# Install Dutch nlp modules on Lisa

**Paul Huygen <paul.huygen@huygen.nl>**

**3rd November 2014**
**09:26 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

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 an Terminology Lab (CLTL [1]) as part of the newsreader [2].

Apart from describing the pipeline set-up, the document actually constructs the pipeline. The described version has been made with an aim to run it on a specific supercomputer (Lisa, Surfsara, Amsterdam [3]), but it can probably be implemented on other unix-like systems without problems.

The installation has been parameterized. 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. Ideally, modules are directly obtained from a public repository, e.g. Github, or from a website of the organisation where the module has been built. However, some of the modules are not yet available in this way and only an informal snapshot is available. Table /reftab:modulesources provides the URL's of the sources that have been obtained from a public repository.

———

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

| module | directory | source | script | Details |
|---|---|---|---|---|
| Tokenizer | `tokenizer-base` | Github | tok | |
| morphosyntactic parser | `morphosyntactic_parser_nl` | Github | `mor` | |
| alpinohack | `clean_hack` | This doc. | alpinohack | [4] |
| NERC | `../modules/jars` | Snap | nerc | |
| WSD | `wsd` | Snap | wsd | |
| Onto | `ontotagger` | Snap | onto | |
| Heidel | `NAF-HeidelTime` | Github | heideltime | |
| SRL | `vua-srl-nl` | Github | srl | |
| NED | `ned` | None | ned | |
| Nom. coref | `/dev/null` | None | nomcoref | |
| Ev. coref | `/dev/null` | None | evcoref | |
| Opinion miner | `/dev/null` | None | opinimin | |
| Framenet sem. role label. | `/dev/null` | None | 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; **script:** Script to be included in a pipeline.

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

| module | directory | source | Details |
|---|---|---|---|
| KafNafParserPy | `python/KafNafParserPy` | Github | |
| Alpino | `Alpino` | RUG | |
| Ticcutils | `ticcutils-0.7` | ILK | |
| Timbl | `timbl-6.4.6` | ILK | |
| Treetagger | | | |

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.

Table 3 lists the source of the modules and utilities that can be installed from an open source.

| module | source | URL |
|---|---|---|
| Tokenizer | Github | https://github.com/opener-project/tokenizer-base.git |
| Morphosynt. p. | Github | `https://github.com/cltl/morphosyntactic_parser_nl.git` |
| heideltime. | Github | `https://github.com/cltl/morphosyntactic_parser_nl.git` |
| Alpino | RUG | `Alpino-x86_64-linux-glibc2.5-20548-sicstus.tar.gz` |
| Ticcutils | ILK | ticcutils-0.7.tar.gz |
| Timble | ILK | timbl-6.4.6.tar.gz |

Table 3: Sources of the modules

The informal snapshots are available in a tarball `nl_pipeline_snapshots` that can be obtained from the author of this document.

## 1.2   File-structure of the pipeline

The files that make up the pipeline are organised in set of directories:

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

**modules:** Contains the program code of each module in a subdirectory. Furthermore, it contains a subdirectory `python` for python software-modules, subdirectory `jars` for jar files and subdirectory /usrlocal/ for binaries and libs that are used by modules.

**bin:** Contains for each of the 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.

**nuweb:** Contains this document, the nuweb source that creates the documents and the sources and a Makefile to perform the actions.

⟨ *directories to create* 4a ⟩ ≡
            ../modules   ⋄
Fragment defined by 4abcdef, 9c, 26c.
Fragment referenced in 32b.

⟨ *directories to create* 4b ⟩ ≡
            ../bin   ⋄
Fragment defined by 4abcdef, 9c, 26c.
Fragment referenced in 32b.

⟨ *directories to create* 4c ⟩ ≡
            ../modules/usrlocal   ⋄
Fragment defined by 4abcdef, 9c, 26c.
Fragment referenced in 32b.

⟨ *directories to create* 4d ⟩ ≡
            ../modules/usrlocal/bin   ⋄
Fragment defined by 4abcdef, 9c, 26c.
Fragment referenced in 32b.

⟨ *directories to create* 4e ⟩ ≡
            ../modules/usrlocal/lib   ⋄
Fragment defined by 4abcdef, 9c, 26c.
Fragment referenced in 32b.

⟨ *directories to create* 4f ⟩ ≡
            ../modules/python /home/paul/dutch-nlp-modules-on-Lisa/modules/jars   ⋄
Fragment defined by 4abcdef, 9c, 26c.
Fragment referenced in 32b.

Make Python utilities findable with the following macro:

⟨ *set pythonpath* 4g ⟩ ≡
            export PYTHONPATH=/home/paul/dutch-nlp-modules-on-Lisa/modules/python:$PYTHONPATH
            ⋄
Fragment referenced in 9a, 14b, 15a.

Similarly, make binaries findable:

⟨ *set local bin directory* 4h ⟩ ≡
            export PATH=../modules/usrlocal/bin:$PATH
            ⋄
Fragment referenced in 15a.

During installation, an extra directory, `snapshot`, that contains modules that are not yet available from public sources, is needed.

So, let us here and now check whether the snapshots are indeed present. The following macro unpacks the tarball if it is present and aborts the installation when the snapshot directory is not present.

⟨ *unpack snapshots or die* 5a ⟩ ≡
```
cd /home/paul/dutch-nlp-modules-on-Lisa
if
  [ -e nl_pipeline_snapshots.tgz ]
then
  tar -zxf nl_pipeline_snapshots.tgz
fi
if
  [ ! -e snapshots ]
then
  echo "No module snapshots"
  exit 1
fi
```
◇

Fragment referenced in 5b.

## 2 Installation

This section describes how the modules are obtained from their open-source and installed. This is performed by script `install-modules`

`"../bin/install-modules"` 5b≡
```
#!/bin/bash
```
⟨ *variables of install-modules* 19d ⟩
⟨ *unpack snapshots or die* 5a ⟩
⟨ *install the tokenizer* 7h ⟩
⟨ *install kafnafparserpy* 16a ⟩
⟨ *install Alpino* 8b ⟩
⟨ *install the morphosyntactic parser* 8h ⟩
⟨ *install the NERC module* 10c ⟩
⟨ *install the WSD module* 11a ⟩
⟨ *install the* NED *module* 12a ⟩
⟨ *install the onto module* 12d ⟩
⟨ *install the heideltime module* 13c ⟩
⟨ *install the srl module* 14d ⟩
⟨ *install the treetagger utility* 16g, . . . ⟩
⟨ *install the ticcutils utility* 18a ⟩
⟨ *install the timbl utility* 18b ⟩

◇

⟨ *make scripts executable* 5c ⟩ ≡
```
chmod 775  /home/paul/dutch-nlp-modules-on-Lisa/bin/install-modules
```
◇

Fragment defined by 5c, 8a, 9b, 10be, 11c, 12c, 13b, 14c, 15f, 16f.
Fragment referenced in 32c.

