2.5 DEVELOPING A NEW MODULE

This part describes how to create a new Python module within the Sparv pipeline.

Usually one does not have to define the format for the input and output data since there are methods defined for reading and writing these.

Below is a minimal example script which could be integrated in the Sparv pipeline:

```
# -*- coding: utf-8 -*-
import sb.util as util

def changecase(word, out, case):
    # Read the annotation "word"
    words = util.read_annotation(word)

    # Transform every word to upper or lower case
    for w in words:
        if case == "upper":
            words[w] = words[w].upper()
        elif case == "lower":
            words[w] = words[w].lower()

    # Write the result (out) to the indicated annotation (words
        )
    util.write_annotation(out, words)

if __name__ == '__main__':
    util.run.main(changecase=changecase)
```

Most modules start with importing `util` which contains functions for reading and writing annotations in the data format used within the pipeline.

Then you can define functions that may be used by Sparv. These functions usually take paths to existing annotation files as arguments as well as a path to the annotation file that should be created by the function.

The function `util.read_annotation()` is used to read existing annotations. It returns a dictionary object with anchors as keys and annotation values as values.

A module function can be concluded with calling `util.write_annotation()` which takes a path to an annotation file and a dictionary as arguments and saves the dictionary contents to the specified file path.

A module should also contain some code that makes its functions available from the command line. This is what the last two lines in the above example are there for. As arguments to the function `util.run.main` you can list all the names of the functions you want to be able to use from the command line in the following format: `alias=name_of_function`. The chosen alias can then be specified as a flag in the command line. The first function may be set as default by omitting the alias.

### 2.2.6 MAKEFILE

When running the Sparv pipeline its different modules are coordinated with makefiles. There are two global makefiles containing configurations and rules for the creation of annotations (`Makefile.config`, `Makefile.rules`). Furthermore, each corpus has its own makefile (`Makefile`) that contains information about the data format and about which annotations should be used. The two global makefiles are imported by the corpus specific makefile.

Documentation for the different corpus specific makefiles can be found in `Makefile.template`. This file is distributed together with the other makefiles in the Sparv pipeline.

`Makefile.rules` is the file containing all rules that create annotations. `Makefile.config` contains variables that can be configured.

A rule in Makefile.rules can for example look like this:

```
%.token.uppercase: %.token.word
    $(python) -m sb.case --changecase --out $@ --word $(1) --
        case upper
```

The above rule creates the annotation *uppercase* from the Python example above. It takes the annotation *word* as input and calls the function *uppercase* in the script *case.py* with *out*, *word* and *case* as arguments.

Adding a new annotation can in most cases be done by adding a rule like in the above example to `Makefile.rules`. The name of the annotation must also be added to the list of annotations i the corpus specific makefile.

## 2.3  Sparv's Import Format

This section describes the XML format you should use when importing text using the Sparv corpus pipeline[34].

### 2.3.1  Format requirements

The input text files must meet the following criteria:

- The files must be encoded in **UTF-8**.
- The format must be **valid XML**, *with the following two exceptions*:
    1. No root element is needed. It is however required that all text in the file must be contained within elements.
       The simplest valid file possible consists of a single element, for example `<text>`, containing nothing but raw text.

    2. Overlapping elements are allowed (`<a><b></a></b>`).
- No file should be larger than **~10-20 MB**, or you might run into memory problems. Split too large files into smaller files.

### 2.3.2  Metadata

Any metadata about the text as a whole or parts of the text (e.g. chapters, sentences, words) must exist in the form of attributes to the elements containing the text. For example:

```
<text title="Title metadata here">
  This is the corpus text.
</text>
```

*All text*, except for element attributes, will be part of the corpus text, and will not be treated as metadata. The following example is *not* a valid way to set title metadata for a text, and the string "Title metadata here" will be treated as corpus text:

```
<text>
  <title>Title metadata here</title>
  This is the corpus text.
</text>
```

#### 2.3.2.1  Headers

One exception to the above is when there is a header in the file, such as:

---

[34]https://spraakbanken.gu.se/eng/research/infrastructure/sparv/distribution/pipeline

```
<text>
  <header>
    <title>Title metadata here</title>
  </header>
  This is the corpus text.
</text>
```

In this case it is possible to treat the whole `header` element (not necessarily named `header`) as not belonging to the text, and instead extract metadata from it. The metadata will be tied to all text contained by the root element. In this case the root element is `text`, and the result will be the same as in the first example.

The following is *not* a valid example where you can extract header metadata, due to the header being disconnected from the text:

```
<header>
  <title>Title metadata here</title>
</header>
<text>
  This is the corpus text.
</text>
```

However, had both `header` and `text` been contained within another element, header metadata extraction would have been possible.

