

Standardised Dutch NLP pipeline

Paul Huygen <paul.huygen@huygen.nl>

29th June 2015
08:34 h.

Abstract

This is a description and documentation of the installation of the current NLP modules on Lisa, so that they can be used in pipelines.

Contents

1	Introduction	2
1.1	List of the modules to be installed	2
1.2	File-structure of the pipeline	3
2	Java and Python environment	5
2.1	Java	5
2.2	Maven	6
2.3	Python	7
2.3.1	Python version	7
2.3.2	Virtual environment	8
2.3.3	KafNafParserPy	9
2.3.4	Python packages	10
3	Installation	10
3.1	Installing vs. updating	10
3.2	Installation from Github	11
3.3	Installation from the snapshot	11
3.4	The installation script	12
3.5	Check availability of resources	14
3.6	Install utilities and resources	14
3.6.1	Alpino	14
3.6.2	Treetagger	15
3.6.3	Timbl and Ticcutils	17
3.6.4	Spotlight	18
3.7	Install modules	19
3.7.1	Install tokenizer	19
3.7.2	Morphosyntactic parser	20
3.7.3	Nominal coreference-base	21
3.7.4	Named entity recognition (NERC)	21
3.7.5	Wordsense-disambiguation	23
3.7.6	Lexical-unit converter	25
3.7.7	NED	25
3.7.8	Ontotagger	26
3.7.9	Framenet SRL	27

3.7.10	Heideltime	28
3.7.11	Semantic Role labelling	29
3.7.12	Event coreference	30
4	Utilities	31
4.1	Test script	31
4.2	Logging	32
4.3	Misc	32
A	How to read and translate this document	33
A.1	Read this document	33
A.2	Process the document	34
A.3	Translate and run	34
A.4	Get Nuweb	35
A.5	Pre-processing	36
A.5.1	Process ‘dollar’ characters	36
A.5.2	Run the M4 pre-processor	37
A.6	Typeset this document	37
A.6.1	Figures	37
A.6.2	Bibliography	38
A.6.3	Create a printable/viewable document	39
A.6.4	Create HTML files	42
B	References	45
B.1	Literature	45
B.2	URL’s	46
C	Indexes	46
C.1	Filenames	46
C.2	Macro’s	46
C.3	Variables	47

1 Introduction

This document describes the current set-up of pipeline that annotates dutch 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.

Apart from describing the pipeline set-up, the document actually constructs the pipeline. Currently, the pipeline has been successfully implemented on a specific supercomputer (Lisa, Surfsara, Amsterdam³) and on computers running Ubuntu and Centos.

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.

1.1 List of the modules to be installed

Table 1 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.

1. <http://wordpress.let.vupr.nl>
2. <http://www.newsreader-project.eu>
3. <https://surfsara.nl/systems/lisa>

Module	Section	Source	Commit	Script
Tokenizer	3.7.1	snapshot	56f83ce4b61680346f15e5d4e6de6293764f7383	tok
morphosyntactic parser	3.7.2	Github	c6cabea2cc37ac3098c5927f5ec5b180ac31246f	mor
NERC	3.7.4	Github	9927fdb32d943f0aa9748a656958af99eeb1f5b7	nerc
WSD	3.7.5	Github	2babeb40a81b3720274a0521ccc2a27c5eff28c9	wsd
Onto-tagger	3.7.8	snapshot		onto
Heideltime	3.7.10	Github	057c93ccc857a427145b9e2ff72fd645172d34df	heideltime
SRL	3.7.11	Github	675d22d361289ede23df11dcdb17195f008c54bf	srl
NED	3.7.7	Github	d35d4df5cb71940bf642bb1a83e2b5b7584010df	ned
Nom. coref	3.7.3	Github	bfa5aec0fa498e57fe14dd4d2c51365dd09a0757	nomcoref
Ev. coref	3.7.12	snapshot		evcoref
Framenet SRL	3.7.9	snapshot		fsrl

Table 1: 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.

- 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 1 this has been indicated as SNAPSHOT.

The modules themselves use other utilities like dependency-taggers and POS taggers. These utilities are listed in table 2.

Module	Version	Section	Source
KafNafParserPy	Feb 1, 2015	2.3.3	Github
Alpino	20706	3.6.1	RUG
Ticcutils	0.7	3.6.3	ILK
Timbl	6.4.6	3.6.3	ILK
Treetagger	3.2	3.6.2	Uni. München

Table 2: 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.

1.2 File-structure of the pipeline

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

$\langle \text{directories to create} \rangle \equiv$
`../modules` ◊

Fragment defined by 3, 4ab, 5d, 6ce, 9c, 40a.

Fragment referenced in 45b.

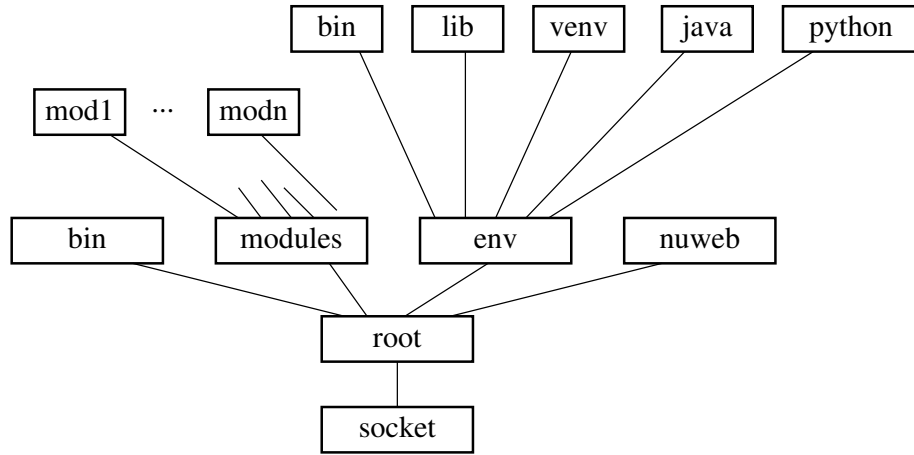


Figure 1: *Directory-structure of the pipeline (see text).*

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

Fragment defined by 3, 4ab, 5d, 6ce, 9c, 40a.
 Fragment referenced in 45b.

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

Fragment defined by 3, 4ab, 5d, 6ce, 9c, 40a.
 Fragment referenced in 45b.

The following macro defines variable `piperoor` 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 `piperoor` 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 4c > ≡
  if
    [ "$piperoor" == "" ]
  then
    export piperoor=/mnt/sdb1/pipelines/testnlpp/nlpp
  fi
  export pipesocket=${piperoor%%/nlpp}
  export envdir=$piperoor/env
  export envbindir=$envdir/bin
  export envlibdir=$envdir/lib
  export modulesdir=$piperoor/modules
  export pipebin=$piperoor/bin
  ◇

```

Fragment defined by 4c, 5ab.

Fragment referenced in 20b, 21ac, 23cd, 24e, 26b, 27, 28a, 29b, 30a, 31c.

Add the environment bin directory to PATH:

```

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

```

Fragment defined by 4c, 5ab.

Fragment referenced in 20b, 21ac, 23cd, 24e, 26b, 27, 28a, 29b, 30a, 31c.

Defines: PATH 6b, 7a.

While setting variables, *source* a scripts that sets variables for directories of which we do not yet know where they are, e.g. paths to Python and Java that we may have to set up dynamically.

```

⟨ set variables that point to the directory-structure 5b ⟩ ≡
    source $envbindir/progenv
    ◇

```

Fragment defined by 4c, 5ab.

Fragment referenced in 20b, 21ac, 23cd, 24e, 26b, 27, 28a, 29b, 30a, 31c.

2 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 file sets up the programming environment in scripts.

```

"../env/bin/progenv" 5c≡
    ⟨ set up java environment in scripts 6b, ... ⟩
    ⟨ activate the python environment 9b, ... ⟩
    ◇

```

2.1 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 5d ⟩ ≡
    ../env/java ◇

```

Fragment defined by 3, 4ab, 5d, 6ce, 9c, 40a.

Fragment referenced in 45b.

```

⟨ set up java 5e ⟩ ≡
    ⟨ get or have (5f server-jre-7u72-linux-x64.tar.gz ) 12b ⟩
    cd $envdir/java
    tar -xzf $pipesocket/server-jre-7u72-linux-x64.tar.gz
    ◇

```

Fragment referenced in 13a.