Installation goes as follows:

1.      If the module exists already, move it to a temporary place.
2.      Try to install the module from the source.
3.      If that is successful, remove the old version. Otherwise, move the old version back to its
        original place.

The following macro's move or remove modules.

⟨ *move module* 6a ⟩ ≡
```
if
  [ -e @1 ]
then
    mv @1 old.@1
fi
```
        ◇
Fragment referenced in 7a, 8b, 20a.

⟨ *remove old module* 6b ⟩ ≡
```
rm -rf old.@1
```
        ◇
Fragment referenced in 7a, 8b, 20a.

⟨ *re-instate old module* 6c ⟩ ≡
```
mv old.@1 @1
MESS="Replaced previous version of @1"
```
        ⟨ *logmess* (6d $MESS ) 19e ⟩

        ◇
Fragment referenced in 7a, 8b, 20a.

The following macro can be used to install a module from github. It needs as parameters:

1.      Name of the module.
2.      Name of the root directory.
3.      URL to clone from.

⟨ *install from github* 7a ⟩ ≡
```
    MODNAM=@1
    DIRN=@2
    GITU=@3
```
    ⟨ *find leave and tree* 7g ⟩
    ⟨ *logmess* (7b `"TREE: $TREE; LEAVE: $LEAVE"` ) 19e ⟩
```
    cd $TREE
```
    ⟨ *move module* (7c `$LEAVE` ) 6a ⟩
```
    git clone $GITU
    if
      [ $? -gt 0 ]
    then
```
      ⟨ *logmess* (7d `Cannot install current $MODNAM version` ) 19e ⟩
      ⟨ *re-instate old module* (7e `$LEAVE` ) 6c ⟩
```
    else
```
      ⟨ *remove old module* (7f `$LEAVE` ) 6b ⟩
```
    fi
```

    ◇

Fragment referenced in 7h, 8h, 13c, 14d, 16a.


Note: Par. 1: Directory; par 2: path to directory; par 3: directory name.

⟨ *find leave and tree* 7g ⟩ ≡
```
    FULLDIR=/home/paul/dutch-nlp-modules-on-Lisa/modules/$DIRN
    LEAVE=${FULLDIR##*/}
    TREE=${FULLDIR%%$LEAVE}
```
    ◇
Fragment referenced in 7a.




## 2.1   Install tokenizer

### 2.1.1   Module

⟨ *install the tokenizer* 7h ⟩ ≡

    ⟨ *install from github* (7i `tokenizer`,7j `tokenizer-base`,7k `https://github.com/opener-project/tokenizer-base.git`
    ◇
Fragment referenced in 5b.




### 2.1.2   Script

The script just runs the tokenizerscript in Perl.

`"../bin/tok"` 7l≡
```
    #!/bin/bash
    ROOT=/home/paul/dutch-nlp-modules-on-Lisa
    TOKBINDIR=/home/paul/dutch-nlp-modules-on-Lisa/modules/tokenizer-base/core
    cat | perl $TOKBINDIR/tokenizer-cli.pl -l nl t
```

    ◇

⟨ *make scripts executable* 8a ⟩ ≡
```
      chmod 775  /home/paul/dutch-nlp-modules-on-Lisa/bin/tok
```
      ◇

Fragment defined by 5c, 8a, 9b, 10be, 11c, 12c, 13b, 14c, 15f, 16f.
Fragment referenced in 32c.

## 2.2   Install Alpino

Install Alpino from the website of Gertjan van Noort.

### 2.2.1   Module

⟨ *install Alpino* 8b ⟩ ≡
```
      SUCCES=0
      cd /home/paul/dutch-nlp-modules-on-Lisa/modules
```
      ⟨ *move module* (8c **Alpino** ) 6a ⟩
```
      wget http://www.let.rug.nl/vannoord/alp/Alpino/binary/versions/Alpino-x86_64-linux-glibc2.5-20548-sic
      SUCCES=$?
      if
        [ $SUCCES -eq 0 ]
      then
        tar -xzf Alpino-x86_64-linux-glibc2.5-20548-sicstus.tar.gz
        SUCCES=$?
        rm -rf Alpino-x86_64-linux-glibc2.5-20548-sicstus.tar.gz
      fi
      if
        [ $SUCCES -eq 0 ]
      then
```
        ⟨ *logmess* (8d `Installed Alpino` ) 19e ⟩
        ⟨ *remove old module* (8e **Alpino** ) 6b ⟩
```
      else
```
        ⟨ *re-instate old module* (8f **Alpino** ) 6c ⟩
```
      fi
```
      ◇
Fragment referenced in 5b.

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* 8g ⟩ ≡
```
      export ALPINO_HOME=/home/paul/dutch-nlp-modules-on-Lisa/modules/Alpino
```
      ◇
Fragment referenced in 9a.
Defines: `ALPINO_HOME` Never used.

## 2.3   Morphosyntactic parser

### 2.3.1   Module

⟨ *install the morphosyntactic parser* 8h ⟩ ≡

      ⟨ *install from github* (8i `morphsynparser`,8j `morphosyntactic_parser_nl`,8k `https://github.com/cltl/morphosyntac`
      ◇
Fragment referenced in 5b.

### 2.3.2   Script

```
"../bin/mor" 9a≡
      #!/bin/bash
      ROOT=/home/paul/dutch-nlp-modules-on-Lisa
      MODDIR=/home/paul/dutch-nlp-modules-on-Lisa/modules/morphosyntactic_parser_nl
      ⟨ set alpinohome 8g ⟩
      ⟨ set pythonpath 4g ⟩
      cat | python $MODDIR/core/morph_syn_parser.py
      ◇
```

⟨ *make scripts executable* 9b ⟩ ≡
```
      chmod 775  /home/paul/dutch-nlp-modules-on-Lisa/bin/mor
      ◇
```
Fragment defined by 5c, 8a, 9b, 10be, 11c, 12c, 13b, 14c, 15f, 16f.
Fragment referenced in 32c.

## 2.4   Alpino hack

Install a hack that removes output from Alpino that cannot be interpreted by following modules.
It is just a small python script.