### 2.3.3 CONFIGURING THE MAKEFILE

When using the corpus pipeline you configure your Makefile to match your XML format, and tell it what elements and attributes to parse. The variables used for this are the following:

```
xml_elements
xml_annotations
```

In `xml_elements` you list (in lowercase, separated by whitespace) the elements of interest, optionally with their attributes (on the form `element:attribute`). `xml_annotations` is a list of the same size, with every item corresponding to an item in `xml_elements`. This list contains the names of the annotation files that will be created during parsing. For example:

```
xml_elements =    text:title  text:author  s
xml_annotations = text.title  text.author  sentence
```

The above will parse elements such as `<text title="X" author="Y">` and save title and author as separate annotations. It will also parse the attribute `<s>`, not saving any attribute data but only the boundaries of the element, marking the beginning and end of every sentence (provided that our example text has sentence segmentation marked with `<s>` tags). In this way you can use existing annotations of a material during the import, rather than letting the pipeline create new annotations for you. By saving the `<s>` element as the annotation named `sentence`, the pipeline will not try to create that annotation itself. There are several hardcoded annotation names, all of which can be overridden like `sentence` above:

```
Positional: token.pos, token.msd, token.baseform, token.lemgram,
   token.saldo, token.prefix, token.suffix
Structural: token, sentence, paragraph
```

13

For positional annotations (annotations for single tokens), use the prefix "token." as shown above.

You can combine several elements into one annotation, if for example sentence segmentation is marked up with <s> in some places and <x> in some. Combining them is done by simply joining both elements with a "+":

```
xml_elements =     s+x
xml_annotations = sentence
```

### 2.3.3.1  HEADERS

For files containing headers like in the previously shown example, three variables has to be set in the Makefile:

```
xml_header
xml_headers
xml_header_annotations
```

The first one should be set to the name of the header element, in lowercase. This tells the parser that this element is a header, and none of its text content will be part of the corpus text.

```
xml_header = header
```

Alternatively you can give the full path to the header, joining elements by ".":

```
xml_header = text.header
```

The second and third variable is similar to the `xml_elements`/`xml_annotations` pair. The first one is a list of headers to parse, and the second a list of corresponding annotation files that will be created:

```
xml_headers
xml_header_annotations
```

`xml_headers` must contain the full path to the element in question, in lowercase. As with `xml_elements` you can specify attributes using a colon, but it is also possible to parse the text contained by the element, by using the keyword "TEXT" instead of an attribute:

```
xml_headers             = text.header.title:TEXT
xml_header_annotations = text.title
```

Any periods in element names have to be escaped using backslash.

### 2.3.4  AN EXAMPLE

The following is an example input file, with existing sentence segmentation and tokenization. Every token is annotated with its baseform.

```
<text title="Test">
  <s>
    <w bf="it">It</w>
    <w bf="be">was</w>
    <w bf="an">an</w>
    <w bf="elephant">elephant</w>
    <w bf=".">.</w>
  </s>
</text>
```

The configuration of the Makefile will look like this:

```
xml_elements =      text:title  s         w      w:bf
xml_annotations = text.title    sentence  token  token.baseform
```

## 2.4 AVAILABLE SPARV ANNOTATIONS

For texts written in contemporary Swedish Sparv can generate the following types of annotations:

- Part of speech tagging:
    - `pos`: part of speech
    - `msd`: morphosyntactic tag

    Tool: Hunpos[35], Model: in-house model trained on SUC 3.0[36] Tag set: MSD tags[37]
- SALDO[38]-based analysis:
    - `baseform`: citation form
    - `lemgram`: lemgram, identifies the inflectional table (using SALDO tags[39])
    - `sense`: identifies a sense in SALDO and its probability
    - (`saldo`: identifies a sense in SALDO - *will be removed soon*)
    - `sentiment`: sentiment score using SenSALDO[40]
- Compound analysis (also based on SALDO):
    - `complemgram`: compound lemgram
    - `compwf`: compound word form
    - (`prefix`: initial part of a compound - *will be removed soon*)
    - (`suffix`: final part of a compound - *will be removed soon*)
- Dependency analysis:
    - `ref`: the position of the word in the sentence
    - `dephead`: dependency head, the ref of the word which the current word modifies or is dependent of
    - `deprel`: dependency relation, the relation of the current word to its dependency head, see http://stp.ling.uu.se/~nivre/swedish_treebank/dep.html

    Tool: MaltParser[41], Model: swemalt, trained on Swedish Treebank
- Named entity recognition:
    - `ne.ex`: named entity (name expression, numerical expression or time expression)
    - `ne.type`: named entity type
    - `ne.subtype`: named entity sub type