Remove the java-ball when cleaning up:

```

⟨ clean up 6a ⟩ ≡
    rm -rf $pipesocket/server-jre-7u72-linux-x64.tar.gz
    ◇

```

Fragment defined by 6a, 7b, 36b.

Fragment referenced in 35b.

```

⟨ set up java environment in scripts 6b ⟩ ≡
    export JAVA_HOME=$envdir/java/jdk1.7.0_72
    export PATH=$JAVA_HOME/bin:$PATH
    ◇

```

Fragment defined by 6bd.

Fragment referenced in 5c, 13a.

Defines: JAVA_HOME Never used.

Uses: PATH 5a.

Put jars in the jar subdirectory of the java directory:

```

⟨ directories to create 6c ⟩ ≡
    ../env/java/jars ◇

```

Fragment defined by 3, 4ab, 5d, 6ce, 9c, 40a.

Fragment referenced in 45b.

```

⟨ set up java environment in scripts 6d ⟩ ≡
    export JARDIR=$envdir/java/jars
    ◇

```

Fragment defined by 6bd.

Fragment referenced in 5c, 13a.

2.2 Maven

```

⟨ directories to create 6e ⟩ ≡
    ../env/apache-maven-3.0.5 ◇

```

Fragment defined by 3, 4ab, 5d, 6ce, 9c, 40a.

Fragment referenced in 45b.

```

⟨ install maven 6f ⟩ ≡
    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 defined by 6f, 7a.

Fragment referenced in 13a.

```

< install maven 7a > ≡
    export MAVEN_HOME=/mnt/sdb1/pipelines/testnlpp/nlpp/env/apache-maven-3.0.5
    export PATH=${MAVEN_HOME}/bin:${PATH}
    ◇

```

Fragment defined by 6f, 7a.
 Fragment referenced in 13a.
 Uses: PATH 5a.

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

```

< clean up 7b > ≡
    rm -rf ../env/apache-maven-3.0.5
    ◇

```

Fragment defined by 6a, 7b, 36b.
 Fragment referenced in 35b.

2.3 Python

Set up the environment for Python. I could not find an easy way to set up Python from scratch, so we have to rely on Python 2.7 being available on the host. However, we can make a virtual environment, so that we are not dependent on the existence of libraries in the right version on the host.

In the virtual environment we will install KafNafParserPy and other Python packages that are needed.

```

< set up python 7c > ≡
    < check/install the correct version of python 8a >
    < create a virtual environment for Python 8e >
    < activate the python environment 9b, ... >
    < install kafnafparserpy 10a >
    < install python packages 10f >
    ◇

```

Fragment referenced in 13a.

2.3.1 Python version

The pipeline relies on Python version 2.7 being available. If possible, the user should provide this version and make sure that the “python” command invokes version 2.7.something of python. However, in some cases (notably in the case of a Centos 6.3 server) this is difficult to achieve. In that case we can use a binary python supplied by ActivePython (<http://www.activestate.com/activepython>). Download in that case the tarball ActivePython-2.7.8.10-linux-x86_64.tar.gz from the ActivePython site and put it in the nlpp directory. The following macro checks whether the python command invokes a correct version of python and, if this is not the case and the ActivePython tarball is present, install ActivePython.

```

< check/install the correct version of python 8a > ≡
pythonok='python --
version 2>&1 | gawk '{if(match($2, "2.7")) print "yes" ; else print "no" }'
if
[ "$pythonok" == "no" ]
then
  < install ActivePython 8b, ... >
fi
◇

```

Fragment referenced in 7c.

Defines: pythonok Never used.

Uses: print 39b.

Check whether we have the ActivePython tarball and quit if this is not the case.

```

< install ActivePython 8b > ≡
  < get or have (8c ActivePython-2.7.8.10-linux-x86_64.tar.gz ) 12b >
◇

```

Fragment defined by 8bd.

Fragment referenced in 8a.

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 8d > ≡
pytinsdir='mktemp -d -t activepyt.XXXXXX'
cd $pytinsdir
tar -xzf $pipesocket/ActivePython-2.7.8.10-linux-x86_64.tar.gz
acdir='ls -1'
cd $acdir
./install.sh -I $envdir
cd $piperoot
rm -rf $pytinsdir
pip install -U pip virtualenv setuptools
◇

```

Fragment defined by 8bd.

Fragment referenced in 8a.

Uses: virtualenv 9a.

2.3.2 Virtual environment

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

```

< create a virtual environment for Python 8e > ≡
  < test whether virtualenv is present on the host 9a >
  cd $envdir
  virtualenv venv
◇

```

Fragment referenced in 7c.

Uses: virtualenv 9a.


```

< test whether virtualenv is present on the host 9a > ≡
    which virtualenv
    if
        [ $? -ne 0 ]
    then
        echo Please install virtualenv
        exit 1
    fi
    ◇

```

Fragment referenced in [8e](#).
 Defines: `virtualenv` [8de](#).

```

< activate the python environment 9b > ≡
    source $envdir/venv/bin/activate
    ◇

```

Fragment defined by [9bd](#).
 Fragment referenced in [5c](#), [7c](#).
 Defines: `activate` Never used.

Subdirectory `$envdir/python` will contain general Python packages like `KafNafParserPy`.

```

< directories to create 9c > ≡
    ../env/python ◇

```

Fragment defined by [3](#), [4ab](#), [5d](#), [6ce](#), [9c](#), [40a](#).
 Fragment referenced in [45b](#).

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

```

< activate the python environment 9d > ≡
    export PYTHONPATH=$envdir/python:$PYTHONPATH
    ◇

```

Fragment defined by [9bd](#).
 Fragment referenced in [5c](#), [7c](#).
 Defines: `PYTHONPATH` Never used.

2.3.3 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 10a> ≡
  cd $envdir/python
  DIRN=KafNafParserPy
  <move module (10b $DIRN) 10g>
  git clone https://github.com/cltl/KafNafParserPy.git
  if
    [ $? -gt 0 ]
  then
    <logmess (10c Cannot install current $DIRN version) 32c>
    <re-instate old module (10d $DIRN) 11b>
  else
    <remove old module (10e $DIRN) 11a>
  fi
  ◇

```

Fragment referenced in 7c.

2.3.4 Python packages

Install python packages:

lxml:

pyyaml: for coreference-graph

```

<install python packages 10f> ≡
  pip install lxml
  pip install pyyaml
  ◇

```

Fragment referenced in 7c.

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

3 Installation

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

3.1 Installing vs. updating

When the install-script installs something that has already been installed, it moves the installed module to a temporary location and then tries to install the module from its source. If that is successful it removes the former version of the module, otherwise it moves the old version back.

The following macro's can be used to move or remove modules, provided they are called when the modules directory is the default directory.

```

<move module 10g> ≡
  if
    [ -e @1 ]
  then
    mv @1 old.@1
  fi
  ◇

```

Fragment referenced in 10a, 11d, 15a, 33a.

```

<remove old module 11a> ≡
    rm -rf old.@1
    ◇

```

Fragment referenced in 10a, 11d, 15a, 33a.

```

<re-instate old module 11b> ≡
    mv old.@1 @1
    MESS="Replaced previous version of @1"
    <logmess (11c $MESS) 32c>
    ◇

```

Fragment referenced in 10a, 11d, 15a, 33a.

3.2 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 11d> ≡
    cd /mnt/sdb1/pipelines/testnlpp/nlpp/modules
    <move module (11e $DIRN) 10g>
    git clone $GITU
    if
        [ $? -gt 0 ]
    then
        <logmess (11f Cannot install current $MODNAM version) 32c>
        <re-instate old module (11g $DIRN) 11b>
    else
        <remove old module (11h $DIRN) 11a>
        cd /mnt/sdb1/pipelines/testnlpp/nlpp/modules/$DIRN
        git checkout $GITC
    fi
    ◇

```

Fragment referenced in 20c, 21b, 24a, 25d, 28b, 29c.

3.3 Installation from the snapshot

For some modules a public repository is not available or not known. They must be installed from a non-public repository. A key to connect to the repository can be requested from the author.

```

< have an SSH key or die 12a > ≡
  if
    [ ! -e /mnt/sdb1/pipelines/testnlpp/nrkey ]
  then
    echo "No key to connect to snapshot!"
    exit 1
  fi
◇

```

Fragment referenced in 12b.