### 2.4.1   Module

⟨ *directories to create* 9c ⟩ ≡
```
      /home/paul/dutch-nlp-modules-on-Lisa/modules/alpinohack ◇
```
Fragment defined by 4abcdef, 9c, 26c.
Fragment referenced in 32b.

```
"../modules/alpinohack/clean_hack.py" 9d≡
      #!/usr/bin/python
      import sys

      input = sys.stdin

      output = ''

      for line in input:
          line = line.replace('"--','"#')
          line = line.replace('--"','#"')
          output += line

      print output

      ◇
```
Uses: `print` 26a.

### 2.4.2  Script

```
"../bin/alpinohack" 10a≡
     #!/bin/bash
     ROOT=/home/paul/dutch-nlp-modules-on-Lisa
     HACKDIR=/home/paul/dutch-nlp-modules-on-Lisa/modules/alpinohack
     cat | python  $HACKDIR/clean_hack.py


     ◊
```

⟨ *make scripts executable* 10b ⟩ ≡
```
     chmod 775  /home/paul/dutch-nlp-modules-on-Lisa/bin/alpinohack
     ◊
```
Fragment defined by 5c, 8a, 9b, 10be, 11c, 12c, 13b, 14c, 15f, 16f.
Fragment referenced in 32c.

## 2.5   Named entity recognition (NERC)

### 2.5.1  Module

We do not (yet) have the source code of the NER module. A snapshot is comprised in a jar library.

⟨ *install the NERC module* 10c ⟩ ≡
```
     cp  /home/paul/dutch-nlp-modules-on-Lisa/snapshots/nerc/ixa-pipe-nerc-1.1.0.jar /home/paul/dutch-nlp-
     ◊
```
Fragment referenced in 5b.

### 2.5.2  Script

```
"../bin/nerc" 10d≡
     #!/bin/bash
     ROOT=/home/paul/dutch-nlp-modules-on-Lisa
     JARDIR=/home/paul/dutch-nlp-modules-on-Lisa/modules/jars
     cat | java -jar $JARDIR/ixa-pipe-nerc-1.1.0.jar tag
     ◊
```

⟨ *make scripts executable* 10e ⟩ ≡
```
     chmod 775  /home/paul/dutch-nlp-modules-on-Lisa/bin/nerc
     ◊
```
Fragment defined by 5c, 8a, 9b, 10be, 11c, 12c, 13b, 14c, 15f, 16f.
Fragment referenced in 32c.

## 2.6   Wordsense-disambiguation

We do not yet have a source-repository of the wsd module. Therefore, install from a snapshot on
Lisa.

### 2.6.1   Module

⟨ *install the WSD module* 11a ⟩ ≡

```
cp -r /home/paul/dutch-nlp-modules-on-Lisa/snapshots/wsd /home/paul/dutch-nlp-modules-on-Lisa/modules
◇
```
Fragment referenced in 5b.

### 2.6.2   Script

`"../bin/wsd"` 11b≡
```
#!/bin/bash
# WSD -- wrapper for word-sense disambiguation
# 8 Jan 2014 Ruben Izquierdo
# 16 sep 2014 Paul Huygen
ROOT=/home/paul/dutch-nlp-modules-on-Lisa
WSDDIR=/home/paul/dutch-nlp-modules-on-Lisa/modules/wsd
WSDSCRIPT=kaf_annotate_senses.pl
UKB=$WSDDIR/ukb_wsd_2.0
POSMAP=$WSDDIR/posmap.NGV.txt

if [ "$1" = "nl" ]
then
  GRAPH=$WSDDIR/cdb2.0-nld-all.infv.0.0.no-allwords.bin
  DICT=$WSDDIR/dictionary
else
  GRAPH=$WSDDIR/wn30g_eng.v20.bin
  DICT=$WSDDIR/wn30_eng_dict.txt
fi

iconv -t utf-8//IGNORE | $WSDDIR/$WSDSCRIPT -x $UKB -M $GRAPH -W $DICT -m $POSMAP
◇
```
Uses: `all` 22a.

⟨ *make scripts executable* 11c ⟩ ≡
```
chmod 775  /home/paul/dutch-nlp-modules-on-Lisa/bin/wsd
◇
```
Fragment defined by 5c, 8a, 9b, 10be, 11c, 12c, 13b, 14c, 15f, 16f.
Fragment referenced in 32c.

## 2.7   NED

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.

### 2.7.1   Installation of the spotlight server

### 2.7.2   Module

⟨ *install the* NED *module* 12a ⟩ ≡

```
cp /home/paul/dutch-nlp-modules-on-Lisa/snapshots/ned/ixa-pipe-ned-1.0.jar /home/paul/dutch-nlp-modul
mkdir -p /home/paul/dutch-nlp-modules-on-Lisa/modules/ned
cd /home/paul/dutch-nlp-modules-on-Lisa/modules/ned
wget http://ixa2.si.ehu.es/ixa-pipes/models/wikipedia-db.v1.tar.gz
tar -xzf wikipedia-db.v1.tar.gz
```
◇

Fragment referenced in 5b.

### 2.7.3   Script

"../bin/ned" 12b≡

```
#!/bin/bash
ROOT=/home/paul/dutch-nlp-modules-on-Lisa
JARDIR=/home/paul/dutch-nlp-modules-on-Lisa/modules/jars
cat | java -jar $JARDIR/ixa-pipe-ned-1.0.jar  -p 2060 -e candidates -i /home/paul/dutch-nlp-modules-o
```
◇

⟨ *make scripts executable* 12c ⟩ ≡

```
chmod 775  /home/paul/dutch-nlp-modules-on-Lisa/bin/ned
```
◇

Fragment defined by 5c, 8a, 9b, 10be, 11c, 12c, 13b, 14c, 15f, 16f.
Fragment referenced in 32c.

## 2.8   Ontotagger

We do not yet have a source-repository of the Ontotagger module. Therefore, install from a snap-shot on Lisa.

### 2.8.1   Module

⟨ *install the onto module* 12d ⟩ ≡

```
cp -r /home/paul/dutch-nlp-modules-on-Lisa/snapshots/ontotagger /home/paul/dutch-nlp-modules-on-Lisa/
```
◇

Fragment referenced in 5b.

### 2.8.2   Script

`"../bin/onto"` 13a≡

```
#!/bin/bash
ROOT=/home/paul/dutch-nlp-modules-on-Lisa
ONTODIR=/home/paul/dutch-nlp-modules-on-Lisa/modules/ontotagger
JARDIR=$ONTODIR/lib
RESOURCESDIR=$ONTODIR/resources
PREDICATEMATRIX="$RESOURCESDIR/PredicateMatrix.v1.1/PredicateMatrix.v1.1.role.nl-1.merged"
GRAMMATICALWORDS="$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


JAVA_ARGS="--mappings \"fn;pb;nb\" "
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

◇
```

⟨ *make scripts executable* 13b ⟩ ≡
```
chmod 775  /home/paul/dutch-nlp-modules-on-Lisa/bin/onto
◇
```
Fragment defined by 5c, 8a, 9b, 10be, 11c, 12c, 13b, 14c, 15f, 16f.
Fragment referenced in 32c.