    Tool: hfst-SweNER[42]

---

[35]http://code.google.com/p/hunpos/
[36]/en/resources/suc3
[37]https://spraakbanken.gu.se/korp/markup/msdtags.html
[38]/en/resources/saldo
[39]https://spraakbanken.gu.se/eng/research/saldo/tagset
[40]/resurser/sensaldo
[41]http://www.maltparser.org/download.html
[42]http://www.ling.helsinki.fi/users/janiemi/finclarin/ner/hfst-swener-0.9.3.tgz

References: HFST-SweNER – A New NER Resource for Swedish[43], Reducing the effect of name explosion[44]

- Readability metrics:
    - `text.lix`: the Swedish readability metric LIX, läsbarhetsindex
    - `text.ovix`: the Swedish readability metric OVIX, ordvariationsindex
    - `text.nk`: the Swedish readability metric nominalkvot
- Lexical classes:
    - `blingbring`: lexical class from the Blingbring resource (on word level)
    - `swefn`: frames from swedish FrameNet (on word level)
    - `text.blingbring`: lexical class from the Blingbring resource (on document level)
    - `text.swefn`: frames from swedish FrameNet (on document level)

Older Swedish texts or texts written in other languages can often be annotated with a sub set of the above annotation types.

The `msd` annotation for non-Swedish languages is based on different tag sets, depending on which language is annotated and what annotation tool is being used. The attribute contains information about the part of speech and in many cases morphosyntactic information. The `pos` annotation contains only part-of-speech information and uses the universal POS tag set[45].

---

[43]http://www.lrec-conf.org/proceedings/lrec2014/pdf/391_Paper.pdf
[44]http://demo.spraakdata.gu.se/svedk/pbl/kokkinakisBNER.pdf
[45]http://universaldependencies.org/u/pos/

# 3 THE SPARV FRONTEND

## 3.1 SETTING UP THE SPARV FRONTEND

This section describes how to get the Sparv frontend up and running on your own machine. Begin by downloading the latest source zip from here[46]. The source code is made available under the MIT license[47].

### 3.1.1 REQUIREMENTS

- A web server, such as Apache[48] or equivalent.

### 3.1.2 INSTALLATION

The frontend comes pre-compiled in the zip and can be installed by simply dropping the `dist` folder into your web server root. Consult your web server documentation for more details. You will also need to set some configuration variables in `dist/config.js` in order to point your frontend to the correct backend address and to the address of the default JSON schema file.

## 3.2 DEVELOPING THE SPARV FRONTEND

This section describes the internal structure of the Sparv frontend and can be used as a developer's guide.

The frontend is available at https://spraakbanken.gu.se/sparv/.

### 3.2.1 REQUIREMENTS

- Node v0.10.x, which should be installed through you package manager. This should by default include NPM, the node package manager.
- Grunt[49], install using `npm install -g grunt-cli`.
- CoffeeScript, install using `npm install -g coffeescript`.
- Sass, install using `npm install sass`.

### 3.2.2 FRONTEND CONFIGURATION

The `app/config.js` file contains the configuration of the backend address, the address to the default settings JSON schema and also the address to Karp[50].

### 3.2.3 RUNNING THE FRONTEND

For running the frontend locally (while developing) run `grunt serve`. In you browser, open `http://localhost:9010` to launch Sparv. While running `grunt serve` the CoffeeScript and Sass files are automatically compiled upon edit, additionally causing the browser window to be reloaded to reflect the new changes.

Before releasing a new version, the scripts are compiled by running `grunt` in the frontend directory. This will create a directory called `dist` which contains all the files necessary to run the frontend.

---

[46]https://spraakbanken.gu.se/pub/sparv.dist/sparv_frontend
[47]https://opensource.org/licenses/MIT
[48]http://httpd.apache.org/download.cgi
[49]http://gruntjs.com/
[50]https://spraakbanken.gu.se/karp/

# 4 THE SPARV BACKEND

## 4.1 SETTING UP THE SPARV BACKEND

This section describes how to get the Sparv backend up and running on your own machine. Begin by downloading the latest source zip from here[51].

Note that the source code is made available in two different versions:

- under the MIT license[52] without usage of FreeLing
- under the AGPL license[53] including calls to FreeLing

### 4.1.1 INFO

Here follow instructions for installation on a UNIX-like environment. For more information on the different configuration variables check the developer's guide[54].

The Sparv backend is configured to be run in combination with the catapult. The catapult source files are distributed together with the Sparv pipeline (see below).

