

Scraper example

Paul Huygen <paul.huygen@huygen.nl>

20th September 2016
16:30 h.

Abstract

In this document a web-scraper is constructed that scrapes the forum <http://web.archive.org/web/20160323073042/http://ragingbull.com>, using Python and BeautifulSoup.

Contents

1	<i>Introduction</i>	2
1.1	Structure of the forum	2
1.2	What are we going to do?	2
1.3	Metadata	2
2	<i>The program</i>	2
2.1	Read the command-line	2
2.2	BeautifulSoup	3
2.3	Make soup from an URL	3
2.4	Extract the topic-title from a topic page	3
2.5	Extract the topics from a board page	4
2.6	Extract the posts from a topic page	5
2.7	Generate the NAF file	6
2.8	Remove mark-up from the text	8
2.9	Scrape a board	10
2.10	The program file	11
A	<i>How to read and translate this document</i>	11
A.1	Read this document	11
A.2	Process the document	12
A.3	Translate and run	13
A.4	Pre-processing	13
A.4.1	Process ‘dollar’ characters	14
A.4.2	Run the M4 pre-processor	14
A.5	Typeset this document	14
A.5.1	Figures	14
A.5.2	Bibliography	16
A.5.3	Create a printable/viewable document	16
A.5.4	Create HTML files	19
B	<i>References</i>	23
B.1	Literature	23
B.2	URL’s	23
C	<i>Indexes</i>	23
C.1	Filenames	23
C.2	Macro’s	24
C.3	Variables	24

1 Introduction

- Scrape a forum on a website.
- In this case <http://web.archive.org/web/20160323073042/http://ragingbull.com>.
- Use Python and BeautifulSoup.

1.1 Structure of the forum

The forum consists of a set of *boards* with different subjects. Each board has an identifying number and a name, e.g. board 14242 is about *Current Events*, abbreviated as CEVT. The main page of that board has as URL: <http://web.archive.org/web/20160323073042/http://ragingbull.com/board/14242>. It contains a table with a list of topics and, when there are too many topics for a single page, references to other URL's that contain lists of older topics. These URL's look like <http://web.archive.org/web/20160323073042/http://ragingbull.com/board/14242/page/2>.

A topic has as url e.g. <http://web.archive.org/web/20160323073042/http://ragingbull.com/topic/1061702> and a title. The page of the topic contains a list of posts.

1.2 What are we going to do?

1. Read the pages of the board and collect the url's of the topics
2. Read the pages of the topics and extract the posts.
3. Wrap each post (text and metadata) in a NAF file.

1.3 Metadata

We need to collect for each post the following metadata:

1. board name and ID.
2. Topic name and ID.
3. Sequence number of the post in the topic.
4. Author ID.
5. Date of the post.

To test whether we have gathered a post with the correct metadata, we can print it as follows:

```
<methods of the main program 2> ≡
def print_post(board_id, board_name, topic, seq, author, post_date, text):
    print( "Board:   {} ({}).format(board_id, board_name))
    print( "Topic:   {}".format(topic))
    print( "Post nr: {}".format(seq))
    print( "Date:     {}".format(post_date))
    print( "Text: {} ".format(text))
```

◇

Fragment defined by 2, 3cd, 4ab, 5b, 6c, 7d, 10abc.

Fragment referenced in 11b.

Defines: print_post 12.

Uses: print 17a.

2 The program

2.1 Read the command-line

In this demo-phase we parse the board “Oil and Natural Gas Investments” (board number 11677). Scrape the Wayback archive of this board (URL: <http://web.archive.org/web/20160323073042/http://ragingbull.com/forum/board/11677>).

⟨ get program options 3a ⟩ ≡

```
boardURL = 'http://web.archive.org/web/20160323073042/http://ragingbull.com/forum/board/11677'
boardDIR = str(11677)
```

◇

Fragment referenced in 11b.
Defines: boardURL 10b, 11b.

2.2 BeautifulSoup

We will use Python's BeautifulSoup module to extract the posts from the forum.

⟨ import modules in main program 3b ⟩ ≡

```
from bs4 import BeautifulSoup
import requests
```

◇

Fragment defined by 3b, 6b, 7c, 9.
Fragment referenced in 11b.
Defines: BeautifulSoup 3c, 10b, bs4 Never used, requests 3c, 10b.