## 2.9   Heideltime

### 2.9.1   Module

⟨ *install the heideltime module* 13c ⟩ ≡

⟨ *install from github* (13d `heideltime`,13e `NAF-HeidelTime`,13f `git@github.com:PaulHuygen/NAF-HeidelTime.git` ) 7⟩
⟨ *adapt heideltime's config.props* 14a ⟩


◇
Fragment referenced in 5b.

⟨ *adapt heideltime's config.props* 14a ⟩ ≡

```
CONFIL=NAF-HeidelTime/config.props
tempfil=`mktemp -t heideltmp.XXXXXX`
mv $CONFIL $tempfil
MODDIR=/home/paul/dutch-nlp-modules-on-Lisa/modules
TREETAGDIR=treetagger
AWKCOMMAND='/^treeTaggerHome/ {$0="treeTaggerHome = /home/paul/dutch-nlp-modules-on-Lisa/modules/tree
gawk "$AWKCOMMAND" $tempfil >$CONFIL
```
                ◇

Fragment referenced in 13c.
Uses: print 26a.

### 2.9.2   Script

`"../bin/heideltime"` 14b≡

```
#!/bin/bash
ROOT=/home/paul/dutch-nlp-modules-on-Lisa
HEIDELDIR=/home/paul/dutch-nlp-modules-on-Lisa/modules/NAF-HeidelTime
TEMPDIR=`mktemp -t -d heideltmp.XXXXXX`
cd $HEIDELDIR
```
        ⟨ *set pythonpath* 4g ⟩
```
iconv -t utf-8//IGNORE | python $HEIDELDIR/HeidelTime_NafKaf.py $HEIDELDIR/heideltime-standalone/ $TE
```
                ◇

⟨ *make scripts executable* 14c ⟩ ≡

```
chmod 775  /home/paul/dutch-nlp-modules-on-Lisa/bin/heideltime
```
                ◇

Fragment defined by 5c, 8a, 9b, 10be, 11c, 12c, 13b, 14c, 15f, 16f.
Fragment referenced in 32c.

## 2.10   Semantic Role labelling

### 2.10.1  Module

⟨ *install the srl module* 14d ⟩ ≡

        ⟨ *install from github* (14e `srl`,14f `vua-srl-nl`,14g `https://github.com/newsreader/vua-srl-nl.git`  ) 7a ⟩
        ◇

Fragment referenced in 5b.

### 2.10.2  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" 15a≡
      #!/bin/bash
      ROOT=/home/paul/dutch-nlp-modules-on-Lisa
      SRLDIR=/home/paul/dutch-nlp-modules-on-Lisa/modules/vua-srl-nl
      TEMPDIR='mktemp -d -t SRLTMP.XXXXXX'
      cd $SRLDIR
```
⟨ *set local bin directory* 4h ⟩
⟨ *set pythonpath* 4g ⟩
```
      INPUTFILE=$TEMPDIR/inputfile
      FEATUREVECTOR=$TEMPDIR/csvfile
      TIMBLOUTPUTFILE=$TEMPDIR/timblpredictions
      ◇
```
File defined by 15abcde.

Create a feature-vector.

```
"../bin/srl" 15b≡
      cat | tee  $INPUTFILE | python nafAlpinoToSRLFeatures.py > $FEATUREVECTOR
      ◇
```
File defined by 15abcde.

Run the trained model on the feature-vector.

```
"../bin/srl" 15c≡

      timbl -mO:I1,2,3,4 -i e-mags_mags_press_newspapers.wgt -t $FEATUREVECTOR -o $TIMBLOUTPUTFILE >/dev/nu
      ◇
```
File defined by 15abcde.

Insert the SRL values into the NAF file.

```
"../bin/srl" 15d≡
      python timblToAlpinoNAF.py $INPUTFILE $TIMBLOUTPUTFILE
      ◇
```
File defined by 15abcde.

Clean up.

```
"../bin/srl" 15e≡
      rm -rf $TEMPDIR
      ◇
```
File defined by 15abcde.

⟨ *make scripts executable* 15f ⟩ ≡
```
      chmod 775  /home/paul/dutch-nlp-modules-on-Lisa/bin/srl
      ◇
```
Fragment defined by 5c, 8a, 9b, 10be, 11c, 12c, 13b, 14c, 15f, 16f.
Fragment referenced in 32c.

## 2.11   KafNafParserPy

Several modules use KafNafParserpy to read and write NAF files.

### 2.11.1 Module

⟨ *install kafnafparserpy* 16a ⟩ ≡

⟨ *install from github* (16b `kafnafparserpy`,16c `python/KafNafParserPy`,16d `https://github.com/cltl/KafNafParserF`
◇
Fragment referenced in 5b.

## 3   Utilities

### 3.1   Test script

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

`"../bin/test"` 16e≡
```
#!/bin/bash
ROOT=/home/paul/dutch-nlp-modules-on-Lisa
BIND=$ROOT/bin
echo "De hond eet jus." | $BIND/tok | $BIND/mor | \
$BIND/alpinohack | $BIND/nerc  | $BIND/wsd | \
$BIND/onto  > $ROOT/test.onto
cat $ROOT/test.onto | $BIND/heideltime  > $ROOT/test.heidel
cat $ROOT/test.heidel | $BIND/srl  > $ROOT/test.srl
cat $ROOT/test.srl | $BIND/srl  > $ROOT/test.srl
```
◇

⟨ *make scripts executable* 16f ⟩ ≡
```
chmod 775  /home/paul/dutch-nlp-modules-on-Lisa/bin/test
```
◇
Fragment defined by 5c, 8a, 9b, 10be, 11c, 12c, 13b, 14c, 15f, 16f.
Fragment referenced in 32c.

### 3.2   Treetagger

### 3.2.1   Module

Installation 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 treetager comes from and the location where it is going to reside:

⟨ *install the treetagger utility* 16g ⟩ ≡
```
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 16g, 17abcde.
Fragment referenced in 5b.

The source tarball, scripts and the installation-script:

⟨ *install the treetagger utility* 17a ⟩ ≡

```
TREETAGSRC=tree-tagger-linux-3.2.tar.gz
TREETAGSCRIPTS=tagger-scripts.tar.gz
TREETAG_INSTALLSCRIPT=install-tagger.sh
```
◇

Fragment defined by 16g, 17abcde.
Fragment referenced in 5b.