### 4.1.2 REQUIREMENTS

- The Sparv corpus pipeline (see here[55] for installation instructions)
- Python 3.4[56] or newer
- A WSGI server (optional but recommended)
- GCC[57] for compiling the `catapult` C extension

### 4.1.3 INSTALLATION

- Set up a Python virtual environment for the backend and install the requirements from `backend/html/app/requirements.txt`.
- Set the backend configuration variables in `backend/html/app/config.py`.
- When running the backend with gunicorn (recommended) set the gunicorn configuration in `backend/html/app/gunicorn_config.py`.
- Set up the catapult, for instance in `backend/data/catapult`.
- Set up a Python virtual environment for the catapult and install the requirements from `catapult/requirements.txt`.
- Set the catapult configuration variables in `catapult/config.sh`.
- From within the `catapult` directory run `make` to build `catalaunch`.
- Set up the Sparv pipeline, for instance in `backend/data/pipeline` and build the pipeline models. You will not need to setup a virtual environment for the pipeline. Set the `VENV_PATH` to the catapult virtual environment.
- Start the catapult by running `./start-server.sh` (from within the `catapult` directory).
- Run the script `index.py` in the `backend` directory. This can be done by running it directly from the python interpreter or by starting a WSGI server using gunicorn (recommended): `html/app/venv/bin/gunicorn -c html/app/gunicorn_config.py index`.
- Set up the cron jobs listed in `catapult/cronjobs` for the automatic maintenance of Sparv.

## 4.2 DEVELOPING THE SPARV BACKEND

This section describes the internal structure of the Sparv backend and its usage of the catapult and can be used as a developer's guide.

---

[51] https://spraakbanken.gu.se/lb/pub/sparv.dist/sparv_backend/?C=M;O=D
[52] https://opensource.org/licenses/MIT
[53] http://www.gnu.org/licenses/agpl.html
[54] https://spraakbanken.gu.se/eng/research/infrastructure/sparv/developersguides
[55] https://spraakbanken.gu.se/eng/research/infrastructure/sparv/distribution/pipeline
[56] http://python.org/
[57] http://gcc.gnu.org/install

### 4.2.1 BACKEND

The backend is run through the WSGI script `index.py`. It is available at [https://spraakbanken.gu.se/ws/sparv](https://spraakbanken.gu.se/ws/sparv).

#### 4.2.1.1 REQUIREMENTS

- the Sparv corpus pipeline, see here[58] for installation instructions
- Python 3.4[59] or newer

#### 4.2.1.2 PYTHON VIRTUAL ENVIRONMENT

Though it is not required, we recommend that you use a Python virtual environment to run the Sparv backend. This is the easiest way to ensure that you have all the Python dependencies needed to run the modules.

Set up a Python virtual environment as a subdirectory to the `backend/html/app` directory:

```
python3 -m venv venv
```

Activate the virtual environment and install the required Python packages:

```
source venv/bin/activate
pip install -r requirements.txt
```

You can then deactivate the virtual environment:

```
deactivate
```

#### 4.2.1.3 BACKEND CONFIGURATION

The configuration variables are stored in `html/app/config.py`:

- `backend`: the address where the backend is hosted
- `sparv_python`: the Python path to the sparv pipeline python directory
- `sparv_backend`: the path to the backend directory
- `builds_dir`: the directory that hosts running and completed builds
- `log_dir`: location of the log files. Can be set to `None` to log to stdout.
- `sparv_models`: location of the sparv pipeline models
- `sparv_makefiles`: location of the sparv pipeline makefiles
- `secret_key`: a string of your choosing. It is needed for queries that may cause deletions of builds.
- `venv_path`: path to the activation script of the Python virtual environment (may be set to None)
- `processes`: number of processes that `make` will run while annotating
- `fileupload_ext`: extension used for builds that contain file uploads
- `socket_file`: path to the socket file used to communicate with the catapult
- `catalaunch_binary`: path to the catalaunch binary file
- `python_interpreter`: "Python" interpreter, replaced with catalaunch

When running the backend with gunicorn (recommended) you may also want to modify the configuration in `html/app/gunicorn_config.py`.

The suggested location for the Sparv pipeline is the directory `data/pipeline`.

---

[58][https://spraakbanken.gu.se/eng/research/infrastructure/sparv/distribution/pipeline](https://spraakbanken.gu.se/eng/research/infrastructure/sparv/distribution/pipeline)
[59][https://www.python.org/](https://www.python.org/)

#### 4.2.1.4 RUNNING THE BACKEND

The backend is set up to be run with gunicorn which is installed automatically inside the Python virtual environment. From the backend directory you can run the following command:

```
html/app/venv/bin/gunicorn -c html/app/gunicorn_config.py index
```

This will start a WSGI server and bind it to the socket defined in `gunicorn_config.py`. Log messages are written to the file specified in the same file (or to the terminal if nothing is specified).