The following macro downloads a resource if it is not already present in the “socket” directory. It turns out that sometimes there is a time-out for unknown reasons. In that case we will try it multiple times.

```

< get or have 12b > ≡
  counter=0
  while
    [ ! -e $pipesocket/@1 ]
  do
    < have an SSH key or die 12a >
    cd $pipesocket
    scp -
  i "/mnt/sdb1/pipelines/testnlpp/nrkey" newsreader@kyoto.let.vu.nl:nlpp_resources/@1 .
    if
      [ $? -gt 0 ]
    then
      counter=$((counter+1))
      if
        [ $counter -gt 3 ]
      then
        echo "Cannot contact snapshot server"
        exit 1
      fi
    fi
  done
◇

```

Fragment referenced in 5e, 8b, 19a, 23a, 24c, 25a, 26c, 31a.

3.4 The installation script

The installation is performed by script `install-modules`

```

"../bin/install-modules" 13a≡
    #!/bin/bash
    echo Set up environment
    < set local bin directory ? >
    < variables of install-modules 32b >
    < check this first 14b >
    echo ... Java
    < set up java 5e >
    < set up java environment in scripts 6b, ... >
    < install maven 6f, ... >
    echo ... Python
    < set up python 7c >
    echo ... Alpino
    < install Alpino 15a >
    echo ... Spotlight
    < install the Spotlight server 19a, ... >
    echo ... Treetagger
    < install the treetagger utility 15g, ... >
    echo ... Ticcutils and Timbl
    < install the ticcutils utility 17b >
    < install the timbl utility 17c >
    ◇

```

File defined by 13ab.

```

"../bin/install-modules" 13b≡
    echo Install modules
    echo ... Tokenizer
    < install the tokenizer 20a >
    echo ... Morphosyntactic parser
    < install the morphosyntactic parser 20c >
    echo ... NERC
    < install the NERC module 21d >
    echo ... Coreference base
    < install coreference-base 21b >
    echo ... WSD
    < install the WSD module 24a >
    echo ... Ontotagger
    < install the onto module 26c >
    echo ... Heideltime
    < install the heideltime module 28b >
    echo ... SRL
    < install the srl module 29c >
    echo ... NED
    < install the NED module 25d >
    echo ... Event-coreference
    < install the event-coreference module 31a >
    echo ... lu2synset
    < install the lu2synset converter 25a >
    echo Final
    ◇

```

File defined by 13ab.

```

< make scripts executable 14a > ≡
    chmod 775 ../bin/install-modules
    ◇

```

Fragment defined by 14a, 45c.

Fragment referenced in 45d.

3.5 Check availability of resources

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

```

< check this first 14b > ≡
    < check whether mercurial is present 14c >
    ◇

```

Fragment referenced in 13a.

```

< check whether mercurial is present 14c > ≡
    which hg
    if
        [ $? -ne 0 ]
    then
        echo Please install Mercurial.
        exit 1
    fi
    ◇

```

Fragment referenced in 14b.

Defines: hg 21b.

3.6 Install utilities and resources

3.6.1 Alpino

Install Alpino from the website of Gertjan van Noort.

Module

```

< install Alpino 15a > ≡
  SUCCES=0
  cd /mnt/sdb1/pipelines/testnlpp/nlpp/modules
  < move module (15b Alpino ) 10g >
  wget http://www.let.rug.nl/vannoord/alp/Alpino/binary/versions/Old/Alpino-x86_64-
  linux-glibc2.5-20706-sicstus.tar.gz
  SUCCES=$?
  if
    [ $SUCCES -eq 0 ]
  then
    tar -xzf Alpino-x86_64-linux-glibc2.5-20706-sicstus.tar.gz
    SUCCES=$?
    rm -rf Alpino-x86_64-linux-glibc2.5-20706-sicstus.tar.gz
  fi
  if
    [ $SUCCES -eq 0 ]
  then
    < logmess (15c Installed Alpino ) 32c >
    < remove old module (15d Alpino ) 11a >
  else
    < re-instate old module (15e Alpino ) 11b >
  fi
  ◇

```

Fragment referenced in 13a.

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 15f > ≡
  export ALPINO_HOME=/mnt/sdb1/pipelines/testnlpp/nlpp/modules/Alpino
  ◇

```

Fragment referenced in 21a.

Defines: ALPINO_HOME Never used.

3.6.2 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 15g > ≡
  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 15g, 16abcde, 17a.

Fragment referenced in 13a.

The source tarball, scripts and the installation-script:

```

< install the treetagger utility 16a > ≡
    TREETAGSRC=tree-tagger-linux-3.2.tar.gz
    TREETAGSCRIPTS=tagger-scripts.tar.gz
    TREETAG_INSTALLSCRIPT=install-tagger.sh
    ◇

```

Fragment defined by 15g, 16abcde, 17a.

Fragment referenced in 13a.

Parametersets:

```

< install the treetagger utility 16b > ≡
    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 15g, 16abcde, 17a.

Fragment referenced in 13a.

Download everything in the target directory:

```

< install the treetagger utility 16c > ≡
    mkdir -p /mnt/sdb1/pipelines/testnlpp/nlpp/modules/$TREETAGDIR
    cd /mnt/sdb1/pipelines/testnlpp/nlpp/modules/$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 15g, 16abcde, 17a.

Fragment referenced in 13a.

Run the install-script:

```

< install the treetagger utility 16d > ≡
    chmod 775 $TREETAG_INSTALLSCRIPT
    ./$TREETAG_INSTALLSCRIPT
    ◇

```

Fragment defined by 15g, 16abcde, 17a.

Fragment referenced in 13a.

Make the treetagger utilities available for everybody.

```

< install the treetagger utility 16e > ≡
    chmod -R o+rx /mnt/sdb1/pipelines/testnlpp/nlpp/modules/$TREETAGDIR/bin
    chmod -R o+rx /mnt/sdb1/pipelines/testnlpp/nlpp/modules/$TREETAGDIR/cmd
    chmod -R o+r /mnt/sdb1/pipelines/testnlpp/nlpp/modules/$TREETAGDIR/doc
    chmod -R o+rx /mnt/sdb1/pipelines/testnlpp/nlpp/modules/$TREETAGDIR/lib
    ◇

```

Fragment defined by 15g, 16abcde, 17a.

Fragment referenced in 13a.

Remove the tarballs:


```

< install the treetagger utility 17a > ≡
    rm $TREETAGSRC
    rm $TREETAGSCRIPTS
    rm $TREETAG_INSTALLSCRIPT
    rm $DUTCHPARS_UTF_GZ
    rm $DUTCH_TAGSET
    rm $DUTCHPARS_2_GZ
    ◇

```

Fragment defined by 15g, 16abcde, 17a.

Fragment referenced in 13a.

3.6.3 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 17b > ≡
    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 18a >
    ◇

```

Fragment referenced in 13a.

```

< install the timbl utility 17c > ≡
    URL=http://software.ticc.uvt.nl/timbl-6.4.6.tar.gz
    TARB=timbl-6.4.6.tar.gz
    DIR=timbl-6.4.6
    < unpack ticcutils or timbl 18a >
    ◇

```

Fragment referenced in 13a.

```

⟨unpack ticcutils or timbl 18a⟩ ≡
  SUCCES=0
  ticbeldir='mktemp -t -d tickbel.XXXXXX'
  cd $ticbeldir
  wget $URL
  SUCCES=$?
  if
    [ $SUCCES -eq 0 ]
  then
    tar -xzf $TARB
    SUCCES=$?
    rm -rf $TARB
  fi
  if
    [ $SUCCES -eq 0 ]
  then
    cd $DIR
    ./configure --prefix=$envdir
    make
    make install
  fi
  cd /mnt/sdb1/pipelines/testnlpp/nlpp
  rm -rf $ticbeldir
  if
    [ $SUCCES -eq 0 ]
  then
    ⟨logmess (18b Installed $DIR ) 32c⟩
  else
    ⟨logmess (18c NOT installed $DIR ) 32c⟩
  fi
  ◇

```

Fragment referenced in 17bc.

3.6.4 Spotlight

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://local
```

Dutch server:

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

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

So, let's do that.

```
< install the Spotlight server 19a > ≡
  < get or have (19b spotlightnl.tgz ) 12b >
  cd $envdir
  tar -xzf $pipesocket/spotlightnl.tgz
  cd /mnt/sdb1/pipelines/testnlpp/nlpp/env/spotlight
  wget http://spotlight.sztaki.hu/downloads/nl.tar.gz
  tar -xzf nl.tar.gz
  rm nl.tar.gz
  ◇
```

Fragment defined by 19ac.

Fragment referenced in 13a.

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

```
< install the Spotlight server 19c > ≡
  cd /mnt/sdb1/pipelines/testnlpp/nlpp/env/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 19ac.

