

Bilingual NLP pipeline

Paul Huygen <paul.huygen@huygen.nl>

18th February 2016
10:54 h.

Abstract

This is a description and documentation of the installation of an instrument to annotate Dutch or English documents with NLP tags.

Contents

1	Introduction	1
1.1	List of the modules to be installed	1
1.2	The things that are not open-source yet	3
1.3	Multi-linguality	3
1.4	File-structure of the pipeline	3
2	How to obtain modules and other material	5
2.1	Location-dependency	5
2.2	Reversible update	5
2.3	Installation from Github	6
2.4	Installation from the snapshot	6
3	Java and Python environment	7
3.1	Java	8
3.2	Maven	9
3.3	Java 1.6	9
3.4	Python	10
3.4.1	Virtual environment	11
3.4.2	Transplant the virtual environment	12
3.4.3	KafNafParserPy	12
3.4.4	Python packages	13
4	Installation of the modules	13
4.1	Conditional installation of the modules	13
4.2	The installation script	14
4.3	Check availability of resources	19
4.4	Install utilities and resources	19
4.4.1	Process synchronisation	19
4.4.2	Prefix of scripts that run modules	20
4.4.3	Language detection	20
4.4.4	Alpino	21
4.4.5	Treetagger	22
4.4.6	Timbl and Ticcutils	23
4.4.7	The Boost library	24

4.4.8	Spotlight	25
4.4.9	VUA-pylib	29
4.4.10	SVMLight	29
4.4.11	CRFsuite	29
4.5	Install modules	30
4.5.1	Install tokenizer	30
4.5.2	Topic analyser	31
4.5.3	Morphosyntactic parser	31
4.5.4	Pos tagger	32
4.5.5	Constituent parser	32
4.5.6	NED-reranker	33
4.5.7	Wikify module	33
4.5.8	UKB	34
4.5.9	IMS-WSD	34
4.5.10	SRL server	35
4.5.11	SRL Dutch nominals	36
4.5.12	FBK-time module	37
4.5.13	FBK-temprel module	39
4.5.14	FBK-causalrel module	40
4.5.15	Factuality module	40
4.5.16	Nominal coreference-base	41
4.5.17	Named entity recognition (NERC)	42
4.5.18	Wordsense-disambiguation	43
4.5.19	Lexical-unit converter	44
4.5.20	NED	45
4.5.21	Ontotagger	46
4.5.22	Framenet SRL	47
4.5.23	Heideltime	48
4.5.24	Semantic Role labelling	50
4.5.25	SRL postprocessing	51
4.5.26	Event coreference	52
4.5.27	Dbpedia-ner	53
4.5.28	Nominal events	53
4.5.29	Opinion miner	54
5	Utilities	55
5.1	Test script	55
5.2	Logging	57
5.3	Misc	57
A	How to read and translate this document	58
A.1	Read this document	58
A.2	Process the document	58
A.3	The Makefile for this project.	59
A.4	Get Nuweb	60
A.5	Pre-processing	61
A.5.1	Process ‘dollar’ characters	61
A.5.2	Run the M4 pre-processor	61
A.6	Typeset this document	62
A.6.1	Figures	62
A.6.2	Bibliography	63
A.6.3	Create a printable/viewable document	63
A.6.4	Create HTML files	66
A.7	Create the program sources	69

A.8 Restore paths after transplantation	70
B References	71
B.1 Literature	71
C Indexes	71
C.1 Filenames	71
C.2 Macro's	72
C.3 Variables	74

1 Introduction

This document describes the current set-up of a pipeline that annotates texts in order to extract knowledge. The pipeline has been set up by the Computational Lexicology and Terminology Lab (CLTL¹) as part of the newsreader² project. It accepts and produces texts in the NAF (Newsreader Annotation Format) format.

Apart from describing the pipeline set-up, the document actually constructs the pipeline. The pipeline has been installed on a (Ubuntu) Linux computer.

The installation has been parameterised. The locations and names that you read (and that will be used to build the pipeline) have been read from variables in file `inst.m4` in the `nuweb` directory.

The pipeline is bi-lingual. It is capable to annotate Dutch and English texts. It recognizes the language from the “lang” attribute of the NAF element of the document.

The aim is, to install the pipeline from open-source modules that can e.g. be obtained from Github. However, that aim is only partially fulfilled. Some of the modules still contain elements that are not open-source or data that are not freely available. Because of lack of time, the current version of the installer installs the English pipeline from a frozen repository of the Newsreader Project.

1.1 List of the modules to be installed

Table 1 lists the modules that are installed. Some of the modules are used for both languages (Dutch and English), some for only one of them.

Table 2 lists the modules in the pipeline. The column *source* indicates the origin of the module. The modules are obtained in one of the following ways:

1. If possible, the module is directly obtained from an open-source repository like Github.
2. Some modules have not been officially published in a repository. These modules have been packed in a tar-ball that can be obtained by the author. In table 2 this has been indicated as SNAPSHOT.

The modules themselves use other utilities like dependency-taggers and POS taggers. These utilities are listed in table 3.

1.2 The things that are not open-source yet

The aim is, that the pipeline-system is completely open-sourced, so that anybody can install it from sources like Github. However, a lot of elements are not yet open-sourced, but need private kludges. The following is a list of not-yet open things.

-
1. <http://wordpress.let.vu.nl>
 2. <http://www.newsreader-project.eu>

Module	NL	EN	EN component
Tokenizer	ixa-pipe-tok	ixa-pipe-tok	
Topic detection		ixa-pipe-topic	EHU-topic.v30
POS/MOR	morphosyntactic_parser_nl	EHU-pos.v30	EHU-pos.v30
Constit. parser		ixa-pipe-parse	EHU-parse.v30
NERC	ixa-pipe-nerc	ixa-pipe-nerc	
UKB		UKB	EHU-ukb.v30
WSD	svm_wsd	ims-wsd	VUA-ims-wsd.v30
NED	ixa-pipe-ned	ixa-pipe-ned	
Heideltime	ixa-pipe-time		
FBK-time		FBK-time.v30	FBK-time.v30
FBK-temprel		FBK-temprel.v30	FBK-temprel.v30
FBK-causalrel		FBK-causalrel.v30	FBK-causalrel.v30
Onto-tagger	onto-tagger		
SRL	vua-srl-nl	EHU-srl-server	\verbEHU-srl-server
Nominal event det.	nominal-event-detection		
NED-reranker		domain_model	VUA-popen-nedreranker.v30
Wikify		ixa-pipe-wikify	EHU-wikify.v30
factuality			VUA-factuality.v30
Corefgraph			EHU-corefgraph.v30
Opinion-miner	opinion-miner	opinion-miner	
Eventcoref	vua-eventcoreference_v2	vua-eventcoreference_v2	

Table 1: List of modules to be installed. **Module**: functional name of the module; **NL**: modules used in Dutch pipeline
; **EN**: modules used in English pipeline; **EN components** Name of the module in the EHU repository.

Module	Section	Source	Commit	Script
Tokenizer	4.5.1	Github	56f83ce4b61680346f15e5d4e6de6293764f7383	tok
morphosyntactic parser	4.5.3	Github	807e938ce4ebb71afd9d7c7f42d9d9ac5f98a184	mor
NERC	4.5.17	Gith./snap	ca02c931bc0b200ccdb8b5795a7552e4cc0d4802	nerc
WSD	4.5.18	Gith./snap	030043903b42f77cd20a9b2443de137e2efe8513	wsd
Onto-tagger	4.5.21	snapshot		onto
Heideltime	4.5.23	Gith./snap.	da4604a7b33975e977017440cbc10f7d59917ddf	heidelttime
SRL	4.5.24	Github	675d22d361289ede23df11dcdb17195f008c54bf	srl
SRL-POST	4.5.25	snapshot		postsrl
NED	4.5.20	Github	d35d4df5cb71940bf642bb1a83e2b5b7584010df	ned
Nom. coref	4.5.16	Github	bfa5aec0fa498e57fe14dd4d2c51365dd09a0757	nomcoref
Ev. coref	4.5.26	snapshot		evcoref
Opinion miner	4.5.29	Github		opinimin
Framenet SRL	4.5.22	snapshot		fsrl
Dbpedia_ner	4.5.27	Github	ab1dcdb860f0ff29bc979f646dc382122a101fc2	dbpner

Table 2: List of the modules to be installed. Column description: **directory**: Name of the subdirectory below subdirectory **modules** in which it is installed; **source**: From where the module has been obtained; **commit**: Commit-name or version-tag **script**: Script to be included in a pipeline. **Note**: The tokenizer module has been temporarily obtained from the snapshot, because the commit that we used has disappeared from the Github repository.

Module	Version	Section	Source
KafNafParserPy	Feb 1, 2015	3.4.3	Github
Alpino	20706	4.4.4	RUG
Ticcutils	0.7	4.4.6	ILK
Timbl	6.4.6	4.4.6	ILK
Treetagger	3.2	4.4.5	Uni. München
Spotlight server	0.7	4.4.8	Spotlight

Table 3: List of the modules to be installed. Column description: **directory**: Name of the subdirectory below *mod* in which it is installed; **Source**: From where the module has been obtained; **script**: Script to be included in a pipeline.

1.3 Multi-linguality

This version of the pipeline is multi-lingual, i.e. it can annotate Dutch as well as English documents. It finds the language of the document in the `language` attribute of the NAF element. Actually, the current version is bi-lingual, because it is only able to process Dutch or English documents.

1.4 File-structure of the pipeline

The files that make up the pipeline are organised in set of directories as shown in figure 1. The

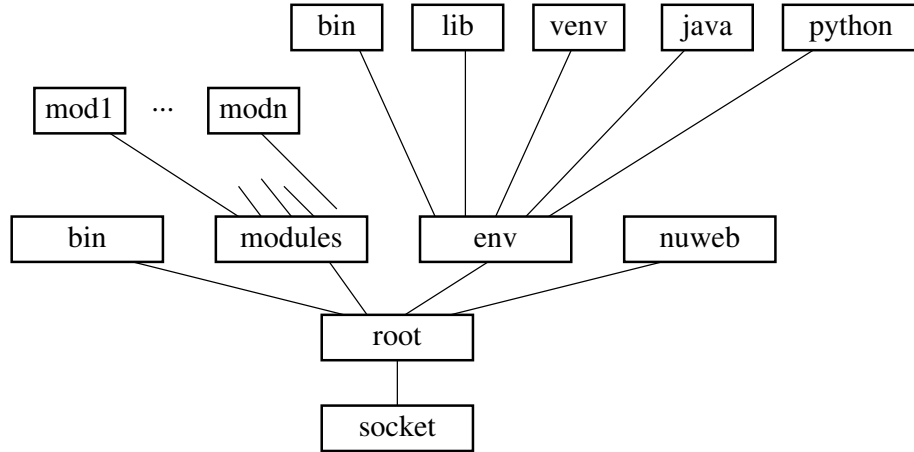


Figure 1: Directory-structure of the pipeline (see text).

directories have the following functions.

socket: The directory in the host where the pipeline is to be implemented.

root: The root of the pipeline directory-structure.

nuweb: This directory contains this document and everything to create the pipeline from the open sources of the modules.

modules: Contains subdirectories with the NLP modules that can be applied in the pipeline.

bin: Contains for each of the applicable modules a script that reads NAF input, passes it to the module in the `modules` directory and produces the output on standard out. Furthermore, the subdirectory contains the script `install-modules` that performs the installation, and a script `test` that shows that the pipeline works in a trivial case.

env: The programming environment. It contains a.o. the Java development kit, Python, the Python virtual environment (`venv`), libraries and binaries.

< directories to create 4a > ≡
../modules ◇

Fragment defined by 4abcd, 8bh, 9a, 11e, 64c.
 Fragment referenced in 70a.

< directories to create 4b > ≡
../bin ../env/bin ◇

Fragment defined by 4abcd, 8bh, 9a, 11e, 64c.
 Fragment referenced in 70a.

< directories to create 4c > ≡
../env/lib ◇

Fragment defined by 4abcd, 8bh, 9a, 11e, 64c.
 Fragment referenced in 70a.

< directories to create 4d > ≡
../env/etc ◇

Fragment defined by 4abcd, 8bh, 9a, 11e, 64c.
 Fragment referenced in 70a.

The following macro defines variable `piperoot` and makes it to point to the root directory in figure 1. Next it defines variables that point to other directories in the figure. The value-setting of `piperoot` can be overruled by defining the variable before running any of the script. In this way the directory tree can be moved to another location, even to another computer, after successful installation.

< set variables that point to the directory-structure 4e > ≡
 if
 ["\$piperoot" == ""]
 then
 export piperoot=/home/phuijgen/nlp/test/nlpp
 fi
 export pipesocket=\${piperoot%/nlpp}
 export nuwebdir=\$piperoot/nuweb
 export envdir=\$piperoot/env
 export envbindir=\$envdir/bin
 export envlibdir=\$envdir/lib
 export modulesdir=\$piperoot/modules
 export pipebin=\$piperoot/bin
 export javadir=\$envdir/java
 export jarsdir=\$javadir/jars
 ◇

Fragment defined by 4e, 5a, 7a, 9c.
 Fragment referenced in 5b, 15a, 71.
 Uses: nuweb 66b.

Add the environment `bin` directory to `PATH`:

```

< set variables that point to the directory-structure 5a > ≡
    export PATH=$envbindir:$PATH
    ◇

```

Fragment defined by 4e, 5a, 7a, 9c.

Fragment referenced in 5b, 15a, 71.

Defines: PATH 8g, 9c, 10a, 41a, 55a.

Put the macro to set variables in a script that can later be sourced by the scripts of the pipeline modules.

```

"../env/bin/progenv" 5b≡
    #!/bin/bash
    < set variables that point to the directory-structure 4e, ... >
    export progenvset=0
    ◇

```

File defined by 5b, 8a.

2 How to obtain modules and other material

As illustrated in tables 2 and 3, most of the modules are obtained as source-code from Github, some of the modules or parts of some modules are downloaded from a snapshot, and some of the utilities are obtained in binary form from the supplier.

This section builds standardised methods to obtain modules and utilities from Github or from the snapshot.

2.1 Location-dependency

The basic way of installation is, to clone this repository from Github on the intended location in the file-system of the target computer and then run the install-scripts. However, it may be advantageous to be able to transplant a complete installation to another location in another computer. This could be done by making all path-descriptions in all scripts relative to anchorpoints within the installation, while it may be hard to find such anchorpoints in advance. Therefore, we take another approach in which we supply a script that repairs paths-descriptions after the transplantation (section A.8).

2.2 Reversible update

This script might be used to update an existing installation. To minimize the risk that the “update” actually ruins an existing installation, move existing modules away before installing the latest version. When the new modules has been installed succesfully, the moved module will be removed. The following macro’s help to achieve this:

```

< move module 5c > ≡
    if
    [ -e @1 ]
    then
        mv @1 old.@1
    fi
    ◇

```

Fragment referenced in 6d, 13a, 57c.

```

< remove old module 6a > ≡
    rm -rf old.@1
    ◇

```

Fragment referenced in 6d, 13a, 57c.

```

< re-instate old module 6b > ≡
    mv old.@1 @1
    MESS="Replaced previous version of @1"
    < logmess (6c $MESS ) 57b >
    ◇

```

Fragment referenced in 6d, 13a, 57c.

2.3 Installation from Github

The following macro can be used to install a module from Github. Before issuing this macro, the following four variables must be set:

MODNAM: Name of the module.

DIRN: Name of the root directory of the module.

GITU: Github URL to clone from.

GITC: Github commit-name or version tag.

```

< install from github 6d > ≡
    cd $modulesdir
    < move module (6e $DIRN ) 5c >
    git clone $GITU
    if
        [ $? -gt 0 ]
    then
        < logmess (6f Cannot install current $MODNAM version ) 57b >
        < re-instate old module (6g $DIRN ) 6b >
    else
        < remove old module (6h $DIRN ) 6a >
        cd $modulesdir/$DIRN
        git checkout $GITC
    fi
    ◇

```

Fragment referenced in 31d, 36c, 41c, 43f, 45c, 48d, 50d, 53a.

2.4 Installation from the snapshot

The sources for the non-open parts of the pipeline are collected in directory `t_nlpp_resources`. They can be accessed via SSH from url `m4_snapshotURL`. Before installing the pipeline download the snapshot on top of directory `snapshotsocket`.


```

< set variables that point to the directory-structure 7a > ≡
  if
    [ ! $snapshotsocket ]
  then
    export snapshotsocket=/home/phuijgen/nlp/test
  fi
  ◇

```

Fragment defined by 4e, 5a, 7a, 9c.

Fragment referenced in 5b, 15a, 71.

The snapshot can be accessed over `scp` on URL `newsreader@kyoto.let.vu.nl`. Access is protected by a public/private key system. So, a private key is needed and this program expects to find the key as `$pipesocket/nrkey`. The key can be obtained from the author. Let us check whether we indeed do have the key:

```

< check this first 7b > ≡
  if
    [ ! -e $pipesocket/nrkey ]
  then
    echo "No key to connect to snapshot!"
    exit 1
  fi
  ◇

```

Fragment defined by 7b, 19e.

Fragment referenced in 15a.

Update the local snapshot repository.

```

< get the snapshot 7c > ≡
  cd $snapshotsocket
  rsync -e "ssh -i $HOME/nrkey" -rLt newsreader@kyoto.let.vu.nl:t_nlpp_resources .
  ◇

```

Fragment referenced in 15a.

3 Java and Python environment

To be independent from the software environment of the host computer and to perform reproducible processing, the pipeline features its own Java and Python environment. The costs of this feature are that the pipeline takes more disk-space by reproducing infra-structure that is already present in the system and that installation takes more time.

The following macro generates a script that specifies the programming environment. Initially it is empty, because we have to create the programming environment first.

```

< create javapython script 7d > ≡
  echo '#!/bin/bash' > /home/phuijgen/nlp/test/nlpp/env/bin/javapython
  ◇

```

Fragment referenced in 15a.

Cause the module scripts to read the javapython script.

```
"../env/bin/progenv" 8a≡
    source $envbindir/javapython
    ◇
```

File defined by 5b, 8a.

3.1 Java

To install Java, download `server-jre-7u72-linux-x64.tar.gz` from <http://www.oracle.com/technetwork/java/javase/downloads/server-jre7-downloads-1931105.html>. Find it in the root directory and unpack it in a subdirectory of `envdir`.

```
< directories to create 8b > ≡
    ../env/java ◇
```

Fragment defined by 4abcd, 8bh, 9a, 11e, 64c.

Fragment referenced in 70a.

```
< set up java 8c > ≡
    < begin conditional install (8d java_installed ) 14b >
        cd $envdir/java
        tar -xzf $snapshotsocket/t_nlpp_resources/server-jre-7u72-linux-x64.tar.gz
    < end conditional install (8e java_installed ) 14d >
    ◇
```

Fragment defined by 8cg.

Fragment referenced in 15a.

Remove the java-ball when cleaning up:

```
< clean up 8f > ≡
    rm -rf $pipesocket/server-jre-7u72-linux-x64.tar.gz
    ◇
```

Fragment defined by 8f, 9d, 22b, 61a.