Parametersets:

⟨ *install the treetagger utility* 17b ⟩ ≡

```
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 16g, 17abcde.
Fragment referenced in 5b.

Download everything in the target directory:

⟨ *install the treetagger utility* 17c ⟩ ≡

```
mkdir -p /home/paul/dutch-nlp-modules-on-Lisa/modules/$TREETAGDIR
cd /home/paul/dutch-nlp-modules-on-Lisa/modules/$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 16g, 17abcde.
Fragment referenced in 5b.

Run the install-script:

⟨ *install the treetagger utility* 17d ⟩ ≡

```
chmod 775 $TREETAG_INSTALLSCRIPT
./$TREETAG_INSTALLSCRIPT
```
◇

Fragment defined by 16g, 17abcde.
Fragment referenced in 5b.

Remove the tarballs:

⟨ *install the treetagger utility* 17e ⟩ ≡

```
rm $TREETAGSRC
rm $TREETAGSCRIPTS
rm $TREETAG_INSTALLSCRIPT
rm $DUTCHPARS_UTF_GZ
rm $DUTCH_TAGSET
rm $DUTCHPARS_2_GZ
```
◇

Fragment defined by 16g, 17abcde.
Fragment referenced in 5b.

### 3.3   Timbl and ticcutils

#### 3.3.1   Module

Timbl and ticcutils are installed from their source-tarballs. The installation is not (yet?) completely reproducibe because it uses the currently available c-compiler. 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 `usrlocal` subdirectory of the modules directory.

⟨ *install the ticcutils utility* 18a ⟩ ≡
```
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* 19a ⟩
◇

Fragment referenced in 5b.

⟨ *install the timbl utility* 18b ⟩ ≡
```
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* 19a ⟩
◇

Fragment referenced in 5b.

⟨ *unpack ticcutils or timbl* 19a ⟩ ≡

```
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=../modules/usrlocal
  make
  make install
fi
cd /home/paul/dutch-nlp-modules-on-Lisa
rm -rf $ticbeldir
if
  [ $SUCCES -eq 0 ]
then
  ⟨ logmess (19b Installed $DIR ) 19e ⟩
else
  ⟨ logmess (19c NOT installed $DIR ) 19e ⟩
fi
◇
```

Fragment referenced in 18ab.

## 3.4   Logging

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

⟨ *variables of install-modules* 19d ⟩ ≡

```
LOGLEVEL=1
◇
```

Fragment referenced in 5b.

⟨ *logmess* 19e ⟩ ≡

```
if
  [ $LOGLEVEL -gt 0 ]
then
 echo @1
fi
◇
```

Fragment referenced in 6c, 7a, 8b, 19a, 20a.

### 3.5   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* 20a ⟩ ≡

```
SUCCES=0
cd /home/paul/dutch-nlp-modules-on-Lisa/modules
⟨ move module (20b $DIR ) 6a ⟩
wget $URL
SUCCES=$?
if
   [ $SUCCES -eq 0 ]
then
   tar -xzf $TARB
   SUCCES=$?
   rm -rf $TARB
fi
if
   [ $SUCCES -eq 0 ]
then
   ⟨ logmess (20c Installed $DIR ) 19e ⟩
   ⟨ remove old module (20d $DIR ) 6b ⟩
else
   ⟨ re-instate old module (20e $DIR ) 6c ⟩
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 macroÂ´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_dutch-nlp-modules-on-Lisa.w`. Figure 1 shows pathways to

Figure 1: 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.

translate it into printable/viewable documents and to extract the program sources. Table 4 lists

| 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

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

## A.3  Translate and run

This chapter assembles the Makefile for this project.

```
"Makefile" 21≡
    ⟨ default target 22a ⟩

    ⟨ parameters in Makefile 22c, … ⟩

    ⟨ impliciete make regels 25a, … ⟩
    ⟨ expliciete make regels 22e, … ⟩
    ⟨ make targets 26a, … ⟩
    ◇
```

The default target of make is `all`.

⟨ *default target* 22a ⟩ ≡
```
        all : ⟨ all targets 22b ⟩
        .PHONY : all
```

        ◇
Fragment referenced in 21.
Defines: all 11b, PHONY 25b.

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

⟨ *all targets* 22b ⟩ ≡
```
        dutch-nlp-modules-on-Lisa.pdf◇
```
Fragment referenced in 22a.
Uses: pdf 26a.

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* 22c ⟩ ≡
```
        .SUFFIXES: .pdf .w .tex .html .aux .log .php
```

        ◇
Fragment defined by 22cd, 24ab, 26d, 29b, 32a.
Fragment referenced in 21.
Defines: SUFFIXES Never used.
Uses: pdf 26a.

## 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.

⟨ *parameters in Makefile* 22d ⟩ ≡
```
        NUWEB=../bin/nuweb
```
        ◇
Fragment defined by 22cd, 24ab, 26d, 29b, 32a.
Fragment referenced in 21.
Defines: NUWEB 22e, 27abc, 28b, 30c, 31c, 32c.
Uses: nuweb 28b.

⟨ *expliciete make regels* 22e ⟩ ≡
```
        $(NUWEB): m4_arojroot/nuweb-1.58
                cd ../nuweb-1.58 && make nuweb
                cp ../nuweb-1.58/nuweb $(NUWEB)
```

        ◇
Fragment defined by 22e, 23abc, 25b, 27a, 29de, 30ab.
Fragment referenced in 21.
Uses: NUWEB 22d, nuweb 28b.

⟨ *expliciete make regels* 23a ⟩ ≡
```
../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 22e, 23abc, 25b, 27a, 29de, 30ab.
Fragment referenced in 21.
Uses: `nuweb` 28b.

## A.5   Pre-processing

To make usable things from the raw input `a_dutch-nlp-modules-on-Lisa.w`, do the following:

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

This results in a L^AT_EX 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" T_EX 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* 23b ⟩ ≡
```
m4_dutch-nlp-modules-on-Lisa.w : a_dutch-nlp-modules-on-Lisa.w
        gawk '{if(match($$0, "@%")) {printf("%s", substr($$0,1,RSTART-1))} else print}' a_dutch-nlp-m
          | gawk '{gsub(/[\\][\$$]/, "$$");print}'  > m4_dutch-nlp-modules-on-Lisa.w
```

◇

Fragment defined by 22e, 23abc, 25b, 27a, 29de, 30ab.
Fragment referenced in 21.
Uses: `print` 26a.