Fragment referenced in 13a.

```
< start the Spotlight server 19d > ≡
  cd /mnt/sdb1/pipelines/testnlpp/nlpp/env/spotlight
  java -jar -Xmx8g dbpedia-spotlight-0.7-jar-with-dependencies-
  candidates.jar nl http://localhost:2060/rest &
  ◇
```

Fragment referenced in 19e.

```
< check/start the Spotlight server 19e > ≡
  spottasks='netstat -an | grep :2060 | wc -l'
  if
  [ $spottasks -eq 0 ]
  then
    < start the Spotlight server 19d >
    sleep 60
  fi
  ◇
```

Fragment referenced in 26b.

3.7 Install modules

3.7.1 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 20a> ≡
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 /mnt/sdb1/pipelines/testnlpp/nlpp/env/java/jars
cd /mnt/sdb1/pipelines/testnlpp/nlpp
rm -rf $tempdir
◇
```

Fragment referenced in 13b.

Script The script runs the tokenizerscript.

```
"../bin/tok" 20b≡
#!/bin/bash
<set variables that point to the directory-structure 4c, ... >
JARFILE=/mnt/sdb1/pipelines/testnlpp/nlpp/env/java/jars/ixa-pipe-tok-1.8.0.jar
java -Xmx1000m -jar $JARFILE tok -l nl --inputkaf
◇
```

3.7.2 Morphosyntactic parser

Module

```
<install the morphosyntactic parser 20c> ≡
MODNAM=morphosynparser
DIRN=morphosyntactic_parser_nl
GITU=https://github.com/cltl/morphosyntactic_parser_nl.git
GITC=c6cabea2cc37ac3098c5927f5ec5b180ac31246f
<install from github 11d>
cd /mnt/sdb1/pipelines/testnlpp/nlpp/modules/morphosyntactic_parser_nl
git checkout c6cabea2cc37ac3098c5927f5ec5b180ac31246f
◇
```

Fragment referenced in 13b.

Script

```

"../bin/mor" 21a≡
    #!/bin/bash
    < set variables that point to the directory-structure 4c, ... >
    ROOT=/mnt/sdb1/pipelines/testnlpp/nlpp
    MODDIR=/mnt/sdb1/pipelines/testnlpp/nlpp/modules/morphosyntactic_parser_nl
    < set alpinohome 15f >
    cat | python $MODDIR/core/morph_syn_parser.py
    ◇

```

3.7.3 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>.

Module

```

< install coreference-base 21b > ≡
    MODNAM=coreference-base
    DIRN=coreference-base
    GITU=https://github.com/opener-project/coreference-base.git
    GIRC=bfa5aec0fa498e57fe14dd4d2c51365dd09a0757
    < install from github 11d >
    pip install --upgrade hg+https://bitbucket.org/Josu/pykaf#egg=pykaf
    pip install --upgrade networkx
    ◇

```

Fragment referenced in 13b.

Uses: hg 14c.

Script

```

"../bin/coreference-base" 21c≡
    #!/bin/bash
    < set variables that point to the directory-structure 4c, ... >
    cd $modulesdir/coreference-base/core
    cat | python -m corefgraph.process.file --language nl --singleton --sieves NO
    ◇

```

3.7.4 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 21d > ≡
    < compile the nerc jar 22 >
    < get the nerc models 23a >
    ◇

```

Fragment referenced in 13b.

The nerc module is a Java program that is contained in a jar. Pul 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 22 > ≡
TEMPDIR==`mktemp -d -t nerc.XXXXXX`
cd $TEMPDIR
git clone https://github.com/ixa-ehu/ixa-pipe-nerc
cd ixa-pipe-nerc/
git checkout 9927fdb32d943f0aa9748a656958af99eeb1f5b7
mvn clean package
mv target/ixa-pipe-nerc-1.3.6.jar /mnt/sdb1/pipelines/testnlpp/nlpp/env/java/jars/
cd /mnt/sdb1/pipelines/testnlpp/nlpp/nuweb
rm -rf $TEMPDIR
◇
```

Fragment referenced in 21d.

Uses: nuweb 41c.

The current version of the pipeline uses the following models, that have been made available by Rodrigo Agerri on march 2, 2015. Rodrigo wrote:

I have recently trained new models for Dutch using both the CoNLL 2002 and the Sonar corpora. These models are better than the one currently being used in the Dutch Newsreader pipeline. They are not yet in the resources of the ixa pipes (no public yet) but in the meantime they might be useful if you plan to do some processing in Dutch.

For CoNLL 2002, the new model obtains 83.46 F1, being the previously best published result 77.05 on that dataset. The Sonar model is trained on the full corpus, and evaluated using random 10 fold cross validation. The only previous result I know of obtains 80.71 F1 wrt to our model which obtains 87.84. However, because it is not evaluated on a separate test partition I do not take these results too seriously.

You will need to update the ixa-pipe-nerc module. The CoNLL 2002 model runs as before but to use the Sonar model you need to add the extra parameter `--clearFeatures` yes, like this:

```
Sonar model: cat file.pos.naf | java -jar ixa-pipe-nerc-1.3.6.jar tag
-m $nermodel --clearFeatures yes
CoNLL model: cat file.pos.naf | java -jar ixa-pipe-nerc-1.3.6.jar tag
-m $nermodel
```

<http://www.lt3.ugent.be/en/publications/fine-grained-dutch-named-entity-recognition/>

[..]

In any case, here are the models.

<http://ixa2.si.ehu.es/ragerri/dutch-nerc-models.tar.gz>

The tarball `dutch-nerc-models.tar.gz` contains the models `nl-clusters-conll02.bin` and `nl-clusters-sonar.bin`. Both models have been placed in subdirectory `/EHU-nerc/nerc-resources/nl` of the snapshot.

```

⟨ get the nerc models 23a ⟩ ≡
  ⟨ get or have (23b EHU-nerc.tgz ) 12b ⟩
  cd /mnt/sdb1/pipelines/testnlpp/nlpp/modules
  tar -xzf $pipesocket/EHU-nerc.tgz
  cp -r /mnt/sdb1/pipelines/testnlpp/nlpp/snapshots/EHU-nerc/nerc-
resources /mnt/sdb1/pipelines/testnlpp/nlpp/modules/EHU-nerc/
  chmod -R 775 /mnt/sdb1/pipelines/testnlpp/nlpp/modules/EHU-nerc
  ◇

```

Fragment referenced in 21d.

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

```

"../bin/nerc_conll02" 23c≡
  #!/bin/bash
  ⟨ set variables that point to the directory-structure 4c, ... ⟩
  MODDIR=$modulesdir/EHU-nerc
  JAR=$JARDIR/ixa-pipe-nerc-1.3.6.jar
  MODEL=nl-clusters-conll02.bin
  cat | java -Xmx1000m -jar $JAR tag -m $MODDIR/nerc-resources/nl/$MODEL
  #cat| java          -jar ixa-pipe-nerc-1.3.6.jar tag -m $nermodel
  ◇

"../bin/nerc_sonar" 23d≡
  #!/bin/bash
  ⟨ set variables that point to the directory-structure 4c, ... ⟩
  MODDIR=$modulesdir/EHU-nerc
  JAR=$JARDIR/ixa-pipe-nerc-1.3.6.jar
  MODEL=nl-clusters-sonar.bin
  cat | java -Xmx1000m -jar $JAR tag -m $MODDIR/nerc-resources/nl/$MODEL --
clearFeatures yes
  #cat| java          -jar ixa-pipe-nerc-1.3.6.jar tag -m $nermodel --
clearFeatures yes
  ◇

```

3.7.5 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 24a > ≡
MODNAM=wsd
DIRN=svm_wsd
GITU=https://github.com/cltl/svm_wsd.git
GITC=2babeb40a81b3720274a0521ccc2a27c5eff28c9
< install from github 11d >
cd /mnt/sdb1/pipelines/testnlpp/nlpp/modules/svm_wsd
< install svm lib 24b >
< download svm models 24c >

```

◇

Fragment referenced in 13b.