Fragment referenced in 60a.

Set variables for Java.

```
< set up java 8g > ≡

    echo 'export JAVA_HOME=$envdir/java/jdk1.7.0_72' >> /home/phuijgen/nlp/test/nlpp/env/bin/javapython
    echo 'export PATH=$JAVA_HOME/bin:$PATH' >> /home/phuijgen/nlp/test/nlpp/env/bin/javapython
    export JAVA_HOME=$envdir/java/jdk1.7.0_72
    export PATH=$JAVA_HOME/bin:$PATH
    ◇
```

Fragment defined by 8cg.

Fragment referenced in 15a.

Uses: PATH 5a.

Put jars in the jar subdirectory of the java directory:

```
< directories to create 8h > ≡
    ../env/java/jars ◇
```

Fragment defined by 4abcd, 8bh, 9a, 11e, 64c.

Fragment referenced in 70a.

3.2 Maven

Some Java-based modules can best be compiled with [Maven](#).

```
< directories to create 9a > ≡
  ../env/apache-maven-3.0.5 ◇
```

Fragment defined by [4abcd](#), [8bh](#), [9a](#), [11e](#), [64c](#).

Fragment referenced in [70a](#).

```
< install maven 9b > ≡
  cd $envdir
  wget http://apache.rediris.es/maven/maven-3/3.0.5/binaries/apache-maven-3.0.5-
  bin.tar.gz
  tar -xzf apache-maven-3.0.5-bin.tar.gz
  rm apache-maven-3.0.5-bin.tar.gz
  ◇
```

Fragment referenced in [15a](#).

```
< set variables that point to the directory-structure 9c > ≡
  export MAVEN_HOME=$envdir/apache-maven-3.0.5
  export PATH=${MAVEN_HOME}/bin:${PATH}
  ◇
```

Fragment defined by [4e](#), [5a](#), [7a](#), [9c](#).

Fragment referenced in [5b](#), [15a](#), [71](#).

Uses: [PATH 5a](#).

When the installation has been done, remove maven, because it is no longer needed.

```
< clean up 9d > ≡
  rm -rf ../env/apache-maven-3.0.5
  < remove installed-variable (9e maven_installed ) 14e >
  ◇
```

Fragment defined by [8f](#), [9d](#), [22b](#), [61a](#).

Fragment referenced in [60a](#).

3.3 Java 1.6

Java 1.7 is able to run nearly all the modules of the pipeline that are based on Java. However, there is one exception, i.e. the `ims-wsd` module, that needs Java version 1.6. So, we have to install that version of Java as well.

```
< install Java 1.6 9f > ≡
  cd $envdir/java
  $snapshotsocket/t_nlpp_resources/jre-6u45-linux-x64.bin
  ◇
```

Fragment referenced in [15a](#).

Insert the following macro in scripts that need to run Java 1.6.

```

< set up Java 1.6 10a > ≡
    export JAVA_HOME=$envdir/java/jre1.6.0_45
    export PATH=$JAVA_HOME/bin:$PATH
    ◇

```

Fragment referenced in 35a.
 Uses: PATH 5a.

3.4 Python

Set up the environment for Python (version 2.7). I could not find an easy way to set up Python from scratch. Therefore we will use Python 2.7 if it has been installed on the host. Otherwise, we will use a binary distribution obtained from [ActiveState](#). A tarball of ActivePython can be obtained from the snapshot.

In order to be independent of the software on the host, we generate a virtual Python environment. In the virtual environment we will install KafNafParserPy and other Python packages that are needed.

```

< set up python 10b > ≡
    < check/install the correct version of python 10c >
    < create a virtual environment for Python 11b >
    < activate the python environment 11d, ... >
    < update pip 12b >
    < install kafnafparserpy 13a >
    < install python packages 13f, ... >
    ◇

```

Fragment referenced in 15a.

```

< check/install the correct version of python 10c > ≡
pythonok='python --
version 2>&1 | gawk '{if(match($2, "2.7")) print "yes" ; else print "no" }'
if
    [ "$pythonok" == "no" ]
then
    < install ActivePython 11a >
fi
◇

```

Fragment referenced in 10b.
 Defines: pythonok Never used.
 Uses: print 64a.

Unpack the tarball in a temporary directory and install active python in the `env` subdirectory of `nlpp`. It turns out that you must upgrade `pip`, `virtualenv` and `setuptools` after the installation (see <https://github.com/ActiveState/activepython-docker/commit/10fff72069e51dbd36330cb8a7c2f0845bcd7b3> and <https://github.com/ActiveState/activepython-docker/issues/1>).

```

< install ActivePython 11a > ≡
    pytinsdir='mktemp -d -t activepyt.XXXXXX'
    cd $pytinsdir
    tar -xzf $snapshotsocket/t_nlpp_resources/ActivePython-2.7.8.10-linux-x86_64.tar.gz
    accdir='ls -1'
    cd $accdir
    ./install.sh -I $envdir
    cd $piperooot
    rm -rf $pytinsdir
    pip install -U pip virtualenv setuptools
    ◇

```

Fragment referenced in 10c.

Uses: virtualenv 11c.

3.4.1 Virtual environment

Create a virtual environment. To begin this, we need the Python module virtualenv on the host.

```

< create a virtual environment for Python 11b > ≡
    < test whether virtualenv is present on the host 11c >
    cd $envdir
    virtualenv venv
    ◇

```

Fragment referenced in 10b.

Uses: virtualenv 11c.

```

< test whether virtualenv is present on the host 11c > ≡
    which virtualenv
    if
        [ $? -ne 0 ]
    then
        echo Please install virtualenv
        exit 1
    fi
    ◇

```

Fragment referenced in 11b.

Defines: virtualenv 11ab.

```

< activate the python environment 11d > ≡
    source $envdir/venv/bin/activate
    echo 'source $en-
vdir/venv/bin/activate' >> /home/phuijgen/nlp/test/nlpp/env/bin/javapython
    ◇

```

Fragment defined by 11d, 12a.

Fragment referenced in 10b, 15a.

Defines: activate 12c.

Subdirectory \$envdir/python will contain general Python packages like KafnafParserPy.

```

< directories to create 11e > ≡
    ../env/python ◇

```

Fragment defined by 4abcd, 8bh, 9a, 11e, 64c.

Fragment referenced in 70a.

Activation of Python include pointing to the place where Python packages are:

```
< activate the python environment 12a > ≡
echo ex-
port 'PYTHONPATH=$envdir/python:$PYTHONPATH' >> /home/phuijgen/nlp/test/nlpp/env/bin/javapython
export PYTHONPATH=$envdir/python:$PYTHONPATH
◇
```

Fragment defined by 11d, 12a.
 Fragment referenced in 10b, 15a.
 Defines: PYTHONPATH Never used.

Update pip in the virtual environment, because otherwise it keeps complaining about outdated versions

```
< update pip 12b > ≡
pip install --upgrade pip
◇
```

Fragment referenced in 10b.

3.4.2 Transplant the virtual environment

It turns out that the script “activate” to engage the virtual environment contains an absolute path, in the definition of VIRTUAL_ENV

```
< set paths after transplantation 12c > ≡
transdir='mktemp -d -t trans.XXXXXX'
cd $transdir
cat <<EOF >redef.awk
#!/usr/bin/gawk -f
BEGIN { envd="$envdir/venv"}

/^VIRTUAL_ENV=/ { print "VIRTUAL_ENV=\"\" envd \"\"\"
                  next
                  }

{print}
EOF

mv $envdir/venv/bin/activate .
gawk -f redef.awk ./activate > $envdir/venv/bin/activate
cd $projroot
rm -rf $transdir
◇
```

Fragment referenced in 71.
 Uses: activate 11d, print 64a.

3.4.3 KafNafParserPy

A cornerstone Pythonmodule for the pipeline is [KafNafParserPy](#). It is a feature of this module that you cannot install it with PIP, but that you can add it to your PYTHONPATH.

```

< install kafnafparserpy 13a > ≡
  cd $envdir/python
  DIRN=KafNafParserPy
  < move module (13b $DIRN ) 5c >
  git clone https://github.com/cltl/KafNafParserPy.git
  if
    [ $? -gt 0 ]
  then
    < logmess (13c Cannot install current $DIRN version ) 57b >
    < re-instate old module (13d $DIRN ) 6b >
  else
    < remove old module (13e $DIRN ) 6a >
  fi
  ◇

```

Fragment referenced in 10b.

3.4.4 Python packages

Install python packages:

lxml:

pyyaml: for coreference-graph

```

< install python packages 13f > ≡
  pip install lxml
  pip install pyyaml
  ◇

```

Fragment defined by 13f, 51e.

Fragment referenced in 10b.

Defines: **lxml** Never used, **pyyaml** Never used.

4 Installation of the modules

This section describes how the modules are obtained from their (open-)source and installed.

4.1 Conditional installation of the modules

Next section generates a script that installs everything.

Installation is very time-intensive. To prevent that everything is re-installed every time that the module-installer is run, there is a list of variables, the *modulelist*, that are set when a module has been installed. To re-install that module, remove the variable from the list and then re-run the installer. It maintains a list of the modules and utilities that it has installed and installs only modules and utilities that are not on the list. So in order to re-install a module that has already been installed, remove it from the list and then re-run the module-installer.

The modulelist is in fact a script named `/home/phuijgen/nlp/test/nlpp/installed_modules` that sets Bash variables. It ought to be sourced if it is present.

Initially the list is not present. When a module or a utility has been installed, an instruction to set a variable is written in or appended to the list.

```

< read the list of installed modules 14a > ≡
    if
        [ -e /home/phuijgen/nlp/test/nlpp/installed_modules ]
    then
        source /home/phuijgen/nlp/test/nlpp/installed_modules
    fi
    ◇

```

Fragment referenced in 15a.

```

< begin conditional install 14b > ≡
    if
        [ ! $@1 ]
    then
        ◇

```

Fragment referenced in 8c, 15a, 16aj, 17aj, 18ahq, 19a.

```

< else conditional install 14c > ≡
    else
        ◇

```

Fragment never referenced.

```

< end conditional install 14d > ≡
    echo "export @1=0" >> /home/phuijgen/nlp/test/nlpp/installed_modules
    fi
    ◇

```

Fragment referenced in 8c, 15a, 16aj, 17aj, 18ahq, 19a.

Remove a variable from the list of installed modules, e.g. after a clean-up.

```

< remove installed-variable 14e > ≡
    cd $piperoot
    mv /home/phuijgen/nlp/test/nlpp/installed_modules old.modulelist
    cat old.modulelist | gawk '/@1/ {next}; {print}' >/home/phuijgen/nlp/test/nlpp/installed_modules
    ◇

```

Fragment referenced in 9d.

Uses: `print` 64a.

4.2 The installation script

The installation is performed by script `install-modules`.

The first part of the script installs the utilities:


```

"../bin/install-modules" 15a≡
    #!/bin/bash
    echo Set up environment
    < set variables that point to the directory-structure 4e, ... >
    < read the list of installed modules 14a >
    < check this first 7b, ... >
    < begin conditional install (15b repo_installed) 14b >
        < get the snapshot 7c >
    < end conditional install (15c repo_installed) 14d >
    < variables of install-modules 57a >
    < create javapython script 7d >
    echo ... Java
    < set up java 8c, ... >
    < begin conditional install (15d maven_installed) 14b >
        < install maven 9b >
    < end conditional install (15e maven_installed) 14d >
    < begin conditional install (15f java16_installed) 14b >
        < install Java 1.6 9f >
    < end conditional install (15g java16_installed) 14d >

    echo ... Python
    if
    [ $python_installed ]
    then
        < activate the python environment 11d, ... >
    fi
    < begin conditional install (15h python_installed) 14b >
        < set up python 10b >
    < end conditional install (15i python_installed) 14d >
    < begin conditional install (15j sematree_installed) 14b >
        < install sematree 20a >
    < end conditional install (15k sematree_installed) 14d >
    echo ... Alpino
    < begin conditional install (15l alpino_installed) 14b >
        < install Alpino 21f >
    < end conditional install (15m alpino_installed) 14d >
    echo ... Spotlight
    < begin conditional install (15n spotlight_installed) 14b >
        < install the Spotlight server 25a, ... >
    < end conditional install (15o spotlight_installed) 14d >
    echo ... Treetagger
    < begin conditional install (15p treetagger_installed) 14b >
        < install the treetagger utility 22c, ... >
    < end conditional install (15q treetagger_installed) 14d >
    echo ... Ticcutils and Timbl
    < begin conditional install (15r ticctimbl_installed) 14b >
        < install the ticcutils utility 24a >
        < install the timbl utility 24b >
    < end conditional install (15s ticctimbl_installed) 14d >
    echo ... Boost
    < begin conditional install (15t boost_installed) 14b >
        < install boost 24e >
    < end conditional install (15u boost_installed) 14d >

    echo ... VUA-pylib, SVMlight, CRFsuite
    < begin conditional install (15v miscutils_installed) 14b >
        < install VUA-pylib 29a >
        < install SVMlight 29b >
        < install CRFsuite 30a >
    < end conditional install (15w miscutils_installed) 14d >

```

◇

File defined by 15a, 16aj, 17aj, 18ahq, 19a.

Next, install the modules:

```
"../bin/install-modules" 16a≡
echo Install modules
  <begin conditional install (16b tokenizer_installed) 14b>
    echo ... Tokenizer
    <install the tokenizer 30b>
  <end conditional install (16c tokenizer_installed) 14d>
  <begin conditional install (16d topic_installed) 14b>
    echo ... Topic detector
    <install the topic analyser 31a>
  <end conditional install (16e topic_installed) 14d>
  <begin conditional install (16f morpar_installed) 14b>
    echo ... Morphosyntactic parser
    <install the morphosyntactic parser 31d>
  <end conditional install (16g morpar_installed) 14d>
  <begin conditional install (16h pos_installed) 14b>
    echo "... Pos tagger (for english docs)"
    <install the pos tagger 32b>
  <end conditional install (16i pos_installed) 14d>
◇
```

File defined by 15a, 16aj, 17aj, 18ahq, 19a.

```
"../bin/install-modules" 16j≡
  <begin conditional install (16k constparse_installed) 14b>
    echo "... Constituent parser (for english docs)"
    <install the constituents parser 32e>
  <end conditional install (16l constparse_installed) 14d>
  <begin conditional install (16m nerc_installed) 14b>
    echo ... NERC
    <install the NERC module 42a>
  <end conditional install (16n nerc_installed) 14d>
  <begin conditional install (16o ned_installed) 14b>
    echo ... NED
    <install the NED module 45c>
  <end conditional install (16p ned_installed) 14d>
  <begin conditional install (16q nedrer_installed) 14b>
    echo ...NED reranker
    <install the NED-reranker module 33c>
  <end conditional install (16r nedrer_installed) 14d>
  <begin conditional install (16s wikify_installed) 14b>
    echo ...WIKIfy module
    <install the wikify module 33f>
  <end conditional install (16t wikify_installed) 14d>
◇
```

File defined by 15a, 16aj, 17aj, 18ahq, 19a.

```

"../bin/install-modules" 17a≡
  < begin conditional install (17b UKB_installed ) 14b >
    echo ... UKB module
    cd $modulesdir
    tar -xzf $snapshotsocket/t_nlpp_resources/20151220_EHU-ukb.v30.tgz
  < end conditional install (17c UKB_installed ) 14d >
  < begin conditional install (17d ims_wsd_installed ) 14b >
    echo ...ims-wsd module
    < install the ims-wsd module 34d >
  < end conditional install (17e ims_wsd_installed ) 14d >
  < begin conditional install (17f srl_server_installed ) 14b >
    echo ...srl-server module
    < install the srl-server module 35c >
  < end conditional install (17g srl_server_installed ) 14d >
  < begin conditional install (17h srl_dutch_nominals_installed ) 14b >
    echo ...srl-dutch-nominal module
    < install the srl-dutch-nominals module 36c >
  < end conditional install (17i srl_dutch_nominals_installed ) 14d >
  ◇

```

File defined by 15a, 16aj, 17aj, 18ahq, 19a.

```

"../bin/install-modules" 17j≡
  < begin conditional install (17k FBK_time_installed ) 14b >
    echo ... FBK-time module
    < install the FBK-time module 37 >
  < end conditional install (17l FBK_time_installed ) 14d >
  < begin conditional install (17m FBK_temprel_installed ) 14b >
    echo ... FBK-temprel module
    < install the FBK-temprel module 39a >
  < end conditional install (17n FBK_temprel_installed ) 14d >
  < begin conditional install (17o FBK_causalrel_installed ) 14b >
    echo ... FBK-causalrel module
    < install the FBK-causalrel module 40a >
  < end conditional install (17p FBK_causalrel_installed ) 14d >
  < begin conditional install (17q factuality_installed ) 14b >
    echo ... factuality module
    < install the factuality module 40f >
  < end conditional install (17r factuality_installed ) 14d >
  ◇

```

File defined by 15a, 16aj, 17aj, 18ahq, 19a.

```

"../bin/install-modules" 18a≡
  < begin conditional install (18b corefb_installed ) 14b >
    echo ... Coreference base
    < install coreference-base 41c >
  < end conditional install (18c corefb_installed ) 14d >
  < begin conditional install (18d wsd_installed ) 14b >
    echo ... WSD
    < install the WSD module 43f >
  < end conditional install (18e wsd_installed ) 14d >
  < begin conditional install (18f onto_installed ) 14b >
    echo ... Ontotagger
    < install the onto module 46d >
  < end conditional install (18g onto_installed ) 14d >
  ◇

```

File defined by 15a, 16aj, 17aj, 18ahq, 19a.

```

"../bin/install-modules" 18h≡
  < begin conditional install (18i heidel_installed ) 14b >
    echo ... Heideltime
    < install the heideltime module 48c >
  < end conditional install (18j heidel_installed ) 14d >
  < begin conditional install (18k SRL_installed ) 14b >
    echo ... SRL
    < install the srl module 50d >
  < end conditional install (18l SRL_installed ) 14d >
  < begin conditional install (18m eventcoref_installed ) 14b >
    echo ... Event-coreference
    < install the event-coreference module 52c >
  < end conditional install (18n eventcoref_installed ) 14d >
  < begin conditional install (18o lu2synset_installed ) 14b >
    echo ... lu2synset
    < install the lu2synset converter 44e >
  < end conditional install (18p lu2synset_installed ) 14d >
  ◇

```

File defined by 15a, 16aj, 17aj, 18ahq, 19a.

```

"../bin/install-modules" 18q≡
  < begin conditional install (18r dbpner_installed ) 14b >
    echo ... dbpedia-ner
    < install the dbpedia-ner module 53a >
  < end conditional install (18s dbpner_installed ) 14d >
  < begin conditional install (18t nomevent_installed ) 14b >
    echo ... nominal event
    < install the nomevent module 53d >
  < end conditional install (18u nomevent_installed ) 14d >
  < begin conditional install (18v post_SRL_installed ) 14b >
    echo ... post-SRL
    < install the post-SRL module 51f >
  < end conditional install (18w post_SRL_installed ) 14d >
  ◇

```

File defined by 15a, 16aj, 17aj, 18ahq, 19a.

```

"../bin/install-modules" 19a≡
  < begin conditional install (19b opimin_installed ) 14b>
    echo ... opinion-miner
    < install the opinion-miner 54a, ... >
  < end conditional install (19c opimin_installed ) 14d>

  echo Final
  ◇

```

File defined by 15a, 16aj, 17aj, 18ahq, 19a.

```

< make scripts executable 19d> ≡
  chmod 775 ../bin/install-modules
  ◇

```

Fragment defined by 19d, 21c, 70b.