2.3 Make soup from an URL

⟨ methods of the main program 3c ⟩ ≡

```
def make_soup_from_url(url):
    r = requests.get(url)
    soup = None
    if r.status_code == 200:
        soup = BeautifulSoup(r.content, 'lxml')
        print("{}: {}".format(r.status_code, url))
    return soup
```

◇

Fragment defined by 2, 3cd, 4ab, 5b, 6c, 7d, 10abc.
Fragment referenced in 11b.
Defines: make_soup_from_url 4a, 5b.
Uses: BeautifulSoup 3b, print 17a, requests 3b.

2.4 Extract the topic-title from a topic page

The title of the topic can be found as the contents of the “title” tag inside the “head” section of the html document:

⟨ methods of the main program 3d ⟩ ≡

```
def get_topic(soup):
    headpart = soup.head
    title = headpart.title.string
    return title
```

◇

Fragment defined by 2, 3cd, 4ab, 5b, 6c, 7d, 10abc.
Fragment referenced in 11b.

2.5 Extract the topics from a board page

A board page contains the data to find the topics that belong to the board. Often the board page has sequel pages with older topics. Sequel pages can be found by appending `/page/<nn>` to the URL of the board page.

The following (recursive) method yields a list of the URL's, the ID's and the titles of the topics in a board page.

```

< methods of the main program 4a > ≡
def next_topic(base_url, pagenumber = 1):
    if pagenumber > 1:
        board_url = "{}/page/{}".format(base_url, pagenumber)
    else:
        board_url = base_url
    soup = make_soup_from_url(board_url)
    if soup == None:
        return
    else:
        < yield topic data from soup 5a >
        pagenumber += 1
        for tdata in next_topic(base_url, pagenumber):
            yield tdata

```

◇

Fragment defined by 2, 3cd, 4ab, 5b, 6c, 7d, 10abc.

Fragment referenced in 11b.

Defines: `next_topic` 11b.

Uses: `make_soup_from_url` 3c.

A board web-page hides the data of a topic in an anchor in a table with class attribute `topics`. So, let us go down to the body of the page, find the table and the anchors.

The method `is_topic_table` determine whether a tag found in the page is the table with the topics. The method `is_topicanchor` does a similar thing to find the anchor that leads to the topic page.

```

< methods of the main program 4b > ≡
def is_topictable(tag):
    if tag.name == 'table':
        if tag.has_attr('class'):
            return tag['class'][0] == 'topics'
    return False

def is_topicanchor(tag):
    if tag.name == 'a':
        if tag.has_attr('class'):
            return tag['class'][0] == 'topic-name'
    return False

```

◇

Fragment defined by 2, 3cd, 4ab, 5b, 6c, 7d, 10abc.

Fragment referenced in 11b.

Find the anchors.

```

<yield topic data from soup 5a> ≡
    sbody = soup.body
    topictabletag = sbody.find(is_topictable)
    for topicanchor in topictabletag.find_all(is_topicanchor):
        url = 'http://web.archive.org' + topicanchor['href']
        m = re.search(topicpattern, topicanchor['href'])
        id = m.group(1)
        title = topicanchor['title'].strip()
        yield [url, id, title]

```

◇

Fragment referenced in 4a.

Uses: sbody 6a.

2.6 Extract the posts from a topic page

A topic page contains a number of posts, wrapped in `<article>/</article>` tags.

When there are many posts in a topic, there will be subsequent pages with posts. We can just try to find such pages (with `/page/<n>`) suffix until we get a “400” result, or we can scrape URL’s.

Between the `<article>` and `</article>` tags we can find:

Post-id: as argument “id” in the `article` tag.

Author name: In a tag “header”, in a div “author-and-time”, in an anchor of class “author-name”.

When we pass the URL to the following function `next_article`, it will yield the texts and metadata of the articles:

```

<methods of the main program 5b> ≡

def next_article(topic_base_url, pagenumber = 1):
    if pagenumber > 1:
        topic_url = "{}/page/{}".format(topic_base_url, pagenumber)
    else:
        topic_url = topic_base_url
    soup = make_soup_from_url(topic_url)
    if soup == None:
        return
    else:
        <yield data from articles in this soup 6a>
        pagenumber += 1
        for tdata in next_article(topic_base_url, pagenumber):
            yield tdata

```

◇

Fragment defined by 2, 3cd, 4ab, 5b, 6c, 7d, 10abc.

Fragment referenced in 11b.

Defines: `nextarticle` Never used.

Uses: `make_soup_from_url` 3c.