The backend can also be started by running `index.py` with the Python interpreter but this is mostly used for development or debugging.

#### 4.2.1.5 MAKEFILE AND SETTINGS JSON SCHEMA

The makefile for each corpus is created from a JSON object that is created by the script `html/app/schema_generator.py`. The frontend builds its form based on the requested schema. New entries can be added and the frontend should render them automatically. The file that creates the makefile is `html/app/make_makefile.py`.

### 4.2.2 CATAPULT

The catapult runs a Python instance that shares the loaded lexicons, keeps malt processes running and lowers the Python interpreter startup time. Scripts are run on the catapult with the tiny c program `catalaunch`.

#### 4.2.2.1 REQUIREMENTS

- GCC[60] for compiling the `catapult` C extension
- Python 3.4[61] or newer

#### 4.2.2.2 CATAPULT SETUP

The catapult can for example be placed inside the `data` directory.

Just as for the Sparv backend we recomment that you use a Python virtual environment for running the catapult. Check the backend requirements for instructions. The standard location for installing the virtual environment for the catapult is inside the `catapult` directory.

After setting up the Python virtual environment you need to adapt the variable `VENV_PATH` in the Sparv pipeline in `/makefiles/Makefile.config` so it points to the catapult virtual environment.

#### 4.2.2.3 CATAPULT CONFIGURATION

The configuration variables are stored in `config.sh`:

- `SPARV_PYTHON`: the path to the `sb` pipeline Python directory
- `SPARV_MODELS`: location of the `sb` pipeline models
- `SPARV_BIN`: location of the `sb` pipeline binaries
- `SPARV_MAKEFILES`: location of the pipeline makefiles
- `CATAPULT_DIR`: location of the catapult directory
- `BUILDS_DIR`: the directory that hosts running and completed builds
- `LOGDIR`: the path to the log file directory
- `CATAPULT_VENV`: the path to the Python virtual environment used by the catapult

---

[60]http://gcc.gnu.org/install
[61]https://www.python.org/

#### 4.2.2.4 Running the Catapult

- Run `make` to build `catalaunch`.
- Run `./start-server.sh` to start the catapult.
- Set up the cron jobs listed in `catapult/cronjobs.` for the automatic maintenance of Sparv.

#### 4.2.2.5 Cron jobs

The following cron jobs are used in Sparv:

- *Cleanup*: Builds that have not been accessed for 7 days are removed every midnight by issuing `https://ws.spraakbanken.gu.se/ws/sparv/cleanup?secret_key=SECRETKEY`.
- *Keep-alive*: The script `catapult/keep-alive.sh` is run every five minutes and restarts the catapult with `catapult/start-server.sh` if it does not respond to ping. Instead of setting up this cron job you can run the catapult using a process control system like supervisord.
- *Update-saldo*: The Saldo lexicon is updated daily with the script `catapult/update-saldo.sh`. This takes some time, and is therefore run during the night. The catapult is restarted afterwards by the `keep-alive.sh` script.

# 5 THE SPARV WEB API

Sparv has an API which is available at the following address:

> https://ws.spraakbanken.gu.se/ws/sparv

This documentation is for API version 2.

## 5.1 QUERIES FOR ANNOTATING TEXTS

Queries to the web service can be sent as a simple GET request:

> https://ws.spraakbanken.gu.se/ws/sparv/?text=En+exempelmening+till+nä
> ttjänsten