Fragment referenced in 70c.

4.3 Check availability of resources

Test for some resources that we need and that may not be available on this host.

```

< check this first 19e> ≡
  < check whether mercurial is present 19f>
  ◇

```

Fragment defined by 7b, 19e.

Fragment referenced in 15a.

```

< check whether mercurial is present 19f> ≡
  which hg
  if
    [ $? -ne 0 ]
  then
    echo Please install Mercurial.
    exit 1
  fi
  ◇

```

Fragment referenced in 19e.

Defines: hg 41c.

4.4 Install utilities and resources

4.4.1 Process synchronisation

We will see that we sometimes have to install server-applications. However, it is possible that multiple processes are running pipeline modules in parallel, and then it may occur that two instances of a module try to install the same server-application. Therefore, we must make sure that only one application at a time is able to start the server.

The program `sematree`, found at <http://www.pixelbeat.org/scripts/sematree/> enables to do this. When invoked with argument “acquire”, the name of a “lockfile” and a time to wait (-1 means “wait an indefinite time”), it checks whether the lockfile exists. If that is the case, it either waits or fails. When the lockfile is not (or no longer) present, `sematree` creates the lockfile.

When installing Sematree, set the default directory for lock-files. We set this as a subdirectory of the `env` tree. However, in some cases, notably when running in a node in Lisa, we need a directory on the filesystem of the node itself.

```
< install sematree 20a > ≡
    cat $snapshotsocket/t_nlpp_resources/sematree | \
    sed "s|/var/run|/home/phuijgen/nlp/test/nlpp/env/etc/sematree|g" \
    > $envbindir/sematree
◇
```

Fragment referenced in 15a.

4.4.2 Prefix of scripts that run modules

Each module will be run by a Bash script located in subdirectory `bin`. The start of these scripts will have similar content. Insert the following macro to include this similar content, with the name of the module-directory as argument:

```
< start of module-script 20b > ≡
    #!/bin/bash
    source /home/phuijgen/nlp/test/nlpp/env/bin/progenenv
    export LC_ALL=en_US.UTF-8
    export LANG=en_US.UTF-8
    export LANGUAGE=en_US.UTF-8
    ROOT=$piperoot
    MODDIR=$modulesdir/@1
◇
```

Fragment referenced in 21a, 30c, 31be, 32c, 33adg, 34b, 35adfh, 36d, 38a, 39d, 40d, 41ad, 43bd, 44c, 45a, 46b, 47a, 48a, 50be, 52ad, 53be, 55a.

4.4.3 Language detection

The following script `../env/bin/langdetect.py` discerns the language of a NAF document. If it cannot find that attribute it prints `unknown`. The macro `set the language variable` uses this script to set variable `lang`. All pipeline modules expect that this variable has been set.

```
"../env/bin/langdetect.py" 20c≡
    #!/usr/bin/env python
    # langdetect -- Detect the language of a NAF document.
    #
    import xml.etree.ElementTree as ET
    import sys
    import re
    xmldoc = sys.stdin.read()
    #print xmldoc
    root = ET.fromstring(xmldoc)
    # print root.attrib['lang']
    lang = "unknown"
    for k in root.attrib:
        if re.match(".*lang$", k):
            language = root.attrib[k]
    print language
◇
```

Uses: print 64a.

```
"../bin/langdetect" 21a≡
  < start of module-script (21b ) 20b >
  echo 'cat | python $envbindir/langdetect.py'
  ◇
```

```
< make scripts executable 21c > ≡
  chmod 775 /home/phuijgen/nlp/test/nlpp/bin/langdetect
  ◇
```

Fragment defined by 19d, 21c, 70b.

Fragment referenced in 70c.

```
< set the language variable 21d > ≡
  naflang='cat @1 | /home/phuijgen/nlp/test/nlpp/bin/langdetect'
  export naflang
  ◇
```

Fragment referenced in 55c.

Defines: naflang 21e, 26bg, 28b, 30c, 33g, 42c, 43d, 46b, 52d, 55ac.

Currently, the pipeline understands only English and Dutch. The following macro aborts pipeline processing when the language is not English or Dutch.

```
< abort when the language is not English or Dutch 21e > ≡
  if
    [ ! "$naflang" == 'nl' ] && [ ! "$naflang" == "en" ]
  then
    echo Language of NAF document not set. >&2
    echo Set variable "naflang" to "en" of "nl" and try again. >&2
    echo Aborting ':-(' >&2
    exit 4
  fi
  ◇
```

Fragment referenced in 30c, 31b.

Uses: naflang 21d.

4.4.4 Alpino

Binary versions of Alpino can be obtained from the [official Alpino website](#) of Gertjan van Noord. However, it seems that older versions are not always retained there, or the location of older versions change. Therefore we have a copy in the snapshot.

Module

```

< install Alpino 21f > ≡
  if
    [ ! $alpino_installed ]
  then
    cd $modulesdir
    tar -xzf $snapshotsocket/t_nlpp_resources/Alpino-x86_64-linux-glibc2.5-20706-
    sicstus.tar.gz
    echo "ex-
port alpino_installed=0" >> /home/phuijgen/nlp/test/nlpp/installed_modules
  fi
  ◇

```

Fragment referenced in [15a](#).

Currently, alpino is not used as a pipeline-module on its own, but it is included in other pipeline-modules. Modules that use Alpino should set the following variables:

```

< set alpinohome 22a > ≡
  export ALPINO_HOME=$modulesdir/Alpino
  ◇

```

Fragment referenced in [31e](#).

Defines: ALPINO_HOME Never used.

Remove the tarball when cleaning up:

```

< clean up 22b > ≡
  rm -rf $snapshotsocket/t_nlpp_resources/Alpino-x86_64-linux-glibc2.5-20706-
  sicstus.tar.gz
  ◇

```

Fragment defined by [8f](#), [9d](#), [22b](#), [61a](#).

Fragment referenced in [60a](#).

4.4.5 Treetagger

Installation of Treetagger goes as follows (See [Treetagger's homepage](#)):

1. Download and unpack the Treetagger tarball. This generates the subdirectories `bin`, `cmd` and `doc`
2. Download and unpack the tagger-scripts tarball

The location where Treetagger comes from and the location where it is going to reside:

```

< install the treetagger utility 22c > ≡
  TREETAGDIR=treetagger
  TREETAG_BASIS_URL=http://www.cis.uni-muenchen.de/%7Eschmid/tools/TreeTagger/data/
  TREETAGURL=http://www.cis.uni-muenchen.de/%7Eschmid/tools/TreeTagger/data/
  ◇

```

Fragment defined by [22cde](#), [23abcd](#).

Fragment referenced in [15a](#).

The source tarball, scripts and the installation-script:


```

< install the treetagger utility 22d > ≡
    TREETAGSRC=tree-tagger-linux-3.2.tar.gz
    TREETAGSCRIPTS=tagger-scripts.tar.gz
    TREETAG_INSTALLSCRIPT=install-tagger.sh
    ◇

```

Fragment defined by 22cde, 23abcd.

Fragment referenced in 15a.

Parametersets:

```

< install the treetagger utility 22e > ≡
    DUTCHPARS_UTF_GZ=dutch-par-linux-3.2-utf8.bin.gz
    DUTCH_TAGSET=dutch-tagset.txt
    DUTCHPARS_2_GZ=dutch2-par-linux-3.2-utf8.bin.gz
    ◇

```

Fragment defined by 22cde, 23abcd.

Fragment referenced in 15a.

Download everything in the target directory:

```

< install the treetagger utility 23a > ≡
    mkdir -p $modulesdir/$TREETAGDIR
    cd $modulesdir/$TREETAGDIR
    wget $TREETAGURL/$TREETAGSRC
    wget $TREETAGURL/$TREETAGSCRIPTS
    wget $TREETAGURL/$TREETAG_INSTALLSCRIPT
    wget $TREETAGURL/$DUTCHPARS_UTF_GZ
    wget $TREETAGURL/$DUTCH_TAGSET
    wget $TREETAGURL/$DUTCHPARS_2_GZ
    ◇

```

Fragment defined by 22cde, 23abcd.

Fragment referenced in 15a.

Run the install-script:

```

< install the treetagger utility 23b > ≡
    chmod 775 $TREETAG_INSTALLSCRIPT
    ./ $TREETAG_INSTALLSCRIPT
    ◇

```

Fragment defined by 22cde, 23abcd.

Fragment referenced in 15a.

Make the treetagger utilities available for everybody.

```

< install the treetagger utility 23c > ≡
    chmod -R o+rx $modulesdir/$TREETAGDIR/bin
    chmod -R o+rx $modulesdir/$TREETAGDIR/cmd
    chmod -R o+r $modulesdir/$TREETAGDIR/doc
    chmod -R o+rx $modulesdir/$TREETAGDIR/lib
    ◇

```

Fragment defined by 22cde, 23abcd.

Fragment referenced in 15a.

Remove the tarballs:

```

< install the treetagger utility 23d > ≡
    rm $TREETAGSRC
    rm $TREETAGSCRIPTS
    rm $TREETAG_INSTALLSCRIPT
    rm $DUTCHPARS_UTF_GZ
    rm $DUTCH_TAGSET
    rm $DUTCHPARS_2_GZ
    ◇

```

Fragment defined by 22cde, 23abcd.

Fragment referenced in 15a.

4.4.6 Timbl and Ticcutils

Timbl and Ticcutils are installed from their source-tarballs. The installation is not (yet?) completely reproducibe because it uses the C-compiler that happens to be available on the host. Installation involves:

1. Download the tarball in a temporary directory.
2. Unpack the tarball.
3. cd to the unpacked directory and perform `./configure`, `make` and `make install`. Note the argument that causes the files to be installed in the `lib` and the `bin` sub-directories of the `env` directory.

```

< install the ticcutils utility 24a > ≡
    URL=http://software.ticc.uvt.nl/ticcutils-0.7.tar.gz
    TARB=ticcutils-0.7.tar.gz
    DIR=ticcutils-0.7
    < unpack ticcutils or timbl 24c >
    ◇

```

Fragment referenced in 15a, 24d.

```

< install the timbl utility 24b > ≡
    TARB=timbl-6.4.6.tar.gz
    DIR=timbl-6.4.6
    < unpack ticcutils or timbl 24c >
    ◇

```

Fragment referenced in 15a, 24d.

```

< unpack ticcutils or timbl 24c > ≡
    SUCCES=0
    ticbeldir='mktemp -t -d tickbel.XXXXXX'
    cd $ticbeldir
    tar -xzf $snapshotsocket/t_nlpp_resources/$TARB
    cd $DIR
    ./configure --prefix=$envdir
    make
    make install
    cd $piperoot
    rm -rf $ticbeldir
    ◇

```

Fragment referenced in 24ab.

When the installation has been transplanted, Timbl and Ticcutils have to be re-installed.

```
⟨ re-install modules after the transplantation 24d ⟩ ≡
  ⟨ install the ticcutils utility 24a ⟩
  ⟨ install the timbl utility 24b ⟩
  ◇
```

Fragment referenced in 71.

4.4.7 The Boost library

Theoretically, it is possible to download a tarball with boost from [it's repository](#) and then install it. However, I did not succeed in doing this. Therefore, I ripped the installed boost from Surfsara's Hadoop installation and put it in the `env` dir.

```
⟨ install boost 24e ⟩ ≡
  cd $envdir
  tar -xzf $snapshotsocket/t_nlpp_resources/20160103_boost_1_54_bin.tgz
  ◇
```

Fragment referenced in 15a.

4.4.8 Spotlight

A Spotlight server occupies a lot of memory and we need two of them, one for each language. We may be lucky and have a spotlight server running somewhere. Otherwise we have to install the server ourselves.

Install Spotlight in the way that Itziar Aldabe (<mailto:itziar.aldabe@ehu.es>) described:

The NED module works for English, Spanish, Dutch and Italian. The module returns multiple candidates and correspondences for all the languages. If you want to integrate it in your Dutch or Italian pipeline, you will need:

1. The jar file with the dbpedia-spotlight server. You need the version that Aitor developed in order to correctly use the "candidates" option. You can copy it from the English VM. The jar file name is `dbpedia-spotlight-0.7-jar-with-dependencies-candidates.jar`
2. The Dutch/Italian model for the dbpedia-spotlight. You can download them from: <http://spotlight.sztaki.hu/downloads/>
3. The jar file with the NED module: `ixa-pipe-ned-1.0.jar`. You can copy it from the English VM too.
4. The file: `wikipedia-db.v1.tar.gz`. You can download it from: <http://ixa2.si.ehu.es/ixa-pipes/models/wikipedia-db.v1.tar.gz>. This file contains the required information to do the mappings between the wikipedia-entries. The zip file contains three files: `wikipedia-db`, `wikipedia-db.p` and `wikipedia-db.t`

To start the dbpedia server: Italian server:

```
java -jar -Xmx8g dbpedia-spotlight-0.7-jar-with-dependencies-candidates.jar \
  it http://localhost:2050/rest
```

Dutch server:

```
java -jar -Xmx8g dbpedia-spotlight-0.7-jar-with-dependencies-candidates.jar nl http://localhost:2
```

We set 8Gb for the English server, but the Italian and Dutch Spotlight will require less memory.

So, let us do that:

```

⟨ install the Spotlight server 25a ⟩ ≡
  cd $envdir
  tar -xzf $snapshotssocket/t_nlpp_resources/spotlightnl.tgz
  cd $envdir/spotlight
  ⟨ get spotlight model ball (25b nl.tar.gz ) 25d ⟩
  ⟨ get spotlight model ball (25c en_2+2.tar.gz ) 25d ⟩
  ◇

```

Fragment defined by 25a, 26a.

Fragment referenced in 15a.

```

⟨ get spotlight model ball 25d ⟩ ≡
  wget http://spotlight.sztaki.hu/downloads/archive/2014/@1
  tar -xzf @1
  rm @1
  ◇

```

Fragment referenced in 25a.

We choose to put the Wikipedia database in the spotlight directory.

```

⟨ install the Spotlight server 26a ⟩ ≡
  cd $envdir/spotlight
  wget http://ixa2.si.ehu.es/ixa-pipes/models/wikipedia-db.v1.tar.gz
  tar -xzf wikipedia-db.v1.tar.gz
  rm wikipedia-db.v1.tar.gz
  ◇

```

Fragment defined by 25a, 26a.

Fragment referenced in 15a.

The macro `check/start spotlight` does the following:

1. Check whether spotlight runs on the default spotlighthost.
2. If that is not the case, and the default host is not `localhost`, check whether Spotlight runs on `localhost`.
3. If a running spotlightserver is still not found, start a spotlightserver on `localhost`.

```

<function to check/start spotlight 26b> ≡
function check_start_spotlight {
    naflang=$1
    <get spotlight language parameters 26g>
    spotlighthost=130.37.53.38
    <check listener on host, port (26c $spotlighthost,26d $spotlightport ) 27a>
    if
        [ $spotlightrunning -ne 0 ]
    then
        if
            [ ! "$spotlighthost" == "localhost" ]
        then
            export spotlighthost="localhost"
            <check listener on host, port (26e $spotlighthost,26f $spotlightport ) 27a>
        fi
    fi
    if
        [ $spotlightrunning -ne 0 ]
    then
        <start the Spotlight server on localhost 28b>
    fi
    export spotlighthost
    export spotlightrunning
}
◇

```

Fragment referenced in 55c.

Set the port-number and the language resource for Spotlight, dependent of the language that the user gave as argument.

```

<get spotlight language parameters 26g> ≡
if
    [ "$naflang" == "nl" ]
then
    spotlightport=2060
else
    spotlightport=2020
fi
◇

```

Fragment referenced in 26b.

Uses: `naflang` 21d.

The following macro has a hostname and a port-number as arguments. It checks whether something in the host listens on the port and sets variable `success` accordingly:

```

<check listener on host, port 27a> ≡
exec 6<>/dev/tcp/01/02
spotlightrunning=$?
exec 6<&-
exec 6>&-
◇

```

Fragment referenced in 26b, 28be.

If variable `spotlighthost` does not exist, set it to `localhost`. Test whether a Spotlightserver runs on `spotlighthost`. If that fails and `spotlighthost` did not point to `localhost`, try `localhost`.

If the previous attempts were not succesfull, start the spotlightserver on localhost.

If some spotlightserver has been contacted, set variable `spotlightrunning`. Otherwise exit. At the end variable `spotlighthost` ought to contain the address of the Spotlight-host.

```

< try to obtain a running spotlightserver 27b > ≡
  < test whether spotlighthost runs (27c $spotlighthost ) 28a >
  if
    [ ! $spotlightrunning ]
  then
    if
      [ "$spotlighthost" != "localhost" ]
    then
      export spotlighthost=localhost
      < test whether spotlighthost runs (27d $spotlighthost ) 28a >
    fi
  fi
  if
    [ ! $spotlightrunning ]
  then
    < start the Spotlight server on localhost 28b >
    < test whether spotlighthost runs (27e $spotlighthost ) 28a >
  fi
  if
    [ ! $spotlightrunning ]
  then
    echo "Cannot start spotlight"
    exit 4
  fi
  ◇

```

Fragment never referenced.

Test whether the Spotlightserver runs on a given host. The “spotlight-test” does not really test Spotlight, but it tests whether something is listening on the port and host where we expect Spotlight. I found the test-construction that is used here on [Stackoverflow](#). If the test is positive, set variable `spotlightrunning` to 0. Otherwise, unset that variable.

```

< test whether spotlighthost runs 28a > ≡
  exec 6<>/dev/tcp/01/2060
  if
    [ $? -eq 0 ]
  then
    export spotlightrunning=0
  else
    spotlightrunning=
  fi
  exec 6<&-
  exec 6>&-
  ◇

```

Fragment referenced in 27b.

When trying to start the Spotlight-server on localhost, take care that only one process does this. So we do this:

1. Acquire a lock.
2. Check that in the mean time Spotlight has not been started by another process.
3. Run the Spotlight java program if Spotlight does still not run.

4. release the lock

```

< start the Spotlight server on localhost 28b > ≡
if
  [ "$naflang" == "nl" ]
then
  spotresource="nl"
else
  spotresource="en_2+2"
fi
cd /home/phuijgen/nlp/test/nlpp/env/spotlight
$envbindir/sematree acquire spotlock -1
< check listener on host, port (28c $spotlighthost,28d $spotlightport ) 27a >
if
  [ ! $spotlightrunning -eq 0 ]
then
  java -jar -Xmx8g dbpedia-spotlight-0.7-jar-with-dependencies-
candidates.jar $spotresource http://localhost:$spotlightport/rest &
  < wait until the spotlight server is up 28e >
fi
$envbindir/sematree release spotlock
◇

```

Fragment referenced in 26b, 27b.

```

< wait until the spotlight server is up 28e > ≡
while
  < check listener on host, port (28f $spotlighthost,28g $spotlightport ) 27a >
  [ $spotlightrunning -ne 0 ]
do
  sleep 10
done
◇

```

Fragment referenced in 28b.