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

```

< install svm lib 24b > ≡
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 24a.

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

```

< download svm models 24c > ≡
< get or have (24d svm_wsd.tgz ) 12b >
cd /mnt/sdb1/pipelines/testnlpp/nlpp/modules
tar -xzf $pipesocket/svm_wsd.tgz

```

◇

Fragment referenced in 24a.

Script

```

"../bin/wsd" 24e≡
#!/bin/bash
# WSD -- wrapper for word-sense disambiguation
# 8 Jan 2014 Ruben Izquierdo
# 16 sep 2014 Paul Huygen
< set variables that point to the directory-structure 4c, ... >
WSDDIR=$modulesdir/svm_wsd
WSDSCRIPT=dsc_wsd_tagger.py
cat | python $WSDDIR/$WSDSCRIPT --naf

```

◇

3.7.6 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 25a> ≡
  <get or have (25b lu2synset.tgz ) 12b>
  cd /mnt/sdb1/pipelines/testnlpp/nlpp/modules
  tar -xzf $pipesocket/lu2synset.tgz
  ◇
```

Fragment referenced in 13b.

Script

```
"../bin/lu2synset" 25c≡
  #!/bin/bash
  ROOT=/mnt/sdb1/pipelines/testnlpp/nlpp
  JAVAILIBDIR=/mnt/sdb1/pipelines/testnlpp/nlpp/modules/lexicalunitconvertor/lib
  RESOURCESDIR=/mnt/sdb1/pipelines/testnlpp/nlpp/modules/lexicalunitconvertor/resources
  JARFILE=WordnetTools-1.0-jar-with-dependencies.jar
  java -Xmx812m -
  cp $JAVAILIBDIR/$JARFILE vu.wntools.util.NafLexicalUnitToSynsetReferences \
    --wn-lmf "$RESOURCESDIR/cornetto2.1.lmf.xml" --format naf
  ◇
```

3.7.7 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 NER 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 25d> ≡
  <put spotlight jar in the Maven repository 26a>
  MODNAM=ned
  DIRN=ixa-pipe-ned
  GITU=https://github.com/ixa-ehu/ixa-pipe-ned.git
  GITC=d35d4df5cb71940bf642bb1a83e2b5b7584010df
  <install from github 11d>
  cd /mnt/sdb1/pipelines/testnlpp/nlpp/modules/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 /mnt/sdb1/pipelines/testnlpp/nlpp/env/java/jars/
  ◇
```

Fragment referenced in 13b.

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 26a > ≡
    echo Put Spotlight jar in the Maven repository.
    tempdir='mktemp -d -t simplespot.XXXXXX'
    cd $tempdir
    wget http://spotlight.sztaki.hu/downloads/dbpedia-spotlight-0.7.jar
    wget http://spotlight.sztaki.hu/downloads/nl.tar.gz
    tar -xzf nl.tar.gz
    MVN_SPOTLIGHT_OPTIONS="-Dfile=dbpedia-spotlight-0.7.jar"
    MVN_SPOTLIGHT_OPTIONS="$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 25d.

Script

```

"../bin/ned" 26b ≡
    #!/bin/bash
    < set variables that point to the directory-structure 4c, ... >
    ROOT=/mnt/sdb1/pipelines/testnlpp/nlpp
    JARDIR=/mnt/sdb1/pipelines/testnlpp/nlpp/env/java/jars
    < check/start the Spotlight server 19e >
    cat | java -Xmx1000m -jar $JARDIR/ixa-pipe-ned-1.1.1.jar -p 2060 -e candidates -
    i /mnt/sdb1/pipelines/testnlpp/nlpp/env/spotlight/wikipedia-db -n nlEn
    ◇

```

3.7.8 Ontotagger

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

Module

```

< install the onto module 26c > ≡
    < get or have (26d vua-ontotagger-v1.0.tar.gz ) 12b >
    cd /mnt/sdb1/pipelines/testnlpp/nlpp/modules
    tar -xzf $pipesocket/vua-ontotagger-v1.0.tar.gz
    rm $pipesocket/vua-ontotagger-v1.0.tar.gz
    chmod -R o+r /mnt/sdb1/pipelines/testnlpp/nlpp/modules
    ◇

```

Fragment referenced in 13b.

Script

```

"../bin/onto" 27≡
    #!/bin/bash
    < set variables that point to the directory-structure 4c, ... >
    ROOT=/mnt/sdb1/pipelines/testnlpp/nlpp
    ONTODIR=$modulesdir/vua-ontotagger-v1.0
    JARDIR=$ONTODIR/lib
    RESOURCESDIR=$ONTODIR/resources
    PREDICATEMATRIX="$RESOURCESDIR/PredicateMatrix_nl_lu_withESO.v0.2.role.txt"
    GRAMMATICALWORDS="$RESOURCESDIR/grammaticals/Grammatical-words.nl"
    TMPFIL='mktemp -t stap6.XXXXXX'
    cat >$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="$JAVA_ARGS --version 1.1"
    JAVA_ARGS="$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

    ◇

```

3.7.9 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" 28a≡
    #!/bin/bash
    < set variables that point to the directory-structure 4c, ... >
    ONTODIR=$modulesdir/vua-ontotagger-v1.0
    JARDIR=$ONTODIR/lib
    RESOURCEDIR=$ONTODIR/resources
    PREDICATEMATRIX="$RESOURCEDIR/PredicateMatrix_nl_lu_withESO.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

```

◇

Uses: `print` 39b.

3.7.10 Heideltime

Module

```

< install the heideltime module 28b > ≡
    MODNAM=heideltime
    DIRN=NAF-HeidelTime
    GITU=https://github.com/cltl/NAF-HeidelTime.git
    GITH=057c93ccc857a427145b9e2ff72fd645172d34df
    < install from github 11d >
    < adapt heideltime's config.props 29a >

```

◇

Fragment referenced in 13b.

```

< adapt heideltime's config.props 29a > ≡
CONFIL=/mnt/sdb1/pipelines/testnlpp/nlpp/modules/NAF-HeidelTime/config.props
tempfil='mktemp -t heideltmp.XXXXXX'
mv $CONFIL $tempfil
MODDIR=/mnt/sdb1/pipelines/testnlpp/nlpp/modules
TREETAGDIR=treetagger
AWKCOMMAND='^treeTaggerHome/ {$0="treeTagger-
Home = /mnt/sdb1/pipelines/testnlpp/nlpp/modules/treetagger"}; {print}'
gawk "$AWKCOMMAND" $tempfil >$CONFIL
rm -rf $tempfil
◇

```

Fragment referenced in 28b.

Uses: `print` 39b.

Script

```

"../bin/heideltime" 29b ≡
#!/bin/bash
< set variables that point to the directory-structure 4c, ... >
HEIDELDIR=/mnt/sdb1/pipelines/testnlpp/nlpp/modules/NAF-HeidelTime
TEMPDIR='mktemp -t -d heideltmp.XXXXXX'
cd $HEIDELDIR
iconv -t utf-
8//IGNORE | python $HEIDELDIR/HeidelTime_NafKaf.py $HEIDELDIR/heideltime-
standalone/ $TEMPDIR
rm -rf $TEMPDIR
◇

```

3.7.11 Semantic Role labelling

Module

```

< install the srl module 29c > ≡
MODNAM=srl
DIRN=vua-srl-nl
GITU=https://github.com/newsreader/vua-srl-nl.git
GITC=675d22d361289ede23df11dcdb17195f008c54bf
< install from github 11d >
◇

```

Fragment referenced in 13b.

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" 30a≡
    #!/bin/bash
    < set variables that point to the directory-structure 4c, ... >
    source $envbindir/progenv
    ROOT=$piperoot
    SRLDIR=$modulesdir/vua-srl-nl
    TEMPDIR='mktemp -d -t SRLTMP.XXXXXX'
    cd $SRLDIR
    INPUTFILE=$TEMPDIR/inputfile
    FEATUREVECTOR=$TEMPDIR/csvfile
    TIMBLOUTPUTFILE=$TEMPDIR/timblpredictions
    ◇

```

File defined by 30abcde.

Create a feature-vector.

```

"../bin/srl" 30b≡
    cat | tee $INPUTFILE | python nafAlpinoToSRLFeatures.py > $FEATUREVECTOR
    ◇

```

File defined by 30abcde.

Run the trained model on the feature-vector.

```

"../bin/srl" 30c≡
    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 30abcde.

Insert the SRL values into the NAF file.

```

"../bin/srl" 30d≡
    python timblToAlpinoNAF.py $INPUTFILE $TIMBLOUTPUTFILE
    ◇

```

File defined by 30abcde.