The response is an XML-objekt which for the above request looks like this:

```
<result>
<build hash="edb2a4717701a65b721d97834a56b129dc3bf6aa"/>
<corpus link="https://ws.spraakbanken.gu.se/ws/sparv/download?hash=
    edb2a4717701a65b721d97834a56b129dc3bf6aa">
<text lix="54.00" ovix="inf" nk="inf">
<paragraph>
<sentence id="8f7-84d">
<w pos="DT" msd="DT.UTR.SIN.IND" lemma="|en|" lex="|en..al.1|"
    sense="|den..1:-1.000|en..2:-1.000|" complemgram="|" compwf="|"
    sentiment="0.6799" ref="1" dephead="2" deprel="DT">En</w>
<w pos="NN" msd="NN.UTR.SIN.IND.NOM" lemma="|exempelmening|" lex
    ="|" sense="|" complemgram="|exempel..nn.1+mening..nn.1:7.044e
    -09|" compwf="|exempel+mening|" ref="2" deprel="ROOT">
    exempelmening</w>
<w pos="PP" msd="PP" lemma="|till|" lex="|till..pp.1|" sense="|till
    ..1:-1.000|" complemgram="|" compwf="|" sentiment="0.5086" ref
    ="3" dephead="2" deprel="ET">till</w>
<w pos="NN" msd="NN.UTR.SIN.DEF.NOM" lemma="|nättjänst|nättjänsten
    |" lex="|" sense="|" complemgram="|nät..nn.1+tjänst..nn.2:6.813e
    -11|nät..nn.1+tjänst..nn.1:6.813e-11|nätt..av.1+tjänst..nn
    .1:1.039e-12|nätt..av.1+tjänst..nn.2:1.039e-12|nät..nn.1+tjäna..
    vb.1+sten..nn.1:1.047e-26|nät..nn.1+tjäna..vb.1+sten..nn.2:1.047
    e-26|nätt..av.1+tjäna..vb.1+sten..nn.2:2.120e-27|nätt..av.1+tjä
    na..vb.1+sten..nn.1:2.120e-27|" compwf="|nät+tjänsten|nätt+tjä
    nsten|nät+tjän+sten|nätt+tjän+sten|" ref="4" dephead="3" deprel
    ="PA">nättjänsten</w>
</sentence>
</paragraph>
</text>
</corpus>
</result>
```

POST requests are also supported using the same address. This can be useful for longer texts. Here is an example using curl:

```
curl -X POST --data-binary text="En exempelmening till nättjänsten"
https://ws.spraakbanken.gu.se/ws/sparv/
```

The response is the same as for the above GET request.

It is also possible to upload text or XML files using curl:

```
curl -X POST -F files[]=@/path/to/file/myfile.txt https://ws.spraakbanken
.gu.se/ws/sparv/upload?
```

The response will be a download link to a zip file containing the annotation.

## 5.2 SETTINGS

The web service supports some costum settings, e.g. it lets you chose between different tokenizers on word, sentence, and paragraph level and you can define which annotations should be generated.

These settings can be provided as a JSON object to the `settings` variable. This object must satisfy the JSON schema available at the following adress:

```
https://ws.spraakbanken.gu.se/ws/sparv/schema
```

The schema holds default values for all the attributes. The use of the settings variable is therefore optional.

A request which only generates a dependency analysis could look like this:

```
https://ws.spraakbanken.gu.se/ws/sparv/?text=Det+trodde+jag+aldrig.&
settings={"positional_attributes":{"dependency_attributes":["ref","
dephead","deprel"],"lexical_attributes":[],"compound_attributes":[]}}
```

Or as a POST request using curl:

```
curl -X POST -g --data-binary text="Det trodde jag aldrig." 'https://
ws.spraakbanken.gu.se/ws/sparv/?settings={"positional_attributes":{"
dependency_attributes":["ref","dephead","deprel"],"lexical_attributes
":[],"compound_attributes":[]}}'
```

If you are not sure how to define the settings variable you can get help from the frontend[62] by clicking `Show JSON Settings` under `Show advanced settings`. This will generate the JSON-objekt for the chosen settings which is sent in the `settings` variable.

The makefile which is generated for a certain set of settings can be viewed by sending a `makefile` query:

```
https://ws.spraakbanken.gu.se/ws/sparv/makefile?settings={"positional_attributes
":{"dependency_attributes":["ref","dephead","deprel"],"lexical_attributes
":[],"compound_attributes":[]}}
```

## 5.3 JOINING A BUILD

At the top of the XML response you can find a hash number inside the `build`-tag. This hash can be used to join a build of an earlier build.

The following request is used for joining the build from the first example of this documentation:

```
https://ws.spraakbanken.gu.se/ws/sparv/join?hash=edb2a4717701a65b721d97834a56b129dc3bf6aa
```

The response contains the chosen settings, the original text and of course the result of the annotation. The entire response looks like this:

---

[62]https://spraakbanken.gu.se/sparv

```
<result>
    <settings>{
      "root": {
        "attributes": [],
        "tag": "text"
      },
      "text_attributes": {
        "readability_metrics": [
          "lix",
          "ovix",
          "nk"
        ]
      },
      "word_segmenter": "default_tokenizer",
      "positional_attributes": {
        "lexical_attributes": [
          "pos",
          "msd",
          "lemma",
          "lex",
          "sense"
        ],
        "compound_attributes": [
          "complemgram",
          "compwf"
        ],
        "dependency_attributes": [
          "ref",
          "dephead",
          "deprel"
        ],
        "sentiment": [
          "sentiment"
        ]
      },
      "sentence_segmentation": {
        "sentence_chunk": "paragraph",
        "sentence_segmenter": "default_tokenizer"
      },
      "paragraph_segmentation": {
        "paragraph_segmenter": "blanklines"
      },
      "lang": "sv",
      "textmode": "plain",
      "named_entity_recognition": [],
      "corpus": "untitled"
}</settings>
```

```
<original>&lt;text&gt;En exempelmening till nättjänsten&lt;/text&gt
    ;</original>
<build hash='edb2a4717701a65b721d97834a56b129dc3bf6aa'/>
<corpus link='https://ws.spraakbanken.gu.se/ws/sparv/download?hash=
    edb2a4717701a65b721d97834a56b129dc3bf6aa'>
    <text lix="54.00" ovix="inf" nk="inf">
        <paragraph>
            <sentence id="8f7-84d">
                <w pos="DT" msd="DT.UTR.SIN.IND" lemma="|en|" lex
                    ="|en..al.1|" sense="|den..1:-1.000|en
                    ..2:-1.000|" complemgram="|" compwf="|"
                    sentiment="0.6799" ref="1" dephead="2" deprel="
                    DT">En</w>
                <w pos="NN" msd="NN.UTR.SIN.IND.NOM" lemma="|
                    exempelmening|" lex="|" sense="|" complemgram="|
                    exempel..nn.1+mening..nn.1:7.044e-09|" compwf="|
                    exempel+mening|" ref="2" deprel="ROOT">
                    exempelmening</w>
                <w pos="PP" msd="PP" lemma="|till|" lex="|till..pp
                    .1|" sense="|till..1:-1.000|" complemgram="|"
                    compwf="|" sentiment="0.5086" ref="3" dephead
                    ="2" deprel="ET">till</w>
                <w pos="NN" msd="NN.UTR.SIN.DEF.NOM" lemma="|nättjä
                    nst|nättjänsten|" lex="|" sense="|" complemgram
                    ="|nät..nn.1+tjänst..nn.2:6.813e-11|nät..nn.1+tj
                    änst..nn.1:6.813e-11|nätt..av.1+tjänst..nn
                    .1:1.039e-12|nätt..av.1+tjänst..nn.2:1.039e-12|n
                    ät..nn.1+tjäna..vb.1+sten..nn.1:1.047e-26|nät..
                    nn.1+tjäna..vb.1+sten..nn.2:1.047e-26|nätt..av
                    .1+tjäna..vb.1+sten..nn.2:2.120e-27|nätt..av.1+
                    tjäna..vb.1+sten..nn.1:2.120e-27|" compwf="|nät+
                    tjänsten|nätt+tjänsten|nät+tjän+sten|nätt+tjän+
                    sten|" ref="4" dephead="3" deprel="PA">nättjä
                    nsten</w>
            </sentence>
        </paragraph>
    </text>
</corpus>
</result>
```

## 5.4 ANALYSING TEXTS IN OTHER LANGUAGES

The default analysis language is Swedish but Sparv also supports other languages. The language is specified through the lang attribute in the settings variable.

The following languages are currently supported:

- Bulgarian (bg)
- Dutch (nl)
- English (en)
- Estonian (et)
- Finish (fi)
- French (fr)
- German (de)
- Italian (it)
- Latin (la)

- Polish (pl)
- Portuguese (pt)
- Russian (ru)
- Slovak (sk)
- Spanish (es)

This is an example of an analysis of a German sentence:

```
https://ws.spraakbanken.gu.se/ws/sparv/?text=Nun+folgt+ein+deutscher+
Beispielsatz.&settings={"lang":"de"}
```

Different kinds of settings are supported for different languages. Please use the frontend[63] if you want to check which options there are for a certain language.

## 5.5 INCREMENTAL PROGRESS INFORMATION

By adding the flag `incremental=true` to your usual query you can receive more information on how your analysis is being processed. An example query could look like this:

```
https://ws.spraakbanken.gu.se/ws/sparv/?text=Nu+med+inkrementell+
information&incremental=true
```