Start the Spotlight if it is not already running. First find out what the host is on which we may expect to find a listening Spotlight.

Variable `spotlighthost` contains the address of the host where we expect to find Spotlight. If the expectation does not come true, and the Spotlighthost was not localhost, test whether Spotlight can be found on localhost. If the spotlight-server cannot be found, start it up on localhost.

4.4.9 VUA-pylib

Module VUA-pylib is needed for the opinion-miner. Install it in the Python library

```

< install VUA-pylib 29a > ≡
cd $envdir/python
git clone https://github.com/cltl/VUA_pylib.git
◇

```

Fragment referenced in 15a.

4.4.10 SVMLight

SVMlight supplies a Support Vector Machine. It is used by the opinion-miner. SVMlight can be obtained from [the site](#) where it is documented.

Installation goes like this:

```
<install SVMlight 29b> ≡
  tempdir='mktemp -d -t SVMlight.XXXXXX'
  cd $tempdir
  wget http://download.joachims.org/svm_light/current/svm_light.tar.gz
  tar -xzf svm_light.tar.gz
  make all
  cp svm_classify /home/phuijgen/nlp/test/nlpp/env/bin/
  cp svm_learn /home/phuijgen/nlp/test/nlpp/env/bin/
  cd /home/phuijgen/nlp/test/nlpp
  rm -rf $tempdir
  ◇
```

Fragment referenced in 15a.

Uses: all 59c.

4.4.11 CRFsuite

CRFsuite is an implementation of Conditional Random Fields (CRF). Module **opinion-miner-deluxe** needs it. It can be installed from it's sources, but I did not manage to this. Therefore, currently we use a pre-compiled ball.

```
<install CRFsuite 30a> ≡
  tempdir='mktemp -d -t crfsuite.XXXXXX'
  cd $tempdir
  tar -xzf $snapshotsocket/t_nlpp_resources/crfsuite-0.12-x86_64.tar.gz
  cd crfsuite-0.12
  cp -r bin/crfsuite $envbindir/
  mkdir -p $envdir/include/
  cp -r include/* $envdir/include/
  mkdir -p $envdir/lib/
  cp -r lib/* $envdir/lib/
  cd /home/phuijgen/nlp/test/nlpp
  rm -rf $tempdir
  ◇
```

Fragment referenced in 15a.

4.5 Install modules

4.5.1 Install tokenizer

Module The tokenizer is just a jar that has to be run in Java. Although the jar is directly available from <http://ixa2.si.ehu.es/ixa-pipes/download.html>, we prefer to compile the package in order to make this thing ready for reproducible set-ups.

To install the tokenizer, we proceed as follows:

1. Clone the source from github into a temporary directory.
2. Compile to produce the jar file with the tokenizer.
3. move the jar file into the jar directory.
4. remove the tempdir with the sourcecode.


```

< install the tokenizer 30b > ≡
    tempdir='mktemp -d -t tok.XXXXXX'
    cd $tempdir
    git clone https://github.com/ixa-ehu/ixa-pipe-tok.git
    cd ixa-pipe-tok
    git checkout 56f83ce4b61680346f15e5d4e6de6293764f7383
    mvn clean package
    mv target/ixa-pipe-tok-1.8.0.jar $jarsdir
    cd $piperoot
    rm -rf $tempdir
    ◇

```

Fragment referenced in 16a.

Script The script runs the tokenizerscript.

```

"../bin/tok" 30c≡
    < start of module-script (30d $jarsdir ) 20b >
    < abort when the language is not English or Dutch 21e >
    JARFILE=$jarsdir/ixa-pipe-tok-1.8.0.jar
    java -Xmx1000m -jar $JARFILE tok -l $naflang --inputkaf
    ◇

```

4.5.2 Topic analyser

The English pipeline contains a topic analyser that seems not yet fit for Dutch. Get it from the Newsreader repo and update the config file.

```

< install the topic analyser 31a > ≡
    cd $modulesdir
    tar -xzf $snapshotsocket/t_nlpp_resources/20151220_EHU-topic.v30.tgz
    cd $modulesdir/EHU-topic.v30
    mv conf.prop old.conf.prop
    gawk '{gsub("/home/newsreader/components", subs); print}' subs=$modulesdir old.conf.prop >conf.prop
    ◇

```

Fragment referenced in 16a.

Uses: print 64a.

Script:

```

"../bin/topic" 31b≡
    < start of module-script (31c EHU-topic.v30 ) 20b >
    < abort when the language is not English or Dutch 21e >
    java -Xmx1000m -jar $MODDIR/ixa-pipe-topic-1.0.1.jar -p $MODDIR/conf.prop
    ◇

```

4.5.3 Morphosyntactic parser

Module

```

< install the morphosyntactic parser 31d > ≡
MODNAM=morphsynparser
DIRN=morphosyntactic_parser_nl
GITU=https://github.com/cltl/morphosyntactic_parser_nl.git
GITC=807e938ce4ebb71afd9d7c7f42d9d9ac5f98a184
< install from github 6d >
cd $modulesdir/morphosyntactic_parser_nl
git checkout 807e938ce4ebb71afd9d7c7f42d9d9ac5f98a184
◇

```

Fragment referenced in 16a.

Script The morpho-syntactic module parses the sentences with Alpino. Alpino takes a lot of time to handle long sentences. Therefore the morpho-syntactic module has an option `-t` to set a time-out (in minutes) for sentence parsing.

```

"../bin/mor" 31e≡
< start of module-script (31f morphosyntactic_parser_nl ) 20b >
< get the mor time-out parameter 32a >
< set alpinohome 22a >
cat | python $MODDIR/core/morph_syn_parser.py $timeoutarg
◇

```

Use `getopts` to read the `-t` option.

```

< get the mor time-out parameter 32a > ≡
OPTIND=1
stimeout=
timeoutarg=
while getopts "t:" opt; do
    case "$opt" in
        t) stimeout=$OPTARG
           ;;
        esac
done
shift $((OPTIND-1))
if
[ $stimeout ]
then
    timeoutarg="-t $stimeout"
fi
◇

```

Fragment referenced in 31e.

4.5.4 Pos tagger

In the Dutch pipeline the morpho-syntactic parser fulfills the role of Pos tagger. In the English pipeline we use the pos-tagger from EHU.

Module

```

< install the pos tagger 32b > ≡
  cd $modulesdir
  tar -xzf $snapshotsocket/t_nlpp_resources/20151220_EHU-pos.v30.tgz
  cd $modulesdir/EHU-topic.v30
  ◇

```

Fragment referenced in 16a.

Script

```

"../bin/pos" 32c ≡
  < start of module-script (32d EHU-pos.v30 ) 20b >
  java -Xmx1000m -jar ${MODDIR}/ixa-pipe-pos-1.4.3.jar tag -m ${MODDIR}/en-maxent-
  100-c5-baseline-dict-penn.bin
  ◇

```

4.5.5 Constituent parser

Module

```

< install the constituents parser 32e > ≡
  cd $modulesdir
  tar -xzf $snapshotsocket/t_nlpp_resources/20151220_EHU-parse.v30.tgz
  cd $modulesdir/conspardir
  chmod 775 *.jar
  chmod 775 *.bin
  ◇

```

Fragment referenced in 16j.

Script

```

"../bin/constpars" 33a ≡
  < start of module-script (33b EHU-parse.v30 ) 20b >
  java -Xmx1000m -jar ${MODDIR}/ixa-pipe-parse-1.1.1.jar parse -g sem -
  m ${MODDIR}/en-parser-chunking.bin
  ◇

```

4.5.6 NED-reranker

Module

```

< install the NED-reranker module 33c > ≡
  cd $modulesdir
  tar -xzf $snapshotsocket/t_nlpp_resources/20151220_VUA-popen-nedreranker.v30.tgz
  ◇

```

Fragment referenced in 16j.

Script

```
"../bin/nedrer" 33d≡
  ⟨ start of module-script (33e VUA-popen-nedreranker.v30 ) 20b ⟩
  cd $MODDIR
  python $MODDIR/domain_model.py
  ◇
```

4.5.7 Wikify module

Module

```
⟨ install the wikify module 33f ⟩ ≡
  cd $modulesdir
  tar -xzf $snapshotsocket/t_nlpp_resources/20151220_EHU-wikify.v30.tgz
  ◇
```

Fragment referenced in 16j.

Script The Wikify module needs DBpedia to generate “markables”.

```
"../bin/wikify" 33g≡
  ⟨ start of module-script (33h EHU-wikify.v30 ) 20b ⟩
  if
    [ "$naflang" == "nl" ]
  then
    spotlightport=2060
  else
    spotlightport=2020
  fi
  [ $spotlightrunning ] || source /home/phuijgen/nlp/test/nlpp/bin/start-spotlight

  cd $MODDIR
  java -Xmx1000m -jar ${MODDIR}/ixa-pipe-wikify-1.2.1.jar -s http://$spotlighthost -
  p $spotlightport
  ◇
```

4.5.8 UKB

UKB needs boost libraries and Perl version 5. For now, we consider them installed.

Module

```
⟨ install the UKB module 34a ⟩ ≡
  ◇
```

Fragment never referenced.

Script Put the path to the boost libraries in the LD_LIBRARY_PATH variable and then run UKB.

```

"../bin/ukb" 34b≡
  < start of module-script (34c EHU-ukb.v30 ) 20b >
  cd $MODDIR
  export LD_LIBRARY_PATH=$LD_LIBRARY_PATH:$envdir/lib:$envdir/boost_1_54_0/stage/lib
  ${MODDIR}/naf_ukb/naf_ukb.pl -x ${MODDIR}/ukb/bin/ukb_wsd -K ${MODDIR}/wn30-
  ili_lkb/wn30g.bin64 -D ${MODDIR}/wn30-ili_lkb/wn30.lex - -- --dict_weight --
  dgraph_dfs --dgraph_rank ppr
  ◇

```

4.5.9 IMS-WSD

Module The package itself supplies an installation script that seems usable. However, today I am in a hurry and just install the module as it comes from the EHU repository.

Although the Hadoop implementation runs this module with Java 1.7, I could only run `ims+wsd` Java 1.6. Using Java 1.7 causes run-time errors “Platform not recognised” and the resulting NAF’s do not contain WordNet references. So, we had to install Java 1.6.

The scripts contain explicit paths that must be corrected:

`ims/testPlain`: Explicit path to Java binary.

`path_to_ims.py`: Set variable `PATH_TO_IMS`.

```

< install the ims-wsd module 34d > ≡
  cd $modulesdir
  tar -xzf $snapshotsocket/t_nlpp_resources/20151220_VUA-ims-wsd.v30.tgz
  cd VUA-ims-wsd.v30
  thisDir='pwd'
  echo PATH_TO_IMS = ""$thisDir/ims"" > path_to_ims.py
  cd ims
  cp testPlain.bash old.testPlain.bash
  sedcommand='s|/usr/lib/jvm/java-1.6.0-openjdk-1.6.0.0.x86_64/jre/bin/java|java|g'
  cat old.testPlain.bash | sed $sedcommand > testPlain.bash
  ◇

```

Fragment referenced in [17a](#).

Script

```

"../bin/ewsd" 35a≡
  < start of module-script (35b VUA-ims-wsd.v30 ) 20b >
  < set up Java 1.6 10a >
  #Setting the output to be ili-wn30 synsets instead of sensekeys
  $MODDIR/call_ims.py -ili30
  ◇

```

4.5.10 SRL server

The EHU SRL-module, that we use for English documents, has been set up as a server/client system. Hence, we have to start the server before we can process something.

We don’t know in advance whether we run the pipeline for a single text or from a whole bunch of text and hence we do not know whether it is advisable that the server keeps running, occupying precious memory. Therefore, currently we just start and stop the server every time that we use it.

Module

```

< install the srl-server module 35c > ≡
    cd $modulesdir
    tar -xzf $snapshotsocket/t_nlpp_resources/20151220_EHU-srl-server.tgz
    cd EHU-srl-server
    mkdir -p /home/phuijgen/nlp/test/nlpp/env/etc/pid
    ◇

```

Fragment referenced in 17a.

Scripts Generate three scripts: `start_eSRL`, `stop_esrl` and `eSRL`, resp. to start the SRL server, to stop it and to process a NAF file.

```

"../bin/start_eSRL" 35d≡
    < start of module-script (35e EHU-srl-server ) 20b >
    < start EHU SRL server if it isn't running 36a >
    ◇

```

```

"../bin/stop_eSRL" 35f≡
    < start of module-script (35g EHU-srl-server ) 20b >
    < stop EHU SRL server 36b >
    ◇

```

```

"../bin/eSRL" 35h≡
    < start of module-script (35i EHU-srl-server ) 20b >
    /home/phuijgen/nlp/test/nlpp/bin/start_eSRL
    java -Xmx1000m -cp $MODDIR/IXA-EHU-srl-3.0.jar ixa.srl.SRLClient en
    ◇

```

```

< start EHU SRL server if it isn't running 36a > ≡
    pidFile=/home/phuijgen/nlp/test/nlpp/env/etc/pid/SRLServer.pid
    portInfo=$(nmap -p 5005 localhost | grep open)
    if [ -z "$portInfo" ]; then
        >&2 echo "Starting srl-server as it is not running"
        java -Xms2500m -cp $MODDIR/IXA-EHU-srl-
        3.0.jar ixa.srl.SRLServer en &> /dev/null &
        pid=$!
        echo $pid > $pidFile
        sleep 60
        >&2 echo "Server running: ${pid}"
    else
        >&2 echo "Server already running.."
    fi
    ◇

```

Fragment referenced in 35d.

```

< stop EHU SRL server 36b > ≡
    pidFile=/home/phuijgen/nlp/test/nlpp/env/etc/pid/SRLServer.pid
    if
        [ -e "$pidFile" ]
    then
        kill 'echo $pidFile'
        rm $pidFile
    fi
    ◇

```

Fragment referenced in 35f.

4.5.11 SRL Dutch nominals

Module

```

< install the srl-dutch-nominals module 36c > ≡
    MODNAM=srl-dutch-nominals
    DIRN=vua-srl-dutch-nominal-events
    GITU=https://github.com/newsreader/vua-srl-dutch-nominal-events
    GITC=6115b3168978acf809916cd2da512295d109d8fb
    < install from github 6d >
    cd $modulesdir/vua-srl-dutch-nominal-events
    chmod 775 vua-srl-dutch-additional-roles.py
    ◇

```

Fragment referenced in 17a.

Script

```

"../bin/srl-dutch-nominals" 36d ≡
    < start of module-script (36e vua-srl-dutch-nominal-events ) 20b >
    cd $MODDIR
    cat | $MODDIR/vua-srl-dutch-additional-roles.py
    ◇

```

4.5.12 FBK-time module

Module

```

< install the FBK-time module 37 > ≡
    cd $modulesdir
    tar -xzf $snapshotsocket/t_nlpp_resources/20151220_FBK-time.v30.tgz
    ◇

```

Fragment referenced in 17j.

Script The script is rather complicated. I just copied it from the original makers, with one exception: Originally at the end of the script there was a pipe consisting of two Java programs. However, that didn't seem to work in one of the computers that we use, therefore we have split the pipe using `mytemp` as temporary storage.

```

"../bin/FBK-time" 38a≡
  < start of module-script (38b FBK-time.v30 ) 20b >
  BEGINTIME='date +%Y-%m-%dT%H:%M:%S%Z'
  YAMCHA=$MODDIR/tools
  timdir='mktemp -d -t time.XXXXXX'
  FILETXP=$timdir/TimePro.txp
  CHUNKIN=$timdir/TimePro.naf
  FILEOUT=$timdir/TimeProOUT.txp
  TIMEPRONORMIN=$timdir/TimeProNormIN.txp
  JAVAMAXHEAP=2g
  mytemp=$timdir/mytemp
  cd $MODDIR
  cat > $CHUNKIN

  JAVAClasspath="lib/jdom-2.0.5.jar:lib/kaflib-naf-1.1.9.jar:lib/NAFtoTXP_v11.jar"
  JAVAMODULE=eu.fbk.newsreader.naf.NAFtoTXP_v11
  cat $CHUNKIN | \
    java -Xmx$JAVAMAXHEAP -cp $JAVAClasspath $JAVAMODULE $FILETXP chunk+entity timex

  #echo "Saving... $FILETXP"
  tail -n +4 $FILETXP | awk -f resources/english-rules > $FILEOUT
  head -n +4 $FILETXP > $TIMEPRONORMIN

  cat $FILEOUT | \
    $YAMCHA/yamcha-0.33/usr/local/bin/yamcha \
      -m models/tempeval3_silver-data.model \
    >> $TIMEPRONORMIN

  JAVAClasspath="lib/scala-library.jar:lib/timenorm-0.9.1-SNAPSHOT.jar"
  JAVAClasspath=$JAVAClasspath:"lib/threetenbp-0.8.1.jar:lib/TimeProNorm_v2.5.jar"
  JAVAMODULE=eu.fbk.timepro.TimeProNormApply
  cat $TIMEPRONORMIN | \
    java -Xmx$JAVAMAXHEAP -cp $JAVAClasspath $JAVAMODULE $FILETXP

  rm $FILEOUT
  rm $TIMEPRONORMIN

  JAVAClasspath="lib/jdom-2.0.5.jar:lib/kaflib-naf-1.1.9.jar:lib/NAFtoTXP_v11.jar"
  JAVAMODULE=eu.fbk.newsreader.naf.NAFtoTXP_v11
  cat $CHUNKIN | java -Xmx$JAVAMAXHEAP -
  cp $JAVAClasspath $JAVAMODULE $FILEOUT chunk+morpho+timex+event eval

  JAVACP1="lib/TXPtoNAF_v5.jar:lib/jdom-2.0.5.jar:lib/kaflib-naf-1.1.9.jar"
  JAVAMOD1=eu.fbk.newsreader.naf.TXPtoNAF_v4
  JAVACP2="lib/kaflib-naf-1.1.9.jar:lib/jdom-2.0.5.jar:lib/TimeProEmptyTimex_v2.jar"
  JAVAMOD2=eu.fbk.timepro.TimeProEmptyTimex
  java -Xmx$JAVAMAXHEAP -Dfile.encoding=UTF8 -
  cp $JAVACP1 $JAVAMOD1 $CHUNKIN $FILETXP "$BEGINTIME" TIMEX3 > $mytemp
  cat $mytemp | java -Dfile.encoding=UTF8 -cp $JAVACP2 $JAVAMOD2 $FILEOUT

  rm $FILETXP
  rm $CHUNKIN
  rm -rf $timdir
  ◇

```


4.5.13 FBK-temprel module

Module

```

⟨ install the FBK-temprel module 39a ⟩ ≡
    cd $modulesdir
    tar -xzf $snapshotsocket/t_nlpp_resources/20151220_FBK-temprel.v30.tgz
    ⟨ repair FBK-*rel's run.sh.hadoop (39b FBK-temprel.v30 ) 39c ⟩
    ◇

```

Fragment referenced in 17j.