Clean up.

```

"../bin/srl" 30e≡
    rm -rf $TEMPDIR
    ◇

```

File defined by 30abcde.

3.7.12 Event coreference

Module Install the module from the snapshot.

```

< install the event-coreference module 31a > ≡
  < get or have (31b vua-eventcoreference_v2.tar.gz ) 12b >
  cd /mnt/sdb1/pipelines/testnlpp/nlpp/modules
  tar -xzf $pipesocket/vua-eventcoreference_v2.tar.gz
  cd vua-eventcoreference_v2
  cp lib/EventCoreference-1.0-SNAPSHOT-jar-with-
  dependencies.jar /mnt/sdb1/pipelines/testnlpp/nlpp/env/java/jars
  ◇

```

Fragment referenced in 13b.

Script

```

"../bin/evcoref" 31c≡
  #!/bin/bash
  < set variables that point to the directory-structure 4c, ... >
  MODROOT=$modulesdir/vua-eventcoreference_v2
  RESOURCEDIR=$MODROOT/resources
  JARFILE=/mnt/sdb1/pipelines/testnlpp/nlpp/env/java/jars/EventCoreference-1.0-
  SNAPSHOT-jar-with-dependencies.jar

  JAVAMODULE=eu.newsreader.eventcoreference.naf.EventCorefWordnetSim
  JAVAOPTIONS="--method leacock-chodorow"
  JAVAOPTIONS="$JAVAOPTIONS --wn-lmf $RESOURCEDIR/cornetto2.1.lmf.xml"
  JAVAOPTIONS="$JAVAOPTIONS --sim 2.0"
  JAVAOPTIONS="$JAVAOPTIONS --
  relations XPOS_NEAR_SYNONYM#HAS_HYPERONYM#HAS_XPOS_HYPERONYM"

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

```

4 Utilities

4.1 Test script

The following script pushes a single sentence through the modules of the pipeline.

```
"../bin/test" 32a≡
    #!/bin/bash
    ROOT=/mnt/sdb1/pipelines/testnlpp/nlpp
    TESTDIR=$ROOT/test
    BIND=$ROOT/bin
    mkdir -p $TESTDIR
    cd $TESTDIR
    cat $ROOT/nuweb/testin.naf | $BIND/tok > $TESTDIR/test.tok.naf
    cat test.tok.naf | $BIND/mor > $TESTDIR/test.mor.naf
    cat test.mor.naf | $BIND/nerc_conll102 > $TESTDIR/test.nerc.naf
    cat $TESTDIR/test.nerc.naf | $BIND/wsd > $TESTDIR/test.wsd.naf
    cat $TESTDIR/test.wsd.naf | $BIND/ned > $TESTDIR/test.ned.naf
    cat $TESTDIR/test.ned.naf | $BIND/onto > $TESTDIR/test.onto.naf
    cat $TESTDIR/test.onto.naf | $BIND/heideltime > $TESTDIR/test.times.naf
    cat $TESTDIR/test.times.naf | $BIND/srl > $TESTDIR/test.srl.naf
    cat $TESTDIR/test.srl.naf | $BIND/evcoref > $TESTDIR/test.ecrf.naf
    cat $TESTDIR/test.ecrf.naf | $BIND/framesrl > $TESTDIR/test.fsrl.naf
```

◇

Uses: **nuweb** 41c.

4.2 Logging

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

```
< variables of install-modules 32b > ≡
    LOGLEVEL=1
```

◇

Fragment referenced in 13a.

```
< logmess 32c > ≡
    if
    [ $LOGLEVEL -gt 0 ]
    then
    echo @1
    fi
```

◇

Fragment referenced in 10a, 11bd, 15a, 18a, 33a.

4.3 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 33a ⟩ ≡
  SUCCES=0
  cd /mnt/sdb1/pipelines/testnlpp/nlpp/modules
  ⟨ move module (33b $DIR ) 10g ⟩
  wget $URL
  SUCCES=$?
  if
    [ $SUCCES -eq 0 ]
  then
    tar -xzf $TARB
    SUCCES=$?
    rm -rf $TARB
  fi
  if
    [ $SUCCES -eq 0 ]
  then
    ⟨ logmess (33c Installed $DIR ) 32c ⟩
    ⟨ remove old module (33d $DIR ) 11a ⟩
  else
    ⟨ re-instate old module (33e $DIR ) 11b ⟩
  fi
  ◇

```

Fragment never referenced.

A How to read and translate this document

This document is an example of *literate programming* [1]. 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 print-

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.

able/viewable documents and to extract the program sources. Table 3 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 T _E X documents into xml/html

Table 3: 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.

```
< parameters in Makefile 34a > ≡
  NUWEB=../env/bin/nuweb
  ◇
```

Fragment defined by 34a, 35d, 37c, 38a, 40b, 42b, 45a.
Fragment referenced in 34b.
Uses: `nuweb` 41c.

A.3 Translate and run

This chapter assembles the Makefile for this project.

```
"Makefile" 34b ≡
  < default target 35a >

  < parameters in Makefile 34a, ... >

  < implicate make regels 38b, ... >
  < expliciete make regels 36a, ... >
  < make targets 35b, ... >
  ◇
```

The default target of make is `all`.

```

< default target 35a > ≡
    all : < all targets 35c >
    .PHONY : all

```

◇

Fragment referenced in 34b.
 Defines: **all** Never used, **PHONY** 39a.

```

< make targets 35b > ≡
    clean:
        < clean up 6a, ... >

```

◇

Fragment defined by 35b, 39bc, 43c, 45bd.
 Fragment referenced in 34b.

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

```

< all targets 35c > ≡
    nlpp.pdf◇

```

Fragment referenced in 35a.
 Uses: **pdf** 39b.

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 35d > ≡
    .SUFFIXES: .pdf .w .tex .html .aux .log .php

```

◇

Fragment defined by 34a, 35d, 37c, 38a, 40b, 42b, 45a.
 Fragment referenced in 34b.
 Defines: **SUFFIXES** Never used.
 Uses: **pdf** 39b.

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.

\langle *expliciete make regels 36a* $\rangle \equiv$

```
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 36ac, 37ab, 39a, 40c, 42c, 43b.

Fragment referenced in 34b.

Uses: nuweb 41c.

\langle *clean up 36b* $\rangle \equiv$

```
rm -rf ../nuweb-1.58
```

◇

Fragment defined by 6a, 7b, 36b.

Fragment referenced in 35b.

Uses: nuweb 41c.

\langle *expliciete make regels 36c* $\rangle \equiv$

```
../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 36ac, 37ab, 39a, 40c, 42c, 43b.

Fragment referenced in 34b.

Uses: nuweb 41c.

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 37a > ≡
    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 36ac, 37ab, 39a, 40c, 42c, 43b.

Fragment referenced in 34b.

Uses: print 39b.

A.5.2 Run the M4 pre-processor

```

< expliciete make regels 37b > ≡
    nlpp.w : m4_nlpp.w inst.m4
            m4 -P m4_nlpp.w > nlpp.w

```

◇

Fragment defined by 36ac, 37ab, 39a, 40c, 42c, 43b.

Fragment referenced in 34b.

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 37c > ≡
    FIGFILES=fileschema directorystructure

```

◇

Fragment defined by 34a, 35d, 37c, 38a, 40b, 42b, 45a.

Fragment referenced in 34b.

Defines: FIGFILES 38a.

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 38a > ≡
    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 34a, 35d, 37c, 38a, 40b, 42b, 45a.

Fragment referenced in 34b.

Defines: FIGFILENAMES Never used, PDFT_NAMES 39c, PDF_FIG_NAMES 39c, PST_NAMES Never used,
PS_FIG_NAMES Never used.

Uses: FIGFILES 37c.

Create the graph files with program fig2dev:

```

< impliciete make regels 38b > ≡
    %.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 38b, 43a.

Fragment referenced in 34b.

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.

```
< explicit make regels 39a > ≡
    bibfile : nlpp.aux /home/paul/bin/mkportbib
              /home/paul/bin/mkportbib nlpp litprog
```

```
.PHONY : bibfile
```

```
◇
```

Fragment defined by 36ac, 37ab, 39a, 40c, 42c, 43b.

Fragment referenced in 34b.

Uses: PHONY 35a.