### A.5.2   Run the M4 pre-processor

⟨ *expliciete make regels* 23c ⟩ ≡
```
dutch-nlp-modules-on-Lisa.w : m4_dutch-nlp-modules-on-Lisa.w inst.m4
        m4 -P m4_dutch-nlp-modules-on-Lisa.w > dutch-nlp-modules-on-Lisa.w
```

◇

Fragment defined by 22e, 23abc, 25b, 27a, 29de, 30ab.
Fragment referenced in 21.

## 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* 24a ⟩ ≡
```
      FIGFILES=fileschema
```

         ◇

Fragment defined by 22cd, 24ab, 26d, 29b, 32a.
Fragment referenced in 21.
Defines: FIGFILES 24b, 29b.


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* 24b ⟩ ≡
```
      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 22cd, 24ab, 26d, 29b, 32a.
Fragment referenced in 21.
Defines: FIGFILENAMES Never used, PDFT_NAMES 26b, PDF_FIG_NAMES 26b, PST_NAMES Never used,
        PS_FIG_NAMES Never used.
Uses: FIGFILES 24a.


Create the graph files with program `fig2dev`:

⟨ *impliciete make regels* 25a ⟩ ≡
```
      %.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 25a, 26b, 29c.
Fragment referenced in 21.
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 `dutch-nlp-modules-on-Lisa.bib`. To create this file, copy the auxiliary file to another file `auxfil.aux`, but replace the argument of the command `\bibdata{dutch-nlp-modules-on-Lisa}` 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* 25b ⟩ ≡
```
      bibfile : dutch-nlp-modules-on-Lisa.aux /home/paul/bin/mkportbib
              /home/paul/bin/mkportbib dutch-nlp-modules-on-Lisa litprog

      .PHONY : bibfile
```
◇

Fragment defined by 22e, 23abc, 25b, 27a, 29de, 30ab.
Fragment referenced in 21.
Uses: `PHONY` 22a.

### A.6.3  Create a printable/viewable document

Make a PDF document for printing and viewing.

⟨ *make targets* 26a ⟩ ≡
```
pdf : dutch-nlp-modules-on-Lisa.pdf

print : dutch-nlp-modules-on-Lisa.pdf
        lpr dutch-nlp-modules-on-Lisa.pdf

view : dutch-nlp-modules-on-Lisa.pdf
        evince dutch-nlp-modules-on-Lisa.pdf
```
⋄

Fragment defined by 26a, 29a, 32bc.
Fragment referenced in 21.
Defines: `pdf` 22bc, 26b, `print` 9d, 14a, 23b, `view` Never used.


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


*The w2pdf script*    The three processors nuweb, LATEX and bibTEX are intertwined. LATEX and
bibTEX 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 LATEX 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`.

Note, that in the following `make` construct, the implicit rule `.w.pdf` is not used. It turned out,
that make did not calculate the dependencies correctly when I did use this rule.

⟨ *impliciete make regels* 26b ⟩ ≡
```
%.pdf : %.w $(W2PDF)  $(PDF_FIG_NAMES) $(PDFT_NAMES)
        chmod 775 $(W2PDF)
        $(W2PDF) $*
```
⋄

Fragment defined by 25a, 26b, 29c.
Fragment referenced in 21.
Uses: `pdf` 26a, `PDFT_NAMES` 24b, `PDF_FIG_NAMES` 24b.


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* 26c ⟩ ≡
```
../nuweb/bin ⋄
```
Fragment defined by 4abcdef, 9c, 26c.
Fragment referenced in 32b.
Uses: `nuweb` 28b.


⟨ *parameters in Makefile* 26d ⟩ ≡
```
W2PDF=../nuweb/bin/w2pdf
```
⋄

Fragment defined by 22cd, 24ab, 26d, 29b, 32a.
Fragment referenced in 21.
Uses: `nuweb` 28b.

⟨ *expliciete make regels* 27a ⟩ ≡
```
      $(W2PDF) : dutch-nlp-modules-on-Lisa.w $(NUWEB)
             $(NUWEB) dutch-nlp-modules-on-Lisa.w
```
    ◇

Fragment defined by 22e, 23abc, 25b, 27a, 29de, 30ab.
Fragment referenced in 21.
Uses: NUWEB 22d.

"../nuweb/bin/w2pdf" 27b≡
```
      #!/bin/bash
      # w2pdf -- compile a nuweb file
      # usage: w2pdf [filename]
      # 20141103 at 0926h: Generated by nuweb from a_dutch-nlp-modules-on-Lisa.w
      NUWEB=/usr/local/bin/nuweb

      LATEXCOMPILER=pdflatex
```
      ⟨ *filenames in nuweb compile script* 27d ⟩
      ⟨ *compile nuweb* 27c ⟩

    ◇

Uses: NUWEB 22d, nuweb 28b.

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

⟨ *compile nuweb* 27c ⟩ ≡
```
      NUWEB=/usr/local/bin/nuweb
```
      ⟨ *run the processors until the aux file remains unchanged* 28c ⟩
      ⟨ *remove the copy of the aux file* 28a ⟩
    ◇

Fragment referenced in 27b.
Uses: NUWEB 22d, nuweb 28b.

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 LaTeX 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* 27d ⟩ ≡
```
      nufil=$1
      trunk=${1%%.*}
      texfil=${trunk}.tex
      auxfil=${trunk}.aux
      oldaux=old.${trunk}.aux
      indexfil=${trunk}.idx
      oldindexfil=old.${trunk}.idx
```
    ◇

Fragment referenced in 27b.
Defines: auxfil 28c, 31ab, indexfil 28c, 31a, nufil 28b, 31ac, oldaux 28ac, 31ab, oldindexfil 28c, 31a,
      texfil 28b, 31ac, trunk 28b, 31acd.

Remove the old copy if it is no longer needed.

⟨ *remove the copy of the aux file* 28a ⟩ ≡

```
rm $oldaux
```
    ◇

Fragment referenced in 27c, 30d.
Uses: `oldaux` 27d, 31a.

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

⟨ *run the three processors* 28b ⟩ ≡

```
$NUWEB $nufil
$LATEXCOMPILER $texfil
makeindex $trunk
bibtex $trunk
```
    ◇

Fragment referenced in 28c.
Defines: `bibtex` 31cd, `makeindex` 31cd, `nuweb` 22de, 23a, 26cd, 27bc, 29a, 30bc.
Uses: `nufil` 27d, 31a, `NUWEB` 22d, `texfil` 27d, 31a, `trunk` 27d, 31a.

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* 28c ⟩ ≡

```
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 28b ⟩
  if [ $LOOPCOUNTER -ge 10 ]
  then
    cp $auxfil $oldaux
  fi;
done
```
    ◇

Fragment referenced in 27c.
Uses: `auxfil` 27d, 31a, `indexfil` 27d, `oldaux` 27d, 31a, `oldindexfil` 27d.