Script run.sh.hadoop seems to be obsolete in the original tarball:

1. The class-path argument in one of the Java statement refers to an obsolete jar (kaflib-naf-1.1.8 instead of kaflib-naf-1.1.9)
2. Another class-path argument refers to PredicateTimeAnchor_tlink.jar instead of PredicateTimeAnchor.jar
3. A “sh” statement is used. The problem is, that in Ubuntu /bin/sh points to bin/dash and the script (temprel-pipeline-per-file-NWR.sh) does not seem to be compatible with dash.

Therefore, we need to repair the script. We will need to repair the script in the FBK-causalrel module in a similar way, and therefore provide the module-directory as argument.

```

⟨ repair FBK-*rel's run.sh.hadoop 39c ⟩ ≡
    cd $modulesdir/@1
    mv run.sh.hadoop old.run.sh.hadoop
    cat old.run.sh.hadoop | \
        sed s/kaflib-naf-1.1.8/kaflib-naf-1.1.9/g | \
        sed s/TimeAnchor_tlink.jar/TimeAnchor.jar/g | \
        sed "s/sh temprel/bash temprel/g" | \
        sed "s/java /java -Xmx2g /g" \
    >run.sh.hadoop
    chmod 775 run.sh.hadoop
    ◇

```

Fragment referenced in 39a, 40a.

Script The original run script seems to not only read the input naf from standard in, but also to obtain the input naf as a file that an argument points to. This constructions makes the pipeline complicated, therefore, we generate the naf file within the script.

The original script generates temporary files in the `temp` directory of the host-computer, and prefixes the names of the temporary files with a random number to prevent confusion between tempfiles of different instances of this module. We generate a temp-directory per instance.

```

"../bin/FBK-temprel" 39d≡
    ⟨ start of module-script (39e FBK-temprel.v30 ) 20b ⟩
    cd $MODDIR
    scratchDir='mktemp -d -t temprel.XXXXXX'
    cat >$scratchDir/in.naf
    ./run.sh.hadoop $MODDIR $scratchDir $scratchDir/in.naf
    rm -rf $scratchDir
    ◇

```

4.5.14 FBK-causalrel module

Module

```

< install the FBK-causalrel module 40a > ≡
  cd $modulesdir
  tar -xzf $snapshotsocket/t_nlpp_resources/20151220_FBK-causalrel.v30.tgz
  < repair FBK-*rel's run.sh.hadoop (40b FBK-causalrel.v30 ) 39c >
  ◇

```

Fragment referenced in 17j.

Like in FBK-temprel, script run.sh.hadoop seems not to work out of the box:

1. The class-path argument in one of the Java statement refers to an obsolete jar (kaflib-naf-1.1.8 instead of kaflib-naf-1.1.9)
2. A “sh” statement is used. The problem is, that in Ubuntu /bin/sh points to bin/dash and the script (temprel-pipeline-per-file-NWR.sh) does not seem to be compatible with dash.

Therefore, we need to repair that script like we did in FBK-temprel.

```

< repair causalrel's run.sh.hadoop 40c > ≡
  cd $modulesdir/FBK-causalrel.v30
  mv run.sh.hadoop old.run.sh.hadoop
  cat old.run.sh.hadoop | \
    sed s/kaflib-naf-1.1.8/kaflib-naf-1.1.9/g | \
    sed s/TimeAnchor_tlink.jar/TimeAnchor.jar/g | \
    sed s/sh temprel/bash temprel/g | \
  >run.sh.hadoop
  chmod 775 run.sh.hadoop
  ◇

```

Fragment never referenced.

Script

```

"../bin/FBK-causalrel" 40d ≡
  < start of module-script (40e FBK-causalrel.v30 ) 20b >
  cd $MODDIR
  scratchDir='mktemp -d -t causalrel.XXXXXX'
  cat >$scratchDir/in.naf
  ./run.sh.hadoop $MODDIR $scratchDir $scratchDir/in.naf
  rm -rf $scratchDir
  ◇

```

4.5.15 Factuality module

Module

```

< install the factuality module 40f > ≡
  cd $modulesdir
  tar -xzf $snapshotsocket/t_nlpp_resources/20151220_VUA-factuality.v30.tgz
  ◇

```

Fragment referenced in 17j.

Script

```
"../bin/factuality" 41a≡
  ⟨ start of module-script (41b VUA-factuality.v30 ) 20b ⟩
  cd $MODDIR
  #local settings to prevent perl from complaining
  export LANGUAGE=en_US.UTF-8
  export LANG=en_US.UTF-8
  export LC_ALL=en_US.UTF-8

  rootDir=${MODDIR}
  tmpDir=$(mktemp -d -t factuality.XXXXXX)

  export PATH=$PATH:${rootDir}:.
  # ex-
  port LD_LIBRARY_PATH=$LD_LIBRARY_PATH:${rootDir}/../opt/lib:${rootDir}/../opt/boost_1_54_0/stage/lib
  export LD_LIBRARY_PATH=$LD_LIBRARY_PATH:/home/phuijgen/nlp/test/nlpp/env/lib/

  #mkdir -p ${scratchDir}/test

  python ${rootDir}/vua_factuality_naf_wrapper.py -
  t /home/phuijgen/nlp/test/nlpp/env/bin/timbl -p ${rootDir} ${tmpDir}/

  ◇
```

4.5.16 Nominal coreference-base

Get this thing from Github (<https://github.com/opener-project/coreference-base/>) and apply the instruction of <https://github.com/opener-project/coreference-base/blob/master/core/README.md>. We implement it, but it does not work yet, because it is too picky on the structure of the NAF format.

Module

```
⟨ install coreference-base 41c ⟩ ≡
  MODNAM=coreference-base
  DIRN=coreference-base
  GITU=https://github.com/opener-project/coreference-base.git
  GITC=bfa5aec0fa498e57fe14dd4d2c51365dd09a0757
  ⟨ install from github 6d ⟩
  pip install --upgrade hg+https://bitbucket.org/Josu/pykaf#egg=pykaf
  pip install --upgrade networkx

  ◇
```

Fragment referenced in 18a.

Uses: hg 19f.

Script

```
"../bin/coreference-base" 41d≡
  ⟨ start of module-script (41e coreference-base ) 20b ⟩
  cd $MODDIR/core
  cat | python -m corefgraph.process.file --language nl --singleton --sieves NO

  ◇
```

4.5.17 Named entity recognition (NERC)

Module The Nerc program can be installed from Github (<https://github.com/ixa-ehu/ixa-pipe-nerc>). However, the model that is needed is not publicly available. Therefore, models have been put in the snapshot-tarball.

```
< install the NERC module 42a > ≡
  < compile the nerc jar 42b >
  < get the nerc models 43a >
```

◇

Fragment referenced in 16j.

The nerc module is a Java program that is contained in a jar. Put the source from Github in a temporary directory, compile the jar with java and move the jar to the jars directory.

```
< compile the nerc jar 42b > ≡
TEMPDIR='mktemp -d -t nerc.XXXXXX'
cd $TEMPDIR
git clone https://github.com/ixa-ehu/ixa-pipe-nerc
cd ixa-pipe-nerc/
git checkout ca02c931bc0b200ccdb8b5795a7552e4cc0d4802
mvn clean package
mv target/ixa-pipe-nerc-1.5.4.jar $jarsdir/
cd $nuwebdir
rm -rf $TEMPDIR
```

◇

Fragment referenced in 42a.

The current version of the pipeline uses the following models, that have been made available by Rodrigo Agerri on december 15, 2015.

The tarball `dutch-nerc-models.tar.gz` contains the models `nl-clusters-conll102.bin` and `nl-clusters-sonar.bin`. Both models have been placed in subdirectory `/m4_nerc_nl_dir/nerc_models/nl` of the snapshot.

The model for English can be found in the newsreader-repository.

Choose a model dependent of the language.

```
< select language-dependent features 42c > ≡
if
  [ "$naflang" == "nl" ]
then
  export nercmodel=nl/nl-clusters-conll102.bin
else
  export nercmodel=en/en-newsreader-clusters-3-class-muc7-conll103-ontonotes-4.0.bin
fi
```

◇

Fragment referenced in 55c.

Uses: `naflang` 21d.

The tarball `20160105_nerc_models.tgz` contains in subdirectories `nl` and `en` a dutch resp. an english nerc-model. They have been randomly selected from a number of models that are available in <http://ixa2.si.ehu.es/ixa-pipes/models/nerc-models-1.5.4.tgz>.

```

< get the nerc models 43a > ≡
    cd $modulesdir
    tar -xzf /home/phuijgen/nlp/test/t_nlpp_resources/20160105_nerc_models.tgz
    ◇

```

Fragment referenced in 42a.

Script Make a script that uses the conll02 model and a script that uses the Sonar model

```

"../bin/nerc_conll02" 43b≡
    < start of module-script (43c m4_nerc_nl_dir ) 20b >
    JAR=$jarsdir/ixa-pipe-nerc-1.5.4.jar
    MODEL=nl-clusters-conll02.bin
    cat | java -Xmx1000m -jar $JAR tag -m $MODDIR/nl/$MODEL
    ◇

"../bin/nerc" 43d≡
    < start of module-script (43e nerc-models ) 20b >
    JAR=$jarsdir/ixa-pipe-nerc-1.5.4.jar
    if
        [ "$naflang" == "nl" ]
    then
        nercmodel=$modulesdir/nerc_models/nl/nl-6-class-clusters-sonar.bin
    else
        nercmodel=$modulesdir/nerc_models/en/en-best-clusters-conll03.bin
    fi
    java -jar $JAR tag -m $nercmodel
    ◇

```

4.5.18 Wordsense-disambiguation

Install WSD from its Github source (https://github.com/cltl/svm_wsd.git). According to the `readme` of that module, the next thing to do is, to execute `install-script install.sh` or `install_naf.sh`. The latter script installs a “Support-Vector-Machine” (SVM) module, “Dutch-SemCor” (DSC) models and `KafNafParserPy`.

Module

```

< install the WSD module 43f > ≡
    MODNAM=wsd
    DIRN=svm_wsd
    GITU=https://github.com/cltl/svm_wsd.git
    GITC=030043903b42f77cd20a9b2443de137e2efe8513
    < install from github 6d >
    cd $modulesdir/svm_wsd
    < install svm lib 44a >
    < download svm models 44b >

    ◇

```

Fragment referenced in 18a.

This part has been copied from `install_naf.sh` in the WSD module.

```
< install svm lib 44a > ≡
    mkdir lib
    cd lib
    wget --no-check-
    certificate https://github.com/cjlin1/libsvm/archive/master.zip 2>/dev/null
    zip_name='ls -1 | head -1'
    unzip $zip_name > /dev/null
    rm $zip_name
    folder_name='ls -1 | head -1'
    mv $folder_name libsvm
    cd libsvm/python
    make > /dev/null 2> /dev/null
    echo LIBSVM installed correctly lib/libsvm
    ◇
```

Fragment referenced in 43f.

This part has also been copied from `install_naf.sh` in the WSD module.

```
< download svm models 44b > ≡
    cd $modulesdir/svm_wsd
    #tar -xzf $pipesocket/m4_wsd_snapball
    wget --user=cltl --
    password='.cltl.' kyoto.let.vu.nl/~izquierdo/models_wsd_svm_dsc.tgz 2> /dev/null
    echo 'Unzipping models...'
    tar xzf models_wsd_svm_dsc.tgz
    rm models_wsd_svm_dsc.tgz
    echo 'Models installed in folder models'
    ◇
```

Fragment referenced in 43f.

Script

```
"../bin/wsd" 44c≡
    < start of module-script (44d svm_wsd ) 20b >
    WSDSCRIPT=dsc_wsd_tagger.py
    cat | python $MODDIR/$WSDSCRIPT --naf -ref odwnSY
    ◇
```

4.5.19 Lexical-unit converter

Module There is not an official repository for this module yet, so copy the module from the tarball.

```
< install the lu2synset converter 44e > ≡
    cd $modulesdir
    tar -xzf $snapshotsocket/t_nlpp_resources/lu2synset.tgz
    ◇
```

Fragment referenced in 18h.

Script

```
"../bin/lu2synset" 45a≡
  ⟨ start of module-script (45b lexicalunitconverter ) 20b⟩
  JAVA_LIBDIR=$MODDIR/lib
  RESOURCESDIR=$MODDIR/resources
  JARFILE=WordnetTools-1.0-jar-with-dependencies.jar
  java -Xmx812m -
  cp $JAVA_LIBDIR/$JARFILE vu.wntools.util.NafLexicalUnitToSynsetReferences \
    --wn-lmf "$RESOURCESDIR/cornetto2.1.lmf.xml" --format naf
  ◇
```

4.5.20 NED

The NED module is rather picky about the structure of the NAF file. In any case, it does not accept a file that has been produced by the *ontotagger*. Hence, in a pipeline NED should be executed before the *ontotagger*.

The NED module wants to consult the Dbpedia Spotlight server, so that one has to be installed somewhere. For this moment, let us suppose that it has been installed on localhost.

Module

```
⟨ install the NED module 45c ⟩ ≡
  ⟨ put spotlight jar in the Maven repository 46a ⟩
  MODNAM=ned
  DIRN=ixa-pipe-ned
  GITU=https://github.com/ixa-ehu/ixa-pipe-ned.git
  GITC=d35d4df5cb71940bf642bb1a83e2b5b7584010df
  ⟨ install from github 6d ⟩
  cd $modulesdir/ixa-pipe-ned
  mvn -Dmaven.compiler.target=1.7 -Dmaven.compiler.source=1.7 clean package
  mv target/ixa-pipe-ned-1.1.1.jar $jarsdir/
  ◇
```

Fragment referenced in 16j.

NED needs to have *dbpedia-spotlight-0.7.jar* in the local Maven repository. That is a different jar than the jar that we use to start Spotlight.

```

<put spotlight jar in the Maven repository 46a> ≡
    echo Put Spotlight jar in the Maven repository.
    tempdir='mktemp -d -t simplespot.XXXXXX'
    cd $tempdir
    wget http://spotlight.sztaki.hu/downloads/archive/2014/dbpedia-spotlight-0.7.jar
    wget http://spotlight.sztaki.hu/downloads/archive/2014/nl.tar.gz
    tar -xzf nl.tar.gz
    MVN_SPOTLIGHT_OPTIONS="-Dfile=dbpedia-spotlight-0.7.jar"
    MVN_SPOTLIGHT_OPTIONS="$MVN_SPOTLIGHT_OPTIONS -DgroupId=ixa"
    MVN_SPOTLIGHT_OPTIONS="$MVN_SPOTLIGHT_OPTIONS -DartifactId=dbpedia-spotlight"
    MVN_SPOTLIGHT_OPTIONS="$MVN_SPOTLIGHT_OPTIONS -Dversion=0.7"
    MVN_SPOTLIGHT_OPTIONS="$MVN_SPOTLIGHT_OPTIONS -Dpackaging=jar"
    MVN_SPOTLIGHT_OPTIONS="$MVN_SPOTLIGHT_OPTIONS -DgeneratePom=true"
    mvn install:install-file $MVN_SPOTLIGHT_OPTIONS

    cd $PROJROOT
    rm -rf $tempdir
    ◇

```

Fragment referenced in 45c.

Script NED needs to contact a Spotlight-server.

```

"../bin/ned" 46b≡
    <start of module-script (46c ) 20b>
    #!/bin/bash
    source /home/phuijgen/nlp/test/nlpp/env/bin/progenv
    ROOT=$piperoot
    JARDIR=$jarsdir
    if
        [ "$naflang" == "nl" ]
    then
        spotlightport=2060
    else
        spotlightport=2020
    fi
    [ $spotlightrunning ] || source /home/phuijgen/nlp/test/nlpp/bin/start-spotlight
    cat | java -Xmx1000m -jar $jarsdir/ixa-pipe-ned-1.1.1.jar -
    H http://$spotlighthost -p $spotlightport -e candidates -
    i $envdir/spotlight/wikipedia-db -n nlEn
    ◇

```

4.5.21 Ontotagger

We do not yet have a source-repository of the Ontotagger module. Therefore, install from a snapshot (20160126_vua-ontotagger-v1.0.tgz).

Module

```

<install the onto module 46d> ≡
    cd $modulesdir
    tar -xzf $snapshotsocket/t_nlpp_resources/20160126_vua-ontotagger-v1.0.tgz
    chmod -R o+r $modulesdir/vua-ontotagger-v1.0
    ◇

```

Fragment referenced in 18a.

Script

```

"../bin/onto" 47a≡
  ( start of module-script (47b vua-ontotagger-v1.0 ) 20b )
  JARDIR=$MODDIR/lib
  RESOURCEDIR=$MODDIR/resources
  PREDICATEMATRIX="$RESOURCEDIR/PredicateMatrix.v1.3.txt.role.odwn"
  GRAMMATICALWORDS="$RESOURCEDIR/grammaticals/Grammatical-words.nl"
  TMPFIL='mktemp -t stap6.XXXXXX'
  cat >$TMPFIL

  CLASSPATH=$JARDIR/ontotagger-1.0-jar-with-dependencies.jar
  JAVASCRIPT=eu.kyotoproject.main.KafPredicateMatrixTagger
  MAPPINGS="fn;mcr;ili;eso"
  JAVA_ARGS="--mappings $MAPPINGS"
  JAVA_ARGS="$JAVA_ARGS --key odwn-eq"
  JAVA_ARGS="$JAVA_ARGS --version 1.2"
  JAVA_ARGS="$JAVA_ARGS --predicate-matrix $PREDICATEMATRIX"
  JAVA_ARGS="$JAVA_ARGS --grammatical-words $GRAMMATICALWORDS"
  JAVA_ARGS="$JAVA_ARGS --naf-file $TMPFIL"
  java -Xmx1812m -cp $CLASSPATH $JAVASCRIPT $JAVA_ARGS
  rm -rf $TMPFIL

  ◇

```

4.5.22 Framenet SRL

The framenet SRL is part of the package that contains the ontotagger. We only need a different script.