A.6.3 Create a printable/viewable document

Make a PDF document for printing and viewing.

```
< make targets 39b > ≡
    pdf : nlpp.pdf

    print : nlpp.pdf
           lpr nlpp.pdf

    view : nlpp.pdf
          evince nlpp.pdf
```

```
◇
```

Fragment defined by 35b, 39bc, 43c, 45bd.

Fragment referenced in 34b.

Defines: pdf 35cd, 39c, print 8a, 28a, 29a, 37a, 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 39c > ≡
    nlpp.pdf : nlpp.w $(W2PDF) $(PDF_FIG_NAMES) $(PDFT_NAMES)
              chmod 775 $(W2PDF)
              $(W2PDF) $*
```

```
◇
```

Fragment defined by 35b, 39bc, 43c, 45bd.

Fragment referenced in 34b.

Uses: pdf 39b, PDFT_NAMES 38a, PDF_FIG_NAMES 38a.

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 40a > \equiv
`../nuweb/bin` \diamond

Fragment defined by 3, 4ab, 5d, 6ce, 9c, 40a.

Fragment referenced in 45b.

Uses: **nuweb** 41c.

< parameters in Makefile 40b > \equiv
`W2PDF=../nuweb/bin/w2pdf`
 \diamond

Fragment defined by 34a, 35d, 37c, 38a, 40b, 42b, 45a.

Fragment referenced in 34b.

Uses: **nuweb** 41c.

< expliciete make regels 40c > \equiv
`$(W2PDF) : nlpp.w $(NUWEB)`
`$(NUWEB) nlpp.w`
 \diamond

Fragment defined by 36ac, 37ab, 39a, 40c, 42c, 43b.

Fragment referenced in 34b.

`"../nuweb/bin/w2pdf" 40d` \equiv
`#!/bin/bash`
`# w2pdf -- compile a nuweb file`
`# usage: w2pdf [filename]`
`# 20150629 at 0834h: Generated by nuweb from a_nlpp.w`
`NUWEB=../env/bin/nuweb`
`LATEXCOMPIER=pdflatex`
`< filenames in nuweb compile script 41a >`
`< compile nuweb 40e >`
 \diamond

Uses: **nuweb** 41c.

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 40e > \equiv
`NUWEB=/mnt/sdb1/pipelines/testnlpp/nlpp/env/bin/nuweb`
`< run the processors until the aux file remains unchanged 42a >`
`< remove the copy of the aux file 41b >`
 \diamond

Fragment referenced in 40d.

Uses: **nuweb** 41c.

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 41a ⟩ ≡
    nufil=$1
    trunk=${1%.*}
    texfil=${trunk}.tex
    auxfil=${trunk}.aux
    oldaux=old.${trunk}.aux
    indexfil=${trunk}.idx
    oldindexfil=old.${trunk}.idx
    ◇

```

Fragment referenced in 40d.

Defines: auxfil 42a, 44ab, indexfil 42a, 44a, nufil 41c, 44ac, oldaux 41b, 42a, 44ab, oldindexfil 42a, 44a, texfil 41c, 44ac, trunk 41c, 44acd.

Remove the old copy if it is no longer needed.

```

⟨ remove the copy of the aux file 41b ⟩ ≡
    rm $oldaux
    ◇

```

Fragment referenced in 40e, 43e.

Uses: oldaux 41a, 44a.

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 41c ⟩ ≡
    $NUWEB $nufil
    $LATEXCOMPILER $texfil
    makeindex $trunk
    bibtex $trunk
    ◇

```

Fragment referenced in 42a.

Defines: bibtex 44cd, makeindex 44cd, nuweb 22, 32a, 34a, 36abc, 40abde, 42b, 43d.

Uses: nufil 41a, 44a, texfil 41a, 44a, trunk 41a, 44a.

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 42a > ≡

```
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 41c >
  if [ $LOOPCOUNTER -ge 10 ]
  then
    cp $auxfil $oldaux
  fi;
done
◇
```

Fragment referenced in 40e.

Uses: auxfil 41a, 44a, indexfil 41a, oldaux 41a, 44a, oldindexfil 41a.

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 42b > ≡

```
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 34a, 35d, 37c, 38a, 40b, 42b, 45a.

Fragment referenced in 34b.

Uses: nuweb 41c.

Make the directory:

< expliciete make regels 42c > ≡

```
$(htmldir) :
  mkdir -p $(htmldir)
```

◇

Fragment defined by 36ac, 37ab, 39a, 40c, 42c, 43b.

Fragment referenced in 34b.

The rule to copy files in it:

```
< implicate make regels 43a > ≡
    $(htmldir)/% : % $(htmldir)
    cp $< $(htmldir)/
```

◇

Fragment defined by 38b, 43a.

Fragment referenced in 34b.

Do the work:

```
< expliciete make regels 43b > ≡
    $(htmltarget) : $(htmlmaterial) $(htmldir)
    cd $(htmldir) && chmod 775 w2html
    cd $(htmldir) && ./w2html nlpp.w
```

◇

Fragment defined by 36ac, 37ab, 39a, 40c, 42c, 43b.

Fragment referenced in 34b.

Invoke:

```
< make targets 43c > ≡
    htm : $(htmldir) $(htmltarget)
```

◇

Fragment defined by 35b, 39bc, 43c, 45bd.

Fragment referenced in 34b.

Create a script that performs the translation.

```
"w2html" 43d≡
#!/bin/bash
# w2html -- make a html file from a nuweb file
# usage: w2html [filename]
# [filename]: Name of the nuweb source file.
# 20150629 at 0834h: Generated by nuweb from a_nlpp.w
echo "translate " $1 >w2html.log
NUWEB=/mnt/sdb1/pipelines/testnlpp/nlpp/env/bin/nuweb
< filenames in w2html 44a >
```

```
< perform the task of w2html 43e >
```

◇

Uses: nuweb 41c.

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 43e > ≡
    < run the html processors until the aux file remains unchanged 44b >
    < remove the copy of the aux file 41b >
```

◇

Fragment referenced in 43d.

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 44a> ≡
  nufil=$1
  trunk=${1%.*}
  texfil=${trunk}.tex
  auxfil=${trunk}.aux
  oldaux=old.${trunk}.aux
  indexfil=${trunk}.idx
  oldindexfil=old.${trunk}.idx
  ◇
```

Fragment referenced in 43d.

Defines: `auxfil` 41a, 42a, 44b, `nufil` 41ac, 44c, `oldaux` 41ab, 42a, 44b, `texfil` 41ac, 44c, `trunk` 41ac, 44cd.

Uses: `indexfil` 41a, `oldindexfil` 41a.

```
<run the html processors until the aux file remains unchanged 44b> ≡
  while
    ! cmp -s $auxfil $oldaux
  do
    if [ -e $auxfil ]
    then
      cp $auxfil $oldaux
    fi
    <run the html processors 44c>
  done
  <run tex4ht 44d>
  ◇
```

Fragment referenced in 43e.

Uses: `auxfil` 41a, 44a, `oldaux` 41a, 44a.

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

```
<run the html processors 44c> ≡
  $NUWEB -o -n $nufil
  latex $texfil
  makeindex $trunk
  bibtex $trunk
  htlatex $trunk
  ◇
```

Fragment referenced in 44b.

Uses: `bibtex` 41c, `makeindex` 41c, `nufil` 41a, 44a, `texfil` 41a, 44a, `trunk` 41a, 44a.

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 44d> ≡
  tex '\def\filename{{nlpp}{idx}{4dx}{ind}} \input idxmake.4ht'
  makeindex -o $trunk.ind $trunk.4dx
  bibtex $trunk
  htlatex $trunk
  ◇
```

Fragment referenced in 44b.

Uses: `bibtex` 41c, `makeindex` 41c, `trunk` 41a, 44a.

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 45a ⟩ ≡
    MKDIR = mkdir -p
```

◇

Fragment defined by 34a, 35d, 37c, 38a, 40b, 42b, 45a.
 Fragment referenced in 34b.
 Defines: MKDIR 45b.

```
⟨ make targets 45b ⟩ ≡
    DIRS = ⟨ directories to create 3, ... ⟩
```

```
$(DIRS) :
    $(MKDIR) $@
```

◇

Fragment defined by 35b, 39bc, 43c, 45bd.
 Fragment referenced in 34b.
 Defines: DIRS 45d.
 Uses: MKDIR 45a.

```
⟨ make scripts executable 45c ⟩ ≡
    chmod -R 775 ../bin/*
```

◇

Fragment defined by 14a, 45c.
 Fragment referenced in 45d.

```
⟨ make targets 45d ⟩ ≡
    sources : nlpp.w $(DIRS) $(NUWEB)
              $(NUWEB) nlpp.w
              ⟨ make scripts executable 14a, ... ⟩
```

◇

Fragment defined by 35b, 39bc, 43c, 45bd.
 Fragment referenced in 34b.
 Uses: DIRS 45b.