### 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`.

Nuweb creates a LATEX file that is suitable for `latex2html` if the source file has `.hw` as suffix instead of `.w`. However, this feature is not compatible with tex4ht.

Make html file:

⟨ *make targets* 29a ⟩ ≡
```
html : ../nuweb/html/dutch-nlp-modules-on-Lisa.html
```
          ◇

Fragment defined by 26a, 29a, 32bc.
Fragment referenced in 21.
Uses: `nuweb` 28b.


The HTML file depends on its source file and the graphics files.

Make lists of the graphics files and copy them.

⟨ *parameters in Makefile* 29b ⟩ ≡
```
HTML_PS_FIG_NAMES=$(foreach fil,$(FIGFILES), m4_htmldocdir/$(fil).pstex)
HTML_PST_NAMES=$(foreach fil,$(FIGFILES), m4_htmldocdir/$(fil).pstex_t)
```
          ◇

Fragment defined by 22cd, 24ab, 26d, 29b, 32a.
Fragment referenced in 21.
Uses: `FIGFILES` 24a.


⟨ *impliciete make regels* 29c ⟩ ≡
```
m4_htmldocdir/%.pstex : %.pstex
        cp  $< $@

m4_htmldocdir/%.pstex_t : %.pstex_t
        cp  $< $@
```
          ◇

Fragment defined by 25a, 26b, 29c.
Fragment referenced in 21.


Copy the nuweb file into the html directory.

⟨ *expliciete make regels* 29d ⟩ ≡
```
m4_htmlsource : dutch-nlp-modules-on-Lisa.w
        cp  dutch-nlp-modules-on-Lisa.w m4_htmlsource
```
          ◇

Fragment defined by 22e, 23abc, 25b, 27a, 29de, 30ab.
Fragment referenced in 21.


We also need a file with the same name as the documentstyle and suffix `.4ht`. Just copy the file `report.4ht` from the tex4ht distribution. Currently this seems to work.

⟨ *expliciete make regels* 29e ⟩ ≡
```
m4_4htfildest : m4_4htfilsource
        cp m4_4htfilsource m4_4htfildest
```
          ◇

Fragment defined by 22e, 23abc, 25b, 27a, 29de, 30ab.
Fragment referenced in 21.

Copy the bibliography.

⟨ *expliciete make regels* 30a ⟩ ≡
```
      m4_htmlbibfil : m4_nuwebdir/dutch-nlp-modules-on-Lisa.bib
              cp m4_nuwebdir/dutch-nlp-modules-on-Lisa.bib m4_htmlbibfil
```

         ◇
Fragment defined by 22e, 23abc, 25b, 27a, 29de, 30ab.
Fragment referenced in 21.

Make a dvi file with `w2html` and then run `htlatex`.

⟨ *expliciete make regels* 30b ⟩ ≡
```
      ../nuweb/html/dutch-nlp-modules-on-Lisa.html : m4_htmlsource m4_4htfildest $(HTML_PS_FIG_NAMES) $(HTM
              cp w2html ../bin
              cd ../bin && chmod 775 w2html
              cd m4_htmldocdir && ../bin/w2html dutch-nlp-modules-on-Lisa.w
```

         ◇
Fragment defined by 22e, 23abc, 25b, 27a, 29de, 30ab.
Fragment referenced in 21.
Uses: nuweb 28b.

Create a script that performs the translation.

"w2html" 30c≡
```
      #!/bin/bash
      # w2html -- make a html file from a nuweb file
      # usage: w2html [filename]
      #  [filename]: Name of the nuweb source file.
      '#' m4_header
      echo "translate " $1 >w2html.log
      NUWEB=/usr/local/bin/nuweb
```

      ⟨ *filenames in w2html* 31a ⟩

      ⟨ *perform the task of w2html* 30d ⟩

         ◇
Uses: NUWEB 22d, nuweb 28b.

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* 30d ⟩ ≡
      ⟨ *run the html processors until the aux file remains unchanged* 31b ⟩
      ⟨ *remove the copy of the aux file* 28a ⟩
         ◇
Fragment referenced in 30c.

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 LaTeX 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* 31a ⟩ ≡
```
nufil=$1
trunk=${1%%.*}
texfil=${trunk}.tex
auxfil=${trunk}.aux
oldaux=old.${trunk}.aux
indexfil=${trunk}.idx
oldindexfil=old.${trunk}.idx
```
          ◇

Fragment referenced in 30c.
Defines: auxfil 27d, 28c, 31b, nufil 27d, 28b, 31c, oldaux 27d, 28ac, 31b, texfil 27d, 28b, 31c, trunk 27d, 28b,
          31cd.
Uses: indexfil 27d, oldindexfil 27d.

⟨ *run the html processors until the aux file remains unchanged* 31b ⟩ ≡
```
while
  ! cmp -s $auxfil $oldaux
do
  if [ -e $auxfil ]
  then
   cp $auxfil $oldaux
  fi
```
  ⟨ *run the html processors* 31c ⟩
```
done
```
  ⟨ *run tex4ht* 31d ⟩

          ◇

Fragment referenced in 30d.
Uses: auxfil 27d, 31a, oldaux 27d, 31a.

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

⟨ *run the html processors* 31c ⟩ ≡
```
$NUWEB -o -n $nufil
latex $texfil
makeindex $trunk
bibtex $trunk
htlatex $trunk
```
          ◇

Fragment referenced in 31b.
Uses: bibtex 28b, makeindex 28b, nufil 27d, 31a, NUWEB 22d, texfil 27d, 31a, trunk 27d, 31a.

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* 31d ⟩ ≡
```
tex '\def\filename{{dutch-nlp-modules-on-Lisa}{idx}{4dx}{ind}} \input idxmake.4ht'
makeindex -o $trunk.ind $trunk.4dx
bibtex $trunk
htlatex $trunk
```
          ◇

Fragment referenced in 31b.
Uses: bibtex 28b, makeindex 28b, trunk 27d, 31a.

*create the program sources*   Run nuweb, but suppress the creation of the LaTeX 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* 32a ⟩ ≡
```
      MKDIR = mkdir -p
```

         ◇

Fragment defined by 22cd, 24ab, 26d, 29b, 32a.
Fragment referenced in 21.
Defines: MKDIR 32b.