Script The script contains a hack, because the framesrl script produces spurious lines containing “frameMap.size()=...”. A GAWK script removes these lines.

```

"../bin/framesrl" 48a≡
  < start of module-script (48b vua-ontotagger-v1.0 ) 20b>
  ONTODIR=$modulesdir/vua-ontotagger-v1.0
  JARDIR=$MODDIR/lib
  RESOURCEDIR=$MODDIR/resources
  PREDICATEMATRIX="$RESOURCEDIR/PredicateMatrix_nl_lu_withES0.v0.2.role.txt"
  GRAMMATICALWORDS="$RESOURCEDIR/grammaticals/Grammatical-words.nl"
  TMPFIL='mktemp -t framesrl.XXXXXX'
  cat >$TMPFIL

  CLASSPATH=$JARDIR/ontotagger-1.0-jar-with-dependencies.jar
  JAVASCRIPT=eu.kyotoproject.main.SrlFrameNetTagger

  JAVA_ARGS="--naf-file $TMPFIL"
  JAVA_ARGS="$JAVA_ARGS --format naf"
  JAVA_ARGS="$JAVA_ARGS --frame-ns fn:"
  JAVA_ARGS="$JAVA_ARGS --role-ns fn-role;;pb-role;;fn-pb-role;;eso-role:"
  JAVA_ARGS="$JAVA_ARGS --ili-ns mcr:ili"
  JAVA_ARGS="$JAVA_ARGS --sense-conf 0.25"
  JAVA_ARGS="$JAVA_ARGS --frame-conf 70"

  java -Xmx1812m -
  cp $CLASSPATH $JAVASCRIPT $JAVA_ARGS | gawk '/^frameMap.size()/ {next}; {print}'
  rm -rf $TMPFIL

◇

```

4.5.23 Heideltime

Module The code for Heideltime can be found in [Github](#). However, we use a compiled Heideltime Jar, compiled by Antske Fokkens, because some bugs have been repaired in that version.

Use Heideltime via a wrapper, `ixa-pipe-time`, obtained from [Github](#).

Heideltime uses treetagger. It expects to find the location of treetagger in a variable `TreetaggerHome` in config-file `config.props`.

```

< install the heideltime module 48c> ≡
  moduledir=/home/phuijgen/nlp/test/nlpp/modules/ixa-pipe-time
  < clone the heideltime wrapper 48d>
  < put Antske's material in the heideltime wrapper 49a>
  < compile the heideltime wrapper 49b>

◇

```

Fragment referenced in [18h](#).

```

< clone the heideltime wrapper 48d> ≡
  MODNAM=heideltime
  DIRN=ixa-pipe-time
  GITU=https://github.com/ixa-ehu/ixa-pipe-time.git
  GITC=da4604a7b33975e977017440cbc10f7d59917ddf
  < install from github (48e ixa-pipe-time ) 6d>
  mkdir $moduledir/lib

◇

```

Fragment referenced in [48c](#).

In the wrapper we need the following extra material:

- A debugged version of the Heidelberg jar.
- A configuration file `config.props`, although it does not seem to be actually used.
- Another configuration file: `alpino-to-treetagger.csv`

The extra material has been provided by Antske Fokkens.

```
< put Antske's material in the heideltime wrapper 49a > ≡
cd $modulesdir/$DIRN
tar -xzf $snapshotsocket/t_nlpp_resources/20151123_antske_heideltime_stuff.tgz
mv antske_heideltime_stuff/de.unihd.dbs.heideltime.standalone.jar lib/
mv antske_heideltime_stuff/config.props .
mv antske_heideltime_stuff/alpino-to-treetagger.csv .
rm -rf antske_heideltime_stuff
◇
```

Fragment referenced in 48c.

Compile the Heideltime wrapper according to the [instruction](#) on Github.

```
< compile the heideltime wrapper 49b > ≡
< get jvntextpro-2.0.jar 49c >
< activate the install-to-project-repo utility 49d >
cd /home/phuijgen/nlp/test/nlpp/modules/$DIRN
mvn clean install
◇
```

Fragment referenced in 48c.

```
< get jvntextpro-2.0.jar 49c > ≡
cd /home/phuijgen/nlp/test/nlpp/modules/$DIRN/lib
wget http://ixa2.si.ehu.es/%7Ejibalar/jvntextpro-2.0.jar
◇
```

Fragment referenced in 49b.

Script `install-to-project-repo.py` generates a library in subdirectory `repo` and copies the jars that it finds in the `lib` subdirectory in this repo in such a way that Maven finds it there. Somewhere in the `install-to-project.py ... mvn` process the jars are copied in your local repository (`~/.m2`) too. As a result, only a Maven Guru understands precisely where Maven obtains its jar from and the best thing to do is to empty the `repo` subdirectory and the local repository before (re-) applying `install-to-project-repo.py`.

```
< activate the install-to-project-repo utility 49d > ≡
< remove outdated heideltime jars 50a >
cd /home/phuijgen/nlp/test/nlpp/modules/$DIRN/
git clone git@github.com:carchrae/install-to-project-repo.git
mv install-to-project-repo/install-to-project-repo.py .
rm -rf install-to-project-repo
python ./install-to-project-repo.py
◇
```

Fragment referenced in 49b.

```

< remove outdated heideltime jars 50a > ≡
  rm -rf /home/phuijgen/nlp/test/nlpp/modules/$DIRN/repo
  mkdir -p /home/phuijgen/nlp/test/nlpp/modules/$DIRN/repo/local
  rm -rf $HOME/.m2/repository/local/de.unihd.dbs.heideltime.standalone
  rm -rf $HOME/.m2/repository/local/jvntextpro-2.0
  ◇

```

Fragment referenced in 49d.

Script

```

"../bin/heideltime" 50b ≡
  < start of module-script (50c ixa-pipe-time ) 20b >
  MODDIR=$modulesdir/ixa-pipe-time
  cd $MODDIR
  iconv -t utf-8//IGNORE | java -Xmx1000m -jar target/ixa.pipe.time.jar -m alpino-to-
  treetagger.csv -c config.props
  ◇

```

4.5.24 Semantic Role labelling

Module

```

< install the srl module 50d > ≡
  MODNAM=srl
  DIRN=vua-srl-nl
  GITU=https://github.com/newsreader/vua-srl-nl.git
  GITH=675d22d361289ede23df11dcdb17195f008c54bf
  < install from github 6d >
  ◇

```

Fragment referenced in 18h.

Script First:

1. set the correct environment. The module needs python and timble.
2. create a tempdir and in that dir a file to store the input and a (scv) file with the feature-vector.

```

"../bin/srl" 50e ≡
  < start of module-script (50f vua-srl-nl ) 20b >
  MODDIR=$modulesdir/vua-srl-nl
  TEMPDIR='mktemp -d -t SRLTMP.XXXXXX'
  cd $MODDIR
  INPUTFILE=$TEMPDIR/inputfile
  FEATUREVECTOR=$TEMPDIR/csvfile
  TIMBLOUTPUTFILE=$TEMPDIR/timblpredictions
  ◇

```

File defined by 50e, 51abcd.

Create a feature-vector.

```
"../bin/srl" 51a≡
    cat | tee $INPUTFILE | python nafAlpinoToSRLFeatures.py > $FEATUREVECTOR
    ◇
```

File defined by 50e, 51abcd.

Run the trained model on the feature-vector.

```
"../bin/srl" 51b≡
    timbl -m0:I1,2,3,4 -i 25Feb2015_e-mags_mags_press_newspapers.wgt -
    t $FEATUREVECTOR -o $TIMBLOUTPUTFILE >/dev/null 2>/dev/null
    ◇
```

File defined by 50e, 51abcd.

Insert the SRL values into the NAF file.

```
"../bin/srl" 51c≡
    python timblToAlpinoNAF.py $INPUTFILE $TIMBLOUTPUTFILE
    ◇
```

File defined by 50e, 51abcd.

Clean up.

```
"../bin/srl" 51d≡
    rm -rf $TEMPDIR
    ◇
```

File defined by 50e, 51abcd.

4.5.25 SRL postprocessing

In addition to the Semantic Role Labeling there is hack that finds additional semantic roles.

Module Get the module from Github. Note that this module needs rdflib

```
<install python packages 51e> ≡
    pip install rdflib
    ◇
```

Fragment defined by 13f, 51e.

Fragment referenced in 10b.

Defines: rdflib Never used.

```
<install the post-SRL module 51f> ≡
    cd $modulesdir
    if
    [ -d vua-srl-postprocess ]
    then
    cd vua-srl-postprocess
    git pull
    else
    git clone https://github.com/newsreader/vua-srl-postprocess.git
    cd vua-srl-postprocess
    fi
    ◇
```

Fragment referenced in 18q.

Script

```
"../bin/postsr1" 52a≡
  ⟨ start of module-script (52b vua-srl-postprocess ) 20b ⟩
  cd $MODDIR
  tempdir='mktemp -d -t postsr1.XXXXX'
  cat >$tempdir/infile
  python $MODDIR/main.py -i $tempdir/infile -o $tempdir/outfile
  cat $tempdir/outfile
  rm -rf $tempdir
  ◇
```

4.5.26 Event coreference

The event-coreference module is language-independent. Although the version in the EHU-repo is 3.0, the version 2.0 used in this pipeline seems to be more recent, so we will use that.

Module Install the module from the snapshot.

```
⟨ install the event-coreference module 52c ⟩ ≡
  cd $modulesdir
  tar -xzf $snapshotsocket/t_nlpp_resources/20151217_vua-eventcoreference_v2.tgz
  cd vua-eventcoreference_v2
  cp lib/EventCoreference-1.0-SNAPSHOT-jar-with-dependencies.jar $jarsdir
  ◇
```

Fragment referenced in 18h.

Script

```
"../bin/evcoref" 52d≡
  ⟨ start of module-script (52e vua-eventcoreference_v2 ) 20b ⟩
  RESOURCESDIR=$MODDIR/resources
  JARFILE=$jarsdir/EventCoreference-1.0-SNAPSHOT-jar-with-dependencies.jar

  if
    [ "$naflang" == 'nl' ]
  then
    lang_resource="odwn_orbn_gwg-LMF_1.3.xml"
  else
    lang_resource="wneng-30.lmf.xml.xpos"
  fi

  JAVAMODULE=eu.newsreader.eventcoreference.naf.EventCorefWordnetSim
  JAVAOPTIONS="--method leacock-chodorow"
  JAVAOPTIONS="$JAVAOPTIONS --wn-lmf $RESOURCESDIR/$lang_resource"
  JAVAOPTIONS="$JAVAOPTIONS --sim 2.0"
  JAVAOPTIONS="$JAVAOPTIONS --wsd 0.8"
  JAVAOPTIONS="$JAVAOPTIONS --
relations XPOS_NEAR_SYNONYM#HAS_HYPERONYM#HAS_XPOS_HYPERONYM#event"

  java -Xmx812m -cp $JARFILE $JAVAMODULE $JAVAOPTIONS
  ◇
```

4.5.27 Dbpedia-ner

Dbpedia-ner finds more named entities than NER, because it checks DBpedia for the candidate NE-'s.

Module

```
<install the dbpedia-ner module 53a> ≡
MODNAM=dbpedia_ner
DIRN=dbpedia_ner
GITU=https://github.com/PaulHuygen/dbpedia_ner.git
GITC=ab1dcdbd860f0ff29bc979f646dc382122a101fc2
<install from github 6d>
◇
```

Fragment referenced in 18q.

Script The main part of the module is a Python script. The README.md file of the Github repo lists the options that can be applied. One of the options is about the URL of the Spotlight server.

```
"../bin/dbpner" 53b≡
<start of module-script (53c dbpedia_ner ) 20b>
cat | iconv -f ISO8859-1 -t UTF-8 | $MODDIR/dbpedia_ner.py -
url http://$spotlighthost:2060/rest/candidates
◇
```

4.5.28 Nominal events

The module “postprocessing-nl” adds nominal events to the srl annotations. It has been obtained directly from the author (Piek Vossen). It is not yet available in a public repo. Probably in future versions the jar from the ontotagger module can be used for this module.

Module

```
<install the nomevent module 53d> ≡
cd $modulesdir
tar -xzf $snapshotsocket/t_nlpp_resources/20160126_vua-nominal-event-detection-
nl.tgz
◇
```

Fragment referenced in 18q.

Script

```
"../bin/nomevent" 53e≡
<start of module-script (53f vua-nominal-event-detection-nl ) 20b>
LIBDIR=$MODDIR/lib
RESOURCESDIR=$MODDIR/resources

JAR=$LIBDIR/ontotagger-1.0-jar-with-dependencies.jar
JAVAMODULE=eu.kyotoproject.main.NominalEventCoreference
cat | java -Xmx812m -cp $JAR $JAVAMODULE --framenet-lu $RESOURCESDIR/nl-luIndex.xml
◇
```

4.5.29 Opinion miner

To run the opinion-miner, the following things are needed:

- SVMlight
- crfsuite
- vua-pylib

Module The module can be cloned from Github. However, currently there are problems with the Github installation. Therefore we borrow the opinion miner from the English NWR pipeline.

```
<install the opinion-miner 54a> ≡
  cd /home/phuijgen/nlp/test/nlpp/modules
  tar -xzf /home/phuijgen/nlp/test/t_nlpp_resources/20151012VUA-opinion-miner.tgz
  ◇
```

Fragment defined by 54ac.

Fragment referenced in 19a.

The opinion-miner needs a configuration file that is located in the directory where the model-data resides. In this pipeline we will use model-data derived from news-articles. An alternative model, derived from hotel evaluations can also be used. Put the configuration file in the `etc` subdir and copy it to its proper location during the installation of the opinion-miner.

```
"../env/etc/opini.cfg" 54b≡
[general]
output_folder = /home/phuijgen/nlp/test/nlpp/modules/VUA-opinion-
miner/final_models/ennl/news_cfg1

[crfsuite]
path_to_binary = /home/phuijgen/nlp/test/nlpp/env/bin/crfsuite

[svmlight]
path_to_binary_learn = /home/phuijgen/nlp/test/nlpp/env/bin/svm_learn
path_to_binary_classify = /home/phuijgen/nlp/test/nlpp/env/bin/svm_classify
  ◇
```

```
<install the opinion-miner 54c> ≡
  cd VUA-opinion-miner
  cat /home/phuijgen/nlp/test/nlpp/env/etc/opini.cfg | \
    sed s/ennl/nl/g > $modulesdir/VUA-opinion-
miner/final_models/nl/news_cfg1/config.cfg
  cat /home/phuijgen/nlp/test/nlpp/env/etc/opini.cfg | \
    sed s/ennl/en/g > $modulesdir/VUA-opinion-
miner/final_models/en/news_cfg1/config.cfg
  ◇
```

Fragment defined by 54ac.

Fragment referenced in 19a.

Script


```

"../bin/opinimin" 55a≡
  < start of module-script (55b VUA-opinion-miner ) 20b >
  cd $MODDIR
  export PATH=$PATH:.
  if
    [ "$naflang" == "nl" ]
  then
    modelconf=$MODDIR/final_models/nl/news_cfg1
  else
    modelconf=$MODDIR/final_models/en/news_cfg1
  fi
  python classify_kaf_naf_file.py -m $modelconf
  ◇

```

5 Utilities

5.1 Test script

The following script pushes a test-document through the modules of the pipeline.

```

"../bin/test" 55c≡
  #!/bin/bash
  ROOT=/home/phuijgen/nlp/test/nlpp
  TESTDIR=$ROOT/test
  < function to check/start spotlight 26b >
  TESTIN=$ROOT/nuweb/test.nl.in.naf
  if
    [ "$1" == "en" ]
  then
    TESTIN=$ROOT/nuweb/test.en.in.naf
  fi
  BIND=$ROOT/bin
  mkdir -p $TESTDIR
  cd $TESTDIR
  < set the language variable (55d $TESTIN ) 21d >
  < select language-dependent features 42c >
  check_start_spotlight $naflang
  if
    [ "$naflang" == "nl" ]
  then
    cat $TESTIN      | $BIND/tok                > tok.naf
    cat tok.naf      | $BIND/mor                > mor.naf
    cat mor.naf      | $BIND/nerc               > nerc.naf
    cat nerc.naf     | $BIND/wsd                > wsd.naf
    cat wsd.naf      | $BIND/ned                > ned.naf
    cat ned.naf      | $BIND/heideltime         > times.naf
    cat times.naf    | $BIND/onto               > onto.naf
    cat onto.naf     | $BIND/srl                > srl.naf
    cat srl.naf      | $BIND/nomevent           > nomev.naf
    cat nomev.naf    | $BIND/srl-dutch-nominals > psrl.naf
    cat psrl.naf     | $BIND/framesrl          > fsrl.naf
    cat fsrl.naf     | $BIND/opinimin           > opin.naf
    cat opin.naf     | $BIND/evcoref            > out.naf
  else
    < annotate english document 56b >
  fi
  ◇

```

Correct sequence of the modules in the Dutch pipeline:

- tok
- mor
- nerc
- wsd
- ned
- heidel
- onto (predicate-matrix-tagger.sh uit vua-ontotagger-v1.0)
- srl
- Nominal event detectie
- vua-srl-extra
- framesrl (srl-framenet-tagger.sh uit vua-ontotagger-v1.0)
- opinion mining
- ecrf

```

< annotate dutch document 56a > ≡
cat $TESTIN      | $BIND/tok                > tok.naf
cat tok.naf      | $BIND/mor                > mor.naf
cat mor.naf      | $BIND/nerc               > nerc.naf
cat nerc.naf     | $BIND/wsd               > wsd.naf
cat wsd.naf      | $BIND/ned               > ned.naf
cat ned.naf      | $BIND/heideltime        > times.naf
cat times.naf    | $BIND/onto              > onto.naf
cat onto.naf     | $BIND/srl               > srl.naf
cat srl.naf      | $BIND/nomevent          > nomev.naf
cat nomev.naf    | $BIND/postsr           > psrl.naf
cat psrl.naf     | $BIND/framesrl          > fsrl.naf
cat fsrl.naf     | $BIND/opinimin          > opin.naf
cat opin.naf     | $BIND/evcoref           > out.naf
◇

```

Fragment never referenced.

```

< annotate english document 56b > ≡
cat $TESTIN      | $BIND/tok                > tok.naf
cat tok.naf      | $BIND/topic             > top.naf
cat top.naf      | $BIND/pos               > pos.naf
cat pos.naf      | $BIND/constpars         > consp.naf
cat consp.naf    | $BIND/nerc              > nerc.naf
cat nerc.naf     | $BIND/nedrer            > nedr.naf
cat nedr.naf     | $BIND/wikify            > wikif.naf
cat wikif.naf    | $BIND/ukb              > ukb.naf
cat ukb.naf      | $BIND/ewsd             > ewsd.naf
cat ewsd.naf     | $BIND/esrl             > esrl.naf
cat esrl.naf     | $BIND/FBK-time          > time.naf
cat time.naf     | $BIND/FBK-temprel       > trel.naf
cat trel.naf     | $BIND/FBK-causalrel     > crel.naf
cat crel.naf     | $BIND/evcoref           > ecrf.naf
cat ecrf.naf     | $BIND/factuality        > fact.naf
cat fact.naf     | $BIND/opinimin          > out.naf
◇

```

Fragment referenced in 55c.

5.2 Logging

Write log messages to standard out if variable LOGLEVEL is equal to 1.

```
< variables of install-modules 57a > ≡
    LOGLEVEL=1
    ◇
```

Fragment referenced in 15a.

```
< logmess 57b > ≡
    if
    [ $LOGLEVEL -gt 0 ]
    then
    echo @1
    fi
    ◇
```

Fragment referenced in 6bd, 13a, 57c.

5.3 Misc

Install a module from a tarball: The macro expects the following three variables to be present:

URL: The URL tfrom where the taball can be downloaded.

TARB: The name of the tarball.

DIR; Name of the directory for the module.