The resulting XML will contain the following extra tags:

```
<increment command="" step="0" steps="27"/>
<increment command="sb.segment" step="1" steps="27"/>
<increment command="sb.segment" step="2" steps="27"/>
<increment command="sb.number --position" step="3" steps="27"/>
<increment command="sb.segment" step="4" steps="27"/>
<increment command="sb.annotate --span_as_value" step="5" steps
    ="27"/>
<increment command="sb.annotate --text_spans" step="6" steps="27"/>
<increment command="sb.parent --children" step="7" steps="27"/>
<increment command="sb.parent --parents" step="8" steps="27"/>
<increment command="sb.parent --parents" step="9" steps="27"/>
<increment command="sb.parent --parents" step="10" steps="27"/>
<increment command="sb.number --position" step="11" steps="27"/>
<increment command="sb.hunpos" step="12" steps="27"/>
<increment command="sb.number --relative" step="13" steps="27"/>
<increment command="sb.annotate --select" step="14" steps="27"/>
<increment command="sb.saldo" step="15" steps="27"/>
<increment command="sb.readability --lix" step="16" steps="27"/>
<increment command="sb.readability --ovix" step="17" steps="27"/>
<increment command="sb.readability --nominal_ratio" step="18" steps
    ="27"/>
<increment command="sb.compound" step="19" steps="27"/>
<increment command="sb.wsd" step="20" steps="27"/>
<increment command="sb.malt" step="21" steps="27"/>
<increment command="sb.sentiment" step="22" steps="27"/>
<increment command="sb.annotate --select" step="23" steps="27"/>
<increment command="sb.annotate --select" step="24" steps="27"/>
<increment command="sb.annotate --chain" step="25" steps="27"/>
<increment command="sb.cwb --export" step="26" steps="27"/>
<increment command="" step="27" steps="27"/>
```

This variable can of course be combined with `settings`, `/join`, and with POST requests.

---

[63]https://spraakbanken.gu.se/sparv

## 5.6 OTHER CALLS

You can get a short description of all existing API calls from the following address:

```
https://ws.spraakbanken.gu.se/ws/sparv/api
```

# 6 INTERACTION BETWEEN THE SPARV COMPONENTS

The following image illustrates how the components involved in the Sparv web interface interact with each other and with the user, both before and during the analysis.
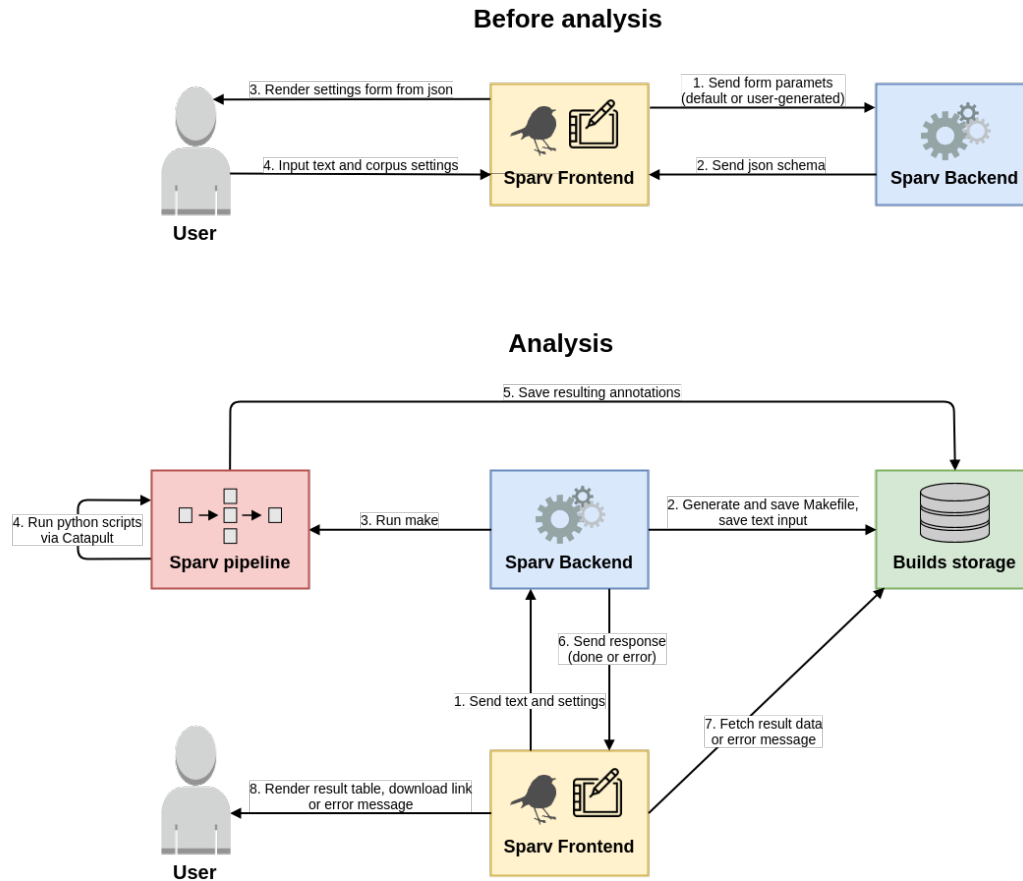
**Before analysis**



**Analysis**



Figure 1: Interaction between the Sparv components