⟨ *make targets* 32b ⟩ ≡
```
      DIRS = ⟨ directories to create 4a, … ⟩

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

         ◇

Fragment defined by 26a, 29a, 32bc.
Fragment referenced in 21.
Defines: DIRS 32c.
Uses: MKDIR 32a.

⟨ *make targets* 32c ⟩ ≡
```
      sources : dutch-nlp-modules-on-Lisa.w $(DIRS) $(NUWEB)
              $(NUWEB) dutch-nlp-modules-on-Lisa.w
              ⟨ make scripts executable 5c, … ⟩
```

         ◇

Fragment defined by 26a, 29a, 32bc.
Fragment referenced in 21.
Uses: DIRS 32b, NUWEB 22d.

# B   References

## B.1   Literature

## 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_velocitylayoutservletURL
**Jetty:** m4_jettycodehausURL
**UserBase javadoc:** m4_userbasejavadocURL
**VU corpus Management development site:** http://code.google.com/p/vucom

# C    Indexes

## C.1    Filenames

`"../bin/alpinohack"` Defined by 10a.
`"../bin/heideltime"` Defined by 14b.
`"../bin/install-modules"` Defined by 5b.
`"../bin/mor"` Defined by 9a.
`"../bin/ned"` Defined by 12b.
`"../bin/nerc"` Defined by 10d.
`"../bin/onto"` Defined by 13a.
`"../bin/srl"` Defined by 15abcde.
`"../bin/test"` Defined by 16e.
`"../bin/tok"` Defined by 7l.
`"../bin/wsd"` Defined by 11b.
`"../modules/alpinohack/clean_hack.py"` Defined by 9d.
`"../nuweb/bin/w2pdf"` Defined by 27b.
`"Makefile"` Defined by 21.
`"w2html"` Defined by 30c.

## C.2    Macro's

⟨ adapt heideltime's config.props 14a ⟩ Referenced in 13c.
⟨ all targets 22b ⟩ Referenced in 22a.
⟨ compile nuweb 27c ⟩ Referenced in 27b.
⟨ default target 22a ⟩ Referenced in 21.
⟨ directories to create 4abcdef, 9c, 26c ⟩ Referenced in 32b.
⟨ expliciete make regels 22e, 23abc, 25b, 27a, 29de, 30ab ⟩ Referenced in 21.
⟨ filenames in nuweb compile script 27d ⟩ Referenced in 27b.
⟨ filenames in w2html 31a ⟩ Referenced in 30c.
⟨ find leave and tree 7g ⟩ Referenced in 7a.
⟨ impliciete make regels 25a, 26b, 29c ⟩ Referenced in 21.
⟨ install Alpino 8b ⟩ Referenced in 5b.
⟨ install from github 7a ⟩ Referenced in 7h, 8h, 13c, 14d, 16a.
⟨ install from tarball 20a ⟩ Not referenced.
⟨ install kafnafparserpy 16a ⟩ Referenced in 5b.
⟨ install the heideltime module 13c ⟩ Referenced in 5b.
⟨ install the morphosyntactic parser 8h ⟩ Referenced in 5b.
⟨ install the NERC module 10c ⟩ Referenced in 5b.
⟨ install the onto module 12d ⟩ Referenced in 5b.
⟨ install the srl module 14d ⟩ Referenced in 5b.
⟨ install the ticcutils utility 18a ⟩ Referenced in 5b.
⟨ install the timbl utility 18b ⟩ Referenced in 5b.
⟨ install the tokenizer 7h ⟩ Referenced in 5b.
⟨ install the treetagger utility 16g, 17abcde ⟩ Referenced in 5b.
⟨ install the WSD module 11a ⟩ Referenced in 5b.
⟨ install the NED module 12a ⟩ Referenced in 5b.
⟨ logmess 19e ⟩ Referenced in 6c, 7a, 8b, 19a, 20a.
⟨ make scripts executable 5c, 8a, 9b, 10be, 11c, 12c, 13b, 14c, 15f, 16f ⟩ Referenced in 32c.
⟨ make targets 26a, 29a, 32bc ⟩ Referenced in 21.
⟨ move module 6a ⟩ Referenced in 7a, 8b, 20a.
⟨ parameters in Makefile 22cd, 24ab, 26d, 29b, 32a ⟩ Referenced in 21.
⟨ perform the task of w2html 30d ⟩ Referenced in 30c.
⟨ re-instate old module 6c ⟩ Referenced in 7a, 8b, 20a.
⟨ remove old module 6b ⟩ Referenced in 7a, 8b, 20a.
⟨ remove the copy of the aux file 28a ⟩ Referenced in 27c, 30d.
⟨ run tex4ht 31d ⟩ Referenced in 31b.
⟨ run the html processors 31c ⟩ Referenced in 31b.
⟨ run the html processors until the aux file remains unchanged 31b ⟩ Referenced in 30d.

⟨run the processors until the aux file remains unchanged 28c⟩ Referenced in 27c.
⟨run the three processors 28b⟩ Referenced in 28c.
⟨set alpinohome 8g⟩ Referenced in 9a.
⟨set local bin directory 4h⟩ Referenced in 15a.
⟨set pythonpath 4g⟩ Referenced in 9a, 14b, 15a.
⟨unpack snapshots or die 5a⟩ Referenced in 5b.
⟨unpack ticcutils or timbl 19a⟩ Referenced in 18ab.
⟨variables of install-modules 19d⟩ Referenced in 5b.

## C.3   Variables

`all`: 11b, 22a.
`ALPINO_HOME`: 8g.
`auxfil`: 27d, 28c, 31a, 31b.
`bibtex`: 28b, 31cd.
`DIRS`: 32b, 32c.
`fig2dev`: 25a.
`FIGFILENAMES`: 24b.
`FIGFILES`: 24a, 24b, 29b.
`indexfil`: 27d, 28c, 31a.
`makeindex`: 28b, 31cd.
`MKDIR`: 32a, 32b.
`nufil`: 27d, 28b, 31a, 31c.
`NUWEB`: 22d, 22e, 27abc, 28b, 30c, 31c, 32c.
`nuweb`: 22de, 23a, 26cd, 27bc, 28b, 29a, 30bc.
`oldaux`: 27d, 28ac, 31a, 31b.
`oldindexfil`: 27d, 28c, 31a.
`pdf`: 22bc, 26a, 26b.
`PDFT_NAMES`: 24b, 26b.
`PDF_FIG_NAMES`: 24b, 26b.
`PHONY`: 22a, 25b.
`print`: 9d, 14a, 23b, 26a.
`PST_NAMES`: 24b.
`PS_FIG_NAMES`: 24b.
`SUCCES`: 8b, 19a, 20a.
`SUFFIXES`: 22c.
`texfil`: 27d, 28b, 31a, 31c.
`trunk`: 27d, 28b, 31a, 31cd.
`view`: 26a.