Arg 1: URL; Arg 2: tarball; Arg 3: directory.

```
< install from tarball 57c > ≡
    SUCCES=0
    cd $modulesdir
    < move module (57d $DIR ) 5c >
    wget $URL
    SUCCES=$?
    if
    [ $SUCCES -eq 0 ]
    then
    tar -xzf $TARB
    SUCCES=$?
    rm -rf $TARB
    fi
    if
    [ $SUCCES -eq 0 ]
    then
    < logmess (57e Installed $DIR ) 57b >
    < remove old module (57f $DIR ) 6a >
    else
    < re-instate old module (57g $DIR ) 6b >
    fi
    ◇
```

Fragment never referenced.

A How to read and translate this document

This document is an example of *literate programming* [?]. It contains the code of all sorts of scripts and programs, combined with explaining texts. In this document the literate programming tool **nuweb** is used, that is currently available from Sourceforge (URL:nuweb.sourceforge.net). The advantages of Nuweb are, that it can be used for every programming language and scripting language, that it can contain multiple program sources and that it is very simple.

A.1 Read this document

The document contains *code scraps* that are collected into output files. An output file (e.g. `output.fil`) shows up in the text as follows:

```
"output.fil" 4a ≡
    # output.fil
    < a macro 4b >
    < another macro 4c >
    ◇
```

The above construction contains text for the file. It is labelled with a code (in this case 4a) The constructions between the < and > brackets are macro's, placeholders for texts that can be found in other places of the document. The test for a macro is found in constructions that look like:

```
< a macro 4b > ≡
    This is a scrap of code inside the macro.
    It is concatenated with other scraps inside the
    macro. The concatenated scraps replace
    the invocation of the macro.
```

Macro defined by 4b, 87e

Macro referenced in 4a

Macro's can be defined on different places. They can contain other macro's.

```
< a scrap 87e > ≡
    This is another scrap in the macro. It is
    concatenated to the text of scrap 4b.
    This scrap contains another macro:
    < another macro 45b >
```

Macro defined by 4b, 87e

Macro referenced in 4a

A.2 Process the document

The raw document is named `a_nlpp.w`. Figure 2 shows pathways to translate it into printable/viewable documents and to extract the program sources. Table 4 lists the tools that are

Tool	Source	Description
gawk	www.gnu.org/software/gawk/	text-processing scripting language
M4	www.gnu.org/software/m4/	Gnu macro processor
nuweb	nuweb.sourceforge.net	Literate programming tool
tex	www.ctan.org	Typesetting system
tex4ht	www.ctan.org	Convert \TeX documents into xml/html

Table 4: Tools to translate this document into readable code and to extract the program sources

needed for a translation. Most of the tools (except Nuweb) are available on a well-equipped Linux system.

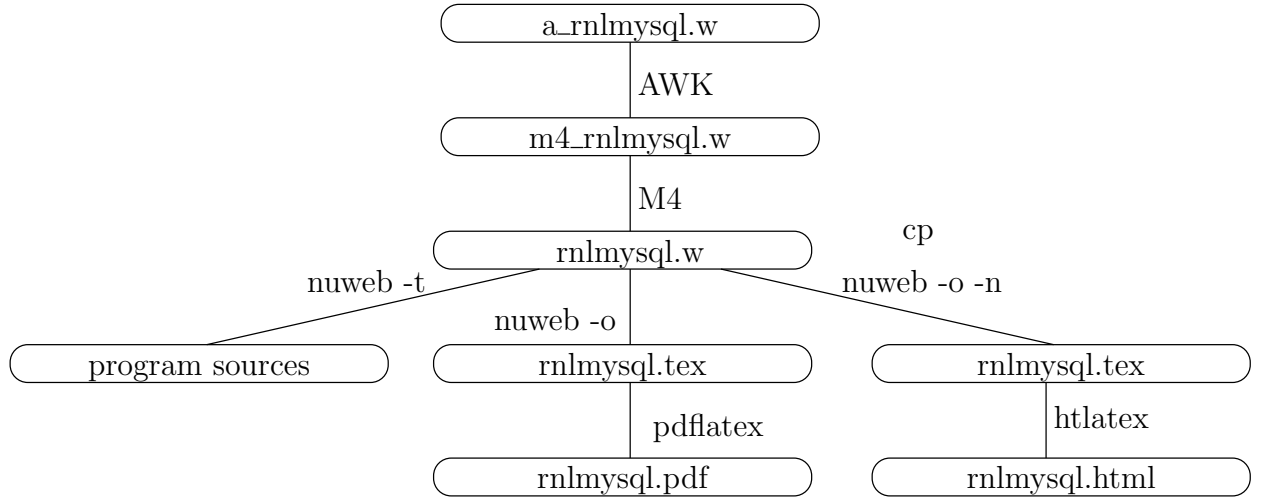


Figure 2: Translation of the raw code of this document into printable/viewable documents and into program sources. The figure shows the pathways and the main files involved.

< parameters in Makefile 59a > \equiv
 NUWEB=../env/bin/nuweb
 \diamond

Fragment defined by 59a, 60c, 62ab, 64d, 67a, 69d.
 Fragment referenced in 59b.
 Uses: nuweb 66b.

A.3 The Makefile for this project.

This chapter assembles the Makefile for this project.

```

"Makefile" 59b  $\equiv$ 
  < default target 59c >

  < parameters in Makefile 59a, ... >

  < impliciete make regels 63a, ... >
  < expliciete make regels 60d, ... >
  < make targets 60a, ... >
   $\diamond$ 

```

The default target of make is **all**.

```

< default target 59c >  $\equiv$ 
  all : < all targets 60b >
  .PHONY : all
   $\diamond$ 

```

Fragment referenced in 59b.
 Defines: all 29b, PHONY 63b.

```

< make targets 60a > ≡
    clean:
        < clean up 8f, ... >

```

◇

Fragment defined by 60a, 64ab, 67e, 70acd.
 Fragment referenced in 59b.

One of the targets is certainly the PDF version of this document.

```

< all targets 60b > ≡
    nlpp.pdf◇

```

Fragment referenced in 59c.
 Uses: pdf 64a.

We use many suffixes that were not known by the C-programmers who constructed the `make` utility. Add these suffixes to the list.

```

< parameters in Makefile 60c > ≡
    .SUFFIXES: .pdf .w .tex .html .aux .log .php

```

◇

Fragment defined by 59a, 60c, 62ab, 64d, 67a, 69d.
 Fragment referenced in 59b.
 Defines: SUFFIXES Never used.
 Uses: pdf 64a.

A.4 Get Nuweb

An annoying problem is, that this program uses nuweb, a utility that is seldom installed on a computer. Therefore, we are going to install that first if it is not present. Unfortunately, nuweb is hosted on sourceforge and it is difficult to achieve automatic downloading from that repository. Therefore I copied one of the versions on a location from where it can be downloaded with a script.

Put the nuweb binary in the nuweb subdirectory, so that it can be used before the directory-structure has been generated.

```

< expliciete make regels 60d > ≡

    nuweb: $(NUWEB)

    $(NUWEB): ../nuweb-1.58
        mkdir -p ../env/bin
        cd ../nuweb-1.58 && make nuweb
        cp ../nuweb-1.58/nuweb $(NUWEB)

```

◇

Fragment defined by 60d, 61bcd, 63b, 65a, 67bd.
 Fragment referenced in 59b.
 Uses: nuweb 66b.

```

⟨ clean up 61a ⟩ ≡
    rm -rf ../nuweb-1.58
    ◇

```

Fragment defined by 8f, 9d, 22b, 61a.
 Fragment referenced in 60a.
 Uses: nuweb 66b.

```

⟨ expliciete make regels 61b ⟩ ≡
    ../nuweb-1.58:
        cd .. && wget http://kyoto.let.vu.nl/~huygen/nuweb-1.58.tgz
        cd .. && tar -xzf nuweb-1.58.tgz
    ◇

```

Fragment defined by 60d, 61bcd, 63b, 65a, 67bd.
 Fragment referenced in 59b.
 Uses: nuweb 66b.

A.5 Pre-processing

To make usable things from the raw input `a_nlpp.w`, do the following:

1. Process `$` characters.
2. Run the m4 pre-processor.
3. Run nuweb.

This results in a \LaTeX file, that can be converted into a PDF or a HTML document, and in the program sources and scripts.

A.5.1 Process ‘dollar’ characters

Many “intelligent” \TeX editors (e.g. the `auctex` utility of Emacs) handle `$` characters as special, to switch into mathematics mode. This is irritating in program texts, that often contain `$` characters as well. Therefore, we make a stub, that translates the two-character sequence `\$` into the single `$` character.

```

⟨ expliciete make regels 61c ⟩ ≡
    m4_nlpp.w : a_nlpp.w
        gawk '{if(match($0, "@%")) {printf("%s", substr($0,1,RSTART-
1))} else print}' a_nlpp.w \
        | gawk '{gsub(/\[\[\] [\$\$]/, "$$");print}' > m4_nlpp.w
    ◇

```

Fragment defined by 60d, 61bcd, 63b, 65a, 67bd.
 Fragment referenced in 59b.
 Uses: print 64a.

A.5.2 Run the M4 pre-processor

```

⟨ expliciete make regels 61d ⟩ ≡
    nlpp.w : m4_nlpp.w inst.m4
        m4 -P m4_nlpp.w > nlpp.w
    ◇

```

Fragment defined by 60d, 61bcd, 63b, 65a, 67bd.
 Fragment referenced in 59b.

A.6 Typeset this document

Enable the following:

1. Create a PDF document.
2. Print the typeset document.
3. View the typeset document with a viewer.
4. Create a HTMLdocument.

In the three items, a typeset PDF document is required or it is the requirement itself.

A.6.1 Figures

This document contains figures that have been made by `xfig`. Post-process the figures to enable inclusion in this document.

The list of figures to be included:

```
<parameters in Makefile 62a> ≡
    FIGFILES=fileschema directorystructure
```

◇

Fragment defined by 59a, 60c, 62ab, 64d, 67a, 69d.

Fragment referenced in 59b.

Defines: FIGFILES 62b.

We use the package `figlatex` to include the pictures. This package expects two files with extensions `.pdftex` and `.pdftex_t` for `pdflatex` and two files with extensions `.pstex` and `.pstex_t` for the `latex/dvips` combination. Probably `tex4ht` uses the latter two formats too.

Make lists of the graphical files that have to be present for `latex/pdflatex`:

```
<parameters in Makefile 62b> ≡
    FIGFILENAMES=$(foreach fil,$(FIGFILES), $(fil).fig)
    PDFT_NAMES=$(foreach fil,$(FIGFILES), $(fil).pdftex_t)
    PDF_FIG_NAMES=$(foreach fil,$(FIGFILES), $(fil).pdftex)
    PST_NAMES=$(foreach fil,$(FIGFILES), $(fil).pstex_t)
    PS_FIG_NAMES=$(foreach fil,$(FIGFILES), $(fil).pstex)
```

◇

Fragment defined by 59a, 60c, 62ab, 64d, 67a, 69d.

Fragment referenced in 59b.

Defines: FIGFILENAMES Never used, PDFT_NAMES 64b, PDF_FIG_NAMES 64b, PST_NAMES Never used,
PS_FIG_NAMES Never used.

Uses: FIGFILES 62a.

Create the graph files with program `fig2dev`:


```

⟨ impliciete make regels 63a ⟩ ≡
    %.eps: %.fig
        fig2dev -L eps $< > $@

    %.pstex: %.fig
        fig2dev -L pstex $< > $@

    .PRECIOUS : %.pstex
    %.pstex_t: %.fig %.pstex
        fig2dev -L pstex_t -p $*.pstex $< > $@

    %.pdftex: %.fig
        fig2dev -L pdftex $< > $@

    .PRECIOUS : %.pdftex
    %.pdftex_t: %.fig %.pstex
        fig2dev -L pdftex_t -p $*.pdftex $< > $@

```

◇

Fragment defined by 63a, 67c.

Fragment referenced in 59b.

Defines: fig2dev Never used.

A.6.2 Bibliography

To keep this document portable, create a portable bibliography file. It works as follows: This document refers in the |bibliography| statement to the local bib-file **nlpp.bib**. To create this file, copy the auxiliary file to another file **auxfil.aux**, but replace the argument of the command **\bibdata{nlpp}** to the names of the bibliography files that contain the actual references (they should exist on the computer on which you try this). This procedure should only be performed on the computer of the author. Therefore, it is dependent of a binary file on his computer.

```

⟨ expliciete make regels 63b ⟩ ≡
    bibfile : nlpp.aux /home/paul/bin/mkportbib
        /home/paul/bin/mkportbib nlpp litprog

    .PHONY : bibfile

```

◇

Fragment defined by 60d, 61bcd, 63b, 65a, 67bd.

Fragment referenced in 59b.

Uses: PHONY 59c.

A.6.3 Create a printable/viewable document

Make a PDF document for printing and viewing.

```

⟨ make targets 64a ⟩ ≡
    pdf : nlpp.pdf

    print : nlpp.pdf
           lpr nlpp.pdf

    view : nlpp.pdf
          evince nlpp.pdf

```

◇

Fragment defined by 60a, 64ab, 67e, 70acd.

Fragment referenced in 59b.

Defines: pdf 60bc, 64b, print 10c, 12c, 14e, 20c, 31a, 48a, 61c, view Never used.

Create the PDF document. This may involve multiple runs of nuweb, the L^AT_EX processor and the bibT_EX processor, and depends on the state of the aux file that the L^AT_EX processor creates as a by-product. Therefore, this is performed in a separate script, w2pdf.

The w2pdf script The three processors nuweb, L^AT_EX and bibT_EX are intertwined. L^AT_EX and bibT_EX create parameters or change the value of parameters, and write them in an auxiliary file. The other processors may need those values to produce the correct output. The L^AT_EX processor may even need the parameters in a second run. Therefore, consider the creation of the (PDF) document finished when none of the processors causes the auxiliary file to change. This is performed by a shell script w2pdf.

```

⟨ make targets 64b ⟩ ≡
    nlpp.pdf : nlpp.w $(W2PDF) $(PDF_FIG_NAMES) $(PDFT_NAMES)
              chmod 775 $(W2PDF)
              $(W2PDF) $*

```

◇

Fragment defined by 60a, 64ab, 67e, 70acd.

Fragment referenced in 59b.

Uses: pdf 64a, PDFT_NAMES 62b, PDF_FIG_NAMES 62b.

The following is an ugly fix of an unsolved problem. Currently I develop this thing, while it resides on a remote computer that is connected via the sshfs filesystem. On my home computer I cannot run executables on this system, but on my work-computer I can. Therefore, place the following script on a local directory.

```

⟨ directories to create 64c ⟩ ≡
    ../nuweb/bin ◇

```

Fragment defined by 4abcd, 8bh, 9a, 11e, 64c.

Fragment referenced in 70a.

Uses: nuweb 66b.

```

⟨ parameters in Makefile 64d ⟩ ≡
    W2PDF=../nuweb/bin/w2pdf

```

◇

Fragment defined by 59a, 60c, 62ab, 64d, 67a, 69d.

Fragment referenced in 59b.

Uses: nuweb 66b.

```

< expliciete make regels 65a > ≡
    $(W2PDF) : nlpp.w $(NUWEB)
              $(NUWEB) nlpp.w

```

◇

Fragment defined by 60d, 61bcd, 63b, 65a, 67bd.
 Fragment referenced in 59b.

```

"../nuweb/bin/w2pdf" 65b≡
    #!/bin/bash
    # w2pdf -- compile a nuweb file
    # usage: w2pdf [filename]
    # 20160218 at 1054h: Generated by nuweb from a_nlpp.w
    NUWEB=../env/bin/nuweb
    LATEXCOMPILER=pdflatex
    < filenames in nuweb compile script 65d >
    < compile nuweb 65c >

```

◇

Uses: nuweb 66b.

The script retains a copy of the latest version of the auxiliary file. Then it runs the four processors nuweb, L^AT_EX, MakeIndex and bibT_EX, until they do not change the auxiliary file or the index.

```

< compile nuweb 65c > ≡
    NUWEB=/home/phuijgen/nlp/test/nlpp/env/bin/nuweb
    < run the processors until the aux file remains unchanged 66c >
    < remove the copy of the aux file 66a >

```

◇

Fragment referenced in 65b.
 Uses: nuweb 66b.

The user provides the name of the nuweb file as argument. Strip the extension (e.g. .w) from the filename and create the names of the L^AT_EX file (ends with .tex), the auxiliary file (ends with .aux) and the copy of the auxiliary file (add old. as a prefix to the auxiliary filename).

```

< filenames in nuweb compile script 65d > ≡
    nufil=$1
    trunk=${1%.*}
    texfil=${trunk}.tex
    auxfil=${trunk}.aux
    oldaux=old.${trunk}.aux
    indexfil=${trunk}.idx
    oldindexfil=old.${trunk}.idx

```

◇

Fragment referenced in 65b.
 Defines: auxfil 66c, 68c, 69a, indexfil 66c, 68c, nufil 66b, 68c, 69b, oldaux 66ac, 68c, 69a, oldindexfil 66c, 68c, texfil 66b, 68c, 69b, trunk 66b, 68c, 69bc.

Remove the old copy if it is no longer needed.

```

⟨ remove the copy of the aux file 66a ⟩ ≡
    rm $oldaux
    ◇

```

Fragment referenced in 65c, 68b.
 Uses: oldaux 65d, 68c.

Run the three processors. Do not use the option `-o` (to suppress generation of program sources) for nuweb, because `w2pdf` must be kept up to date as well.

```

⟨ run the three processors 66b ⟩ ≡
    $NUWEB $nufil
    $LATEXCOMPILER $texfil
    makeindex $trunk
    bibtex $trunk
    ◇

```

Fragment referenced in 66c.
 Defines: bibtex 69bc, makeindex 69bc, nuweb 4e, 55c, 59a, 60d, 61ab, 64cd, 65bc, 67a, 68a.
 Uses: nufil 65d, 68c, texfil 65d, 68c, trunk 65d, 68c.

Repeat to copy the auxiliary file and the index file and run the processors until the auxiliary file and the index file are equal to their copies. However, since I have not yet been able to test the `aux` file and the `idx` in the same test statement, currently only the `aux` file is tested.

It turns out, that sometimes a strange loop occurs in which the `aux` file will keep to change. Therefore, with a counter we prevent the loop to occur more than 10 times.

```

⟨ run the processors until the aux file remains unchanged 66c ⟩ ≡
    LOOPCOUNTER=0
    while
        ! cmp -s $auxfil $oldaux
    do
        if [ -e $auxfil ]
        then
            cp $auxfil $oldaux
        fi
        if [ -e $indexfil ]
        then
            cp $indexfil $oldindexfil
        fi
        ⟨ run the three processors 66b ⟩
        if [ $LOOPCOUNTER -ge 10 ]
        then
            cp $auxfil $oldaux
        fi;
    done
    ◇

```

Fragment referenced in 65c.
 Uses: auxfil 65d, 68c, indexfil 65d, oldaux 65d, 68c, oldindexfil 65d.

A.6.4 Create HTML files

HTML is easier to read on-line than a PDF document that was made for printing. We use `tex4ht` to generate HTML code. An advantage of this system is, that we can include figures in the same way as we do for `pdflatex`.