B References

B.1 Literature

References

- [1] Donald E. Knuth. Literate programming. Technical report STAN-CS-83-981, Stanford University, Department of Computer Science, 1983.

B.2 URL's

Nuweb: nuweb.sourceforge.net

Apache Velocity: [m4_velocityURL](#)

Velocitytools: [m4_velocitytoolsURL](#)

Parameterparser tool: [m4_parameterparserdocURL](#)

Cookietool: [m4_cookietooldocURL](#)

VelocityView: [m4_velocityviewURL](#)

VelocityLayoutServlet: [m4_velocitylayoutervletURL](#)

Jetty: [m4_jettycodehausURL](#)

UserBase javadoc: [m4_userbasejavadocURL](#)

VU corpus Management development site: <http://code.google.com/p/vucom>

C Indexes

C.1 Filenames

"../bin/coreference-base" Defined by 21c.

"../bin/evcoref" Defined by 31c.

"../bin/framesrl" Defined by 28a.

"../bin/heideltime" Defined by 29b.

"../bin/install-modules" Defined by 13ab.

"../bin/lu2synset" Defined by 25c.

"../bin/mor" Defined by 21a.

"../bin/ned" Defined by 26b.

"../bin/nerc_conll02" Defined by 23c.

"../bin/nerc_sonar" Defined by 23d.

"../bin/onto" Defined by 27.

"../bin/srl" Defined by 30abcde.

"../bin/test" Defined by 32a.

"../bin/tok" Defined by 20b.

"../bin/wsd" Defined by 24e.

"../env/bin/progenv" Defined by 5c.

"../nuweb/bin/w2pdf" Defined by 40d.

"Makefile" Defined by 34b.

"w2html" Defined by 43d.

C.2 Macro's

<activate the python environment 9bd> Referenced in 5c, 7c.

<adapt heideltime's config.props 29a> Referenced in 28b.

<all targets 35c> Referenced in 35a.

<check this first 14b> Referenced in 13a.

<check whether mercurial is present 14c> Referenced in 14b.

<check/install the correct version of python 8a> Referenced in 7c.

<check/start the Spotlight server 19e> Referenced in 26b.

<clean up 6a, 7b, 36b> Referenced in 35b.

<compile nuweb 40e> Referenced in 40d.

<compile the nerc jar 22> Referenced in 21d.

<create a virtual environment for Python 8e> Referenced in 7c.

<default target 35a> Referenced in 34b.

<directories to create 3, 4ab, 5d, 6ce, 9c, 40a> Referenced in 45b.

<download svm models 24c> Referenced in 24a.

<explicitly make regels 36ac, 37ab, 39a, 40c, 42c, 43b> Referenced in 34b.

<filenames in nuweb compile script 41a> Referenced in 40d.

<filenames in w2html 44a> Referenced in 43d.

<get or have 12b> Referenced in 5e, 8b, 19a, 23a, 24c, 25a, 26c, 31a.

<get the nerc models 23a> Referenced in 21d.
 <have an SSH key or die 12a> Referenced in 12b.
 <impliciete make regels 38b, 43a> Referenced in 34b.
 <install ActivePython 8bd> Referenced in 8a.
 <install Alpino 15a> Referenced in 13a.
 <install coreference-base 21b> Referenced in 13b.
 <install from github 11d> Referenced in 20c, 21b, 24a, 25d, 28b, 29c.
 <install from tarball 33a> Not referenced.
 <install kafnafparserpy 10a> Referenced in 7c.
 <install maven 6f, 7a> Referenced in 13a.
 <install python packages 10f> Referenced in 7c.
 <install svm lib 24b> Referenced in 24a.
 <install the event-coreference module 31a> Referenced in 13b.
 <install the heideltime module 28b> Referenced in 13b.
 <install the lu2synset converter 25a> Referenced in 13b.
 <install the morphosyntactic parser 20c> Referenced in 13b.
 <install the NERC module 21d> Referenced in 13b.
 <install the onto module 26c> Referenced in 13b.
 <install the Spotlight server 19ac> Referenced in 13a.
 <install the srl module 29c> Referenced in 13b.
 <install the ticcutils utility 17b> Referenced in 13a.
 <install the timbl utility 17c> Referenced in 13a.
 <install the tokenizer 20a> Referenced in 13b.
 <install the treetagger utility 15g, 16abcde, 17a> Referenced in 13a.
 <install the WSD module 24a> Referenced in 13b.
 <install the NED module 25d> Referenced in 13b.
 <logmess 32c> Referenced in 10a, 11bd, 15a, 18a, 33a.
 <make scripts executable 14a, 45c> Referenced in 45d.
 <make targets 35b, 39bc, 43c, 45bd> Referenced in 34b.
 <move module 10g> Referenced in 10a, 11d, 15a, 33a.
 <parameters in Makefile 34a, 35d, 37c, 38a, 40b, 42b, 45a> Referenced in 34b.
 <perform the task of w2html 43e> Referenced in 43d.
 <put spotlight jar in the Maven repository 26a> Referenced in 25d.
 <re-instate old module 11b> Referenced in 10a, 11d, 15a, 33a.
 <remove old module 11a> Referenced in 10a, 11d, 15a, 33a.
 <remove the copy of the aux file 41b> Referenced in 40e, 43e.
 <run tex4ht 44d> Referenced in 44b.
 <run the html processors 44c> Referenced in 44b.
 <run the html processors until the aux file remains unchanged 44b> Referenced in 43e.
 <run the processors until the aux file remains unchanged 42a> Referenced in 40e.
 <run the three processors 41c> Referenced in 42a.
 <set alpinohome 15f> Referenced in 21a.
 <set local bin directory ?> Referenced in 13a.
 <set up java 5e> Referenced in 13a.
 <set up java environment in scripts 6bd> Referenced in 5c, 13a.
 <set up python 7c> Referenced in 13a.
 <set variables that point to the directory-structure 4c, 5ab> Referenced in 20b, 21ac, 23cd, 24e, 26b, 27, 28a, 29b, 30a, 31c.
 <start the Spotlight server 19d> Referenced in 19e.
 <test whether virtualenv is present on the host 9a> Referenced in 8e.
 <unpack ticcutils or timbl 18a> Referenced in 17bc.
 <variables of install-modules 32b> Referenced in 13a.

C.3 Variables

activate: 9b.

all: 35a.

ALPINO_HOME: 15f.

auxfil: [41a](#), [42a](#), [44a](#), [44b](#).
bibtex: [41c](#), [44cd](#).
DIRS: [45b](#), [45d](#).
fig2dev: [38b](#).
FIGFILENAMES: [38a](#).
FIGFILES: [37c](#), [38a](#).
hg: [14c](#), [21b](#).
indexfil: [41a](#), [42a](#), [44a](#).
JAVA_HOME: [6b](#).
lxml: [10f](#).
makeindex: [41c](#), [44cd](#).
MKDIR: [45a](#), [45b](#).
nufil: [41a](#), [41c](#), [44a](#), [44c](#).
nuweb: [22](#), [32a](#), [34a](#), [36abc](#), [40abde](#), [41c](#), [42b](#), [43d](#).
oldaux: [41a](#), [41b](#), [42a](#), [44a](#), [44b](#).
oldindexfil: [41a](#), [42a](#), [44a](#).
PATH: [5a](#), [6b](#), [7a](#).
pdf: [35cd](#), [39b](#), [39c](#).
PDFT_NAMES: [38a](#), [39c](#).
PDF_FIG_NAMES: [38a](#), [39c](#).
PHONY: [35a](#), [39a](#).
print: [8a](#), [28a](#), [29a](#), [37a](#), [39b](#).
PST_NAMES: [38a](#).
PS_FIG_NAMES: [38a](#).
pythonok: [8a](#).
PYTHONPATH: [9d](#).
pyyaml: [10f](#).
SUCCES: [15a](#), [18a](#), [33a](#).
SUFFIXES: [35d](#).
texfil: [41a](#), [41c](#), [44a](#), [44c](#).
trunk: [41a](#), [41c](#), [44a](#), [44cd](#).
view: [39b](#).
virtualenv: [8de](#), [9a](#).