To create a HTML doc, we do the following:

1. Create a directory `../nuweb/html` for the HTML document.
2. Put the nuweb source in it, together with style-files that are needed (see variable `HTMLSOURCE`).
3. Put the script `w2html` in it and make it executable.
4. Execute the script `w2html`.

Make a list of the entities that we mentioned above:

```
<parameters in Makefile 67a> ≡
    htmldir=../nuweb/html
    htmlsource=nlpp.w nlpp.bib html.sty artikel3.4ht w2html
    htmlmaterial=$(foreach fil, $(htmlsource), $(htmldir)/$(fil))
    htmltarget=$(htmldir)/nlpp.html
◇
```

Fragment defined by 59a, 60c, 62ab, 64d, 67a, 69d.

Fragment referenced in 59b.

Uses: nuweb 66b.

Make the directory:

```
<expliciete make regels 67b> ≡
    $(htmldir) :
        mkdir -p $(htmldir)
◇
```

Fragment defined by 60d, 61bcd, 63b, 65a, 67bd.

Fragment referenced in 59b.

The rule to copy files in it:

```
<impliciete make regels 67c> ≡
    $(htmldir)/% : % $(htmldir)
        cp $< $(htmldir)/
◇
```

Fragment defined by 63a, 67c.

Fragment referenced in 59b.

Do the work:

```
<expliciete make regels 67d> ≡
    $(htmltarget) : $(htmlmaterial) $(htmldir)
        cd $(htmldir) && chmod 775 w2html
        cd $(htmldir) && ./w2html nlpp.w
◇
```

Fragment defined by 60d, 61bcd, 63b, 65a, 67bd.

Fragment referenced in 59b.

Invoke:

```
<make targets 67e> ≡
    htm : $(htmldir) $(htmltarget)
◇
```

Fragment defined by 60a, 64ab, 67e, 70acd.

Fragment referenced in 59b.

Create a script that performs the translation.

```
"w2html" 68a≡
#!/bin/bash
# w2html -- make a html file from a nuweb file
# usage: w2html [filename]
# [filename]: Name of the nuweb source file.
# 20160218 at 1054h: Generated by nuweb from a_nlpp.w
echo "translate " $1 >w2html.log
NUWEB=/home/phuijgen/nlp/test/nlpp/env/bin/nuweb
⟨filenames in w2html 68c⟩

⟨perform the task of w2html 68b⟩
```

◇

Uses: **nuweb** 66b.

The script is very much like the **w2pdf** script, but at this moment I have still difficulties to compile the source smoothly into HTML and that is why I make a separate file and do not recycle parts from the other file. However, the file works similar.

```
⟨perform the task of w2html 68b⟩ ≡
  ⟨run the html processors until the aux file remains unchanged 69a⟩
  ⟨remove the copy of the aux file 66a⟩
```

◇

Fragment referenced in 68a.

The user provides the name of the nuweb file as argument. Strip the extension (e.g. **.w**) from the filename and create the names of the L^AT_EX file (ends with **.tex**), the auxiliary file (ends with **.aux**) and the copy of the auxiliary file (add **old.** as a prefix to the auxiliary filename).

```
⟨filenames in w2html 68c⟩ ≡
nufil=$1
trunk=${1%.*}
texfil=${trunk}.tex
auxfil=${trunk}.aux
oldaux=old.${trunk}.aux
indexfil=${trunk}.idx
oldindexfil=old.${trunk}.idx
```

◇

Fragment referenced in 68a.

Defines: **auxfil** 65d, 66c, 69a, **nufil** 65d, 66b, 69b, **oldaux** 65d, 66ac, 69a, **texfil** 65d, 66b, 69b, **trunk** 65d, 66b, 69bc.

Uses: **indexfil** 65d, **oldindexfil** 65d.

```

⟨run the html processors until the aux file remains unchanged 69a⟩ ≡
    while
        ! cmp -s $auxfil $oldaux
    do
        if [ -e $auxfil ]
        then
            cp $auxfil $oldaux
        fi
        ⟨run the html processors 69b⟩
    done
    ⟨run tex4ht 69c⟩

```

◇

Fragment referenced in 68b.

Uses: auxfil 65d, 68c, oldaux 65d, 68c.

To work for HTML, nuweb *must* be run with the `-n` option, because there are no page numbers.

```

⟨run the html processors 69b⟩ ≡
    $NUWEB -o -n $nufil
    latex $texfil
    makeindex $trunk
    bibtex $trunk
    htlatex $trunk

```

◇

Fragment referenced in 69a.

Uses: bibtex 66b, makeindex 66b, nufil 65d, 68c, texfil 65d, 68c, trunk 65d, 68c.

When the compilation has been satisfied, run makeindex in a special way, run bibtex again (I don't know why this is necessary) and then run htlatex another time.

```

⟨run tex4ht 69c⟩ ≡
    tex '\def\filename{{\nlpp}{\idx}{4dx}{ind}} \input idxmake.4ht'
    makeindex -o $trunk.ind $trunk.4dx
    bibtex $trunk
    htlatex $trunk

```

◇

Fragment referenced in 69a.

Uses: bibtex 66b, makeindex 66b, trunk 65d, 68c.

A.7 Create the program sources

Run nuweb, but suppress the creation of the L^AT_EX documentation. Nuweb creates only sources that do not yet exist or that have been modified. Therefore make does not have to check this. However, “make” has to create the directories for the sources if they do not yet exist. So, let's create the directories first.

```

⟨parameters in Makefile 69d⟩ ≡
    MKDIR = mkdir -p

```

◇

Fragment defined by 59a, 60c, 62ab, 64d, 67a, 69d.

Fragment referenced in 59b.

Defines: MKDIR 70a.

```

< make targets 70a > ≡
    DIRS = < directories to create 4a, ... >

```

```

    $(DIRS) :
        $(MKDIR) $@

```

◇

Fragment defined by 60a, 64ab, 67e, 70acd.

Fragment referenced in 59b.

Defines: DIRS 70c.

Uses: MKDIR 69d.

```

< make scripts executable 70b > ≡
    chmod -R 775 ../bin/*
    chmod -R 775 ../env/bin/*

```

◇

Fragment defined by 19d, 21c, 70b.

Fragment referenced in 70c.

```

< make targets 70c > ≡
    sources : nlpp.w $(DIRS) $(NUWEB)
              $(NUWEB) nlpp.w
              < make scripts executable 19d, ... >

```

◇

Fragment defined by 60a, 64ab, 67e, 70acd.

Fragment referenced in 59b.

Uses: DIRS 70a.

A.8 Restore paths after transplantation

When an existing installation has been transplanted to another location, many path indications have to be adapted to the new situation. The scripts that are generated by nuweb can be repaired by re-running nuweb. After that, configuration files of some modules must be modified.

```

< make targets 70d > ≡
    transplant :
        touch a_nlpp.w
        $(MAKE) sources
        ../env/bin/transplant

```

◇

Fragment defined by 60a, 64ab, 67e, 70acd.

Fragment referenced in 59b.

In order to work as expected, the following script must be re-made after a transplantation.


```

"../env/bin/transplant" 71≡
    #!/bin/bash
    LOGLEVEL=1
    < set variables that point to the directory-structure 4e, ... >
    < set paths after transplantation 12c >
    < re-install modules after the transplantation 24d >

◇

```

B References

B.1 Literature

C Indexes

C.1 Filenames

```

"../bin/constpars" Defined by 33a.
"../bin/coreference-base" Defined by 41d.
"../bin/dbpner" Defined by 53b.
"../bin/eSRL" Defined by 35h.
"../bin/evcoref" Defined by 52d.
"../bin/ewsd" Defined by 35a.
"../bin/factuality" Defined by 41a.
"../bin/FBK-causalrel" Defined by 40d.
"../bin/FBK-temprel" Defined by 39d.
"../bin/FBK-time" Defined by 38a.
"../bin/framesrl" Defined by 48a.
"../bin/heideltime" Defined by 50b.
"../bin/install-modules" Defined by 15a, 16aj, 17aj, 18ahq, 19a.
"../bin/langdetect" Defined by 21a.
"../bin/lu2synset" Defined by 45a.
"../bin/mor" Defined by 31e.
"../bin/ned" Defined by 46b.
"../bin/nedrer" Defined by 33d.
"../bin/nerc" Defined by 43d.
"../bin/nerc_conll02" Defined by 43b.
"../bin/nomevent" Defined by 53e.
"../bin/onto" Defined by 47a.
"../bin/opinimin" Defined by 55a.
"../bin/pos" Defined by 32c.
"../bin/postsrl" Defined by 52a.
"../bin/srl" Defined by 50e, 51abcd.
"../bin/srl-dutch-nominals" Defined by 36d.
"../bin/start_eSRL" Defined by 35d.
"../bin/stop_eSRL" Defined by 35f.
"../bin/test" Defined by 55c.
"../bin/tok" Defined by 30c.
"../bin/topic" Defined by 31b.
"../bin/ukb" Defined by 34b.
"../bin/wikify" Defined by 33g.
"../bin/wsd" Defined by 44c.
"../env/bin/langdetect.py" Defined by 20c.
"../env/bin/progenv" Defined by 5b, 8a.
"../env/bin/transplant" Defined by 71.
"../env/etc/opini.cfg" Defined by 54b.

```

"../nuweb/bin/w2pdf" Defined by 65b.

"Makefile" Defined by 59b.

"w2html" Defined by 68a.

C.2 Macro's

<abort when the language is not English or Dutch 21e> Referenced in 30c, 31b.
 <activate the install-to-project-repo utility 49d> Referenced in 49b.
 <activate the python environment 11d, 12a> Referenced in 10b, 15a.
 <all targets 60b> Referenced in 59c.
 <annotate dutch document 56a> Not referenced.
 <annotate english document 56b> Referenced in 55c.
 <begin conditional install 14b> Referenced in 8c, 15a, 16aj, 17aj, 18ahq, 19a.
 <check listener on host, port 27a> Referenced in 26b, 28be.
 <check this first 7b, 19e> Referenced in 15a.
 <check whether mercurial is present 19f> Referenced in 19e.
 <check/install the correct version of python 10c> Referenced in 10b.
 <clean up 8f, 9d, 22b, 61a> Referenced in 60a.
 <clone the heideltime wrapper 48d> Referenced in 48c.
 <compile nuweb 65c> Referenced in 65b.
 <compile the heideltime wrapper 49b> Referenced in 48c.
 <compile the nerc jar 42b> Referenced in 42a.
 <create a virtual environment for Python 11b> Referenced in 10b.
 <create javapython script 7d> Referenced in 15a.
 <default target 59c> Referenced in 59b.
 <directories to create 4abcd, 8bh, 9a, 11e, 64c> Referenced in 70a.
 <download svm models 44b> Referenced in 43f.
 <else conditional install 14c> Not referenced.
 <end conditional install 14d> Referenced in 8c, 15a, 16aj, 17aj, 18ahq, 19a.
 <expliciete make regels 60d, 61bcd, 63b, 65a, 67bd> Referenced in 59b.
 <filenames in nuweb compile script 65d> Referenced in 65b.
 <filenames in w2html 68c> Referenced in 68a.
 <function to check/start spotlight 26b> Referenced in 55c.
 <get jvntextpro-2.0.jar 49c> Referenced in 49b.
 <get spotlight language parameters 26g> Referenced in 26b.
 <get spotlight model ball 25d> Referenced in 25a.
 <get the mor time-out parameter 32a> Referenced in 31e.
 <get the nerc models 43a> Referenced in 42a.
 <get the snapshot 7c> Referenced in 15a.
 <impliciete make regels 63a, 67c> Referenced in 59b.
 <install ActivePython 11a> Referenced in 10c.
 <install Alpino 21f> Referenced in 15a.
 <install boost 24e> Referenced in 15a.
 <install coreference-base 41c> Referenced in 18a.
 <install CRFsuite 30a> Referenced in 15a.
 <install from github 6d> Referenced in 31d, 36c, 41c, 43f, 45c, 48d, 50d, 53a.
 <install from tarball 57c> Not referenced.
 <install Java 1.6 9f> Referenced in 15a.
 <install kafnafparserpy 13a> Referenced in 10b.
 <install maven 9b> Referenced in 15a.
 <install python packages 13f, 51e> Referenced in 10b.
 <install sematree 20a> Referenced in 15a.
 <install svm lib 44a> Referenced in 43f.
 <install SVMLight 29b> Referenced in 15a.
 <install the constituents parser 32e> Referenced in 16j.
 <install the dbpedia-ner module 53a> Referenced in 18q.
 <install the event-coreference module 52c> Referenced in 18h.
 <install the factuality module 40f> Referenced in 17j.

<install the FBK-causalrel module 40a> Referenced in 17j.
 <install the FBK-temprel module 39a> Referenced in 17j.
 <install the FBK-time module 37> Referenced in 17j.
 <install the heideltime module 48c> Referenced in 18h.
 <install the ims-wsd module 34d> Referenced in 17a.
 <install the lu2synset converter 44e> Referenced in 18h.
 <install the morphosyntactic parser 31d> Referenced in 16a.
 <install the NERC module 42a> Referenced in 16j.
 <install the nomevent module 53d> Referenced in 18q.
 <install the onto module 46d> Referenced in 18a.
 <install the opinion-miner 54ac> Referenced in 19a.
 <install the pos tagger 32b> Referenced in 16a.
 <install the post-SRL module 51f> Referenced in 18q.
 <install the Spotlight server 25a, 26a> Referenced in 15a.
 <install the srl module 50d> Referenced in 18h.
 <install the srl-dutch-nominals module 36c> Referenced in 17a.
 <install the srl-server module 35c> Referenced in 17a.
 <install the ticcutils utility 24a> Referenced in 15a, 24d.
 <install the timbl utility 24b> Referenced in 15a, 24d.
 <install the tokenizer 30b> Referenced in 16a.
 <install the topic analyser 31a> Referenced in 16a.
 <install the treetagger utility 22cde, 23abcd> Referenced in 15a.
 <install the UKB module 34a> Not referenced.
 <install the wikify module 33f> Referenced in 16j.
 <install the WSD module 43f> Referenced in 18a.
 <install the NED-reranker module 33c> Referenced in 16j.
 <install the NED module 45c> Referenced in 16j.
 <install VUA-pylib 29a> Referenced in 15a.
 <logmess 57b> Referenced in 6bd, 13a, 57c.
 <make scripts executable 19d, 21c, 70b> Referenced in 70c.
 <make targets 60a, 64ab, 67e, 70acd> Referenced in 59b.
 <move module 5c> Referenced in 6d, 13a, 57c.
 <parameters in Makefile 59a, 60c, 62ab, 64d, 67a, 69d> Referenced in 59b.
 <perform the task of w2html 68b> Referenced in 68a.
 <put Antske's material in the heideltime wrapper 49a> Referenced in 48c.
 <put spotlight jar in the Maven repository 46a> Referenced in 45c.
 <re-install modules after the transplantation 24d> Referenced in 71.
 <re-instate old module 6b> Referenced in 6d, 13a, 57c.
 <read the list of installed modules 14a> Referenced in 15a.
 <remove installed-variable 14e> Referenced in 9d.
 <remove old module 6a> Referenced in 6d, 13a, 57c.
 <remove outdated heideltime jars 50a> Referenced in 49d.
 <remove the copy of the aux file 66a> Referenced in 65c, 68b.
 <repair causalrel's run.sh.hadoop 40c> Not referenced.
 <repair FBK-*rel's run.sh.hadoop 39c> Referenced in 39a, 40a.
 <run tex4ht 69c> Referenced in 69a.
 <run the html processors 69b> Referenced in 69a.
 <run the html processors until the aux file remains unchanged 69a> Referenced in 68b.
 <run the processors until the aux file remains unchanged 66c> Referenced in 65c.
 <run the three processors 66b> Referenced in 66c.
 <select language-dependent features 42c> Referenced in 55c.
 <set alpinohome 22a> Referenced in 31e.
 <set paths after transplantation 12c> Referenced in 71.
 <set the language variable 21d> Referenced in 55c.
 <set up java 8cg> Referenced in 15a.
 <set up Java 1.6 10a> Referenced in 35a.
 <set up python 10b> Referenced in 15a.
 <set variables that point to the directory-structure 4e, 5a, 7a, 9c> Referenced in 5b, 15a, 71.

⟨start EHU SRL server if it isn't running 36a⟩ Referenced in 35d.
 ⟨start of module-script 20b⟩ Referenced in 21a, 30c, 31be, 32c, 33adg, 34b, 35adfh, 36d, 38a, 39d, 40d, 41ad, 43bd, 44c, 45a, 46b, 47a, 48a, 50be, 52ad, 53be, 55a.
 ⟨start the Spotlight server on localhost 28b⟩ Referenced in 26b, 27b.
 ⟨stop EHU SRL server 36b⟩ Referenced in 35f.
 ⟨test whether spotlightserver runs 28a⟩ Referenced in 27b.
 ⟨test whether virtualenv is present on the host 11c⟩ Referenced in 11b.
 ⟨try to obtain a running spotlightserver 27b⟩ Not referenced.
 ⟨unpack ticcutils or timbl 24c⟩ Referenced in 24ab.
 ⟨update pip 12b⟩ Referenced in 10b.
 ⟨variables of install-modules 57a⟩ Referenced in 15a.
 ⟨wait until the spotlight server is up 28e⟩ Referenced in 28b.

C.3 Variables

activate: 11d, 12c.
 all: 29b, 59c.
 ALPINO_HOME: 22a.
 auxfil: 65d, 66c, 68c, 69a.
 bibtex: 66b, 69bc.
 DIRS: 70a, 70c.
 fig2dev: 63a.
 FIGFILENAMES: 62b.
 FIGFILES: 62a, 62b.
 hg: 19f, 41c.
 indexfil: 65d, 66c, 68c.
 lxml: 13f.
 makeindex: 66b, 69bc.
 MKDIR: 69d, 70a.
 naflang: 21d, 21e, 26bg, 28b, 30c, 33g, 42c, 43d, 46b, 52d, 55ac.
 nufil: 65d, 66b, 68c, 69b.
 nuweb: 4e, 55c, 59a, 60d, 61ab, 64cd, 65bc, 66b, 67a, 68a.
 oldaux: 65d, 66ac, 68c, 69a.
 oldindexfil: 65d, 66c, 68c.
 PATH: 5a, 8g, 9c, 10a, 41a, 55a.
 pdf: 60bc, 64a, 64b.
 PDFT_NAMES: 62b, 64b.
 PDF_FIG_NAMES: 62b, 64b.
 PHONY: 59c, 63b.
 print: 10c, 12c, 14e, 20c, 31a, 48a, 61c, 64a.
 PST_NAMES: 62b.
 PS_FIG_NAMES: 62b.
 pythonok: 10c.
 PYTHONPATH: 12a.
 pyyaml: 13f.
 rdflib: 51e.
 SUFFIXES: 60c.
 texfil: 65d, 66b, 68c, 69b.
 trunk: 65d, 66b, 68c, 69bc.
 view: 64a.
 virtualenv: 11ab, 11c.