The posts of the topic can be found in `article` tags.

```

<yield data from articles in this soup 6a> ≡
sbody = soup.body
postnum = 0
for article in soup.find_all("article"):
    postnum += 1
    header = article.header
    for sp in header.find_all("span"):
        if sp['class'][0] == "postId":
            postid = sp.string
        elif sp['class'][0] == "time":
            posttime = sp.string
    for div in header.find_all("div"):
        if div['class'][0] == "author-and-time":
            for anchor in div.find_all("a"):
                if anchor['class'][0] == "author-name":
                    author=anchor.string
                    author_url = anchor.href
    if author == None:
        author = "Anonymus"
    text = article.textarea.string
    yield [ postid, posttime, postnum, author, author_url, text ]

```

◇

Fragment referenced in 5b.

Defines: postnum 11b, sbody 5a, 10c.

2.7 Generate the NAF file

Generate the NAF file with the [KafNafParserPy](#) package.

```

<import modules in main program 6b> ≡
import KafNafParserPy

```

◇

Fragment defined by 3b, 6b, 7c, 9.

Fragment referenced in 11b.

If you construct a NAF from scratch, it doesn't have a header section. To work around this, we read in a template of a NAF file that contains an empty header. Fill in the header, add a **raw** tag with the text of the post and write out to a file that is named after the ID of the post:

```

<methods of the main program 6c> ≡
def printnaf(nafpath, topic, author, post_date, text):
    naf = KafNafParserPy.KafNafParser(filename = 'template.naf')
    naf.set_language("en")
    outtext = Contents_block(text)
    naf.set_raw(outtext.without_bbcode())
    <create the naf header 7b>
    print("To write naf in {}".format(nafpath))
    naf.dump(filename = nafpath)
    print("Wrote {}".format(nafpath))

```

◇

Fragment defined by 2, 3cd, 4ab, 5b, 6c, 7d, 10abc.

Fragment referenced in 11b.

Defines: printnaf 11b.

Uses: print 17a.

```

"../template.naf" 7a≡
  <?xml version="1.0" encoding="UTF-8"?>
  <NAF>
    <nafHeader></nafHeader>
  </NAF>

  ◇

```

The following metadata goes in the NAF header:

- Topic
- Author
- Date of the post.

```

⟨ create the naf header 7b ⟩ ≡
  header = naf.get_header()
  fileDesc = KafNafParserPy.CfileDesc()
  header.set_fileDesc(fileDesc)
  fileDesc.set_title(topic)
  fileDesc.set_author(author)
  fileDesc.set_creationtime(convert_timestring(post_date))

  ◇

```

Fragment referenced in 6c.

Uses: `convert_timestring` 7d.

Find the time of the post. Sometimes the time-stamp is a string like 2013-09-09 16:04, but in other instances it is expressed like Mar 22 22:48. We must find out what kind of string it is and then convert the time-stamp to the ISO 8601 format. It turns out that the `python-dateutil` parser can read in both formats. So:

```

⟨ import modules in main program 7c ⟩ ≡
  import dateutil.parser

  ◇

```

Fragment defined by 3b, 6b, 7c, 9.

Fragment referenced in 11b.

Defines: `dateutil.parser` 7d.

```

⟨ methods of the main program 7d ⟩ ≡
  def convert_timestring(post_string):
    pubtime = dateutil.parser.parse(post_string)
    return pubtime.isoformat()

  ◇

```

Fragment defined by 2, 3cd, 4ab, 5b, 6c, 7d, 10abc.

Fragment referenced in 11b.

Defines: `convert_timestring` 7b.

Uses: `dateutil.parser` 7c.

To convert month-names (e.g. “Jan”) to month-numbers (e.g. 1), use the following dictionary.

(variables of the main program 8) \equiv
`monthnums = {v: k for k,v in enumerate(calendar.month_abbrev)}`
 \diamond

Fragment defined by [8](#), [11a](#).
 Fragment referenced in [11b](#).
 Defines: `monthnums` Never used.
 Uses: `calendar` [9](#).

(import modules in main program 9) \equiv
`import datetime`
`import calendar`
 \diamond

Fragment defined by [3b](#), [6b](#), [7c](#), [9](#).
 Fragment referenced in [11b](#).
 Defines: `calendar` [8](#), `datetime` [12](#).

2.8 Remove mark-up from the text

The HTML pages of Ragingbull contain the text of the posts as HTML code or as “bb-code”. A concise guide for bb-code can be found [here](#).

tag	description	action
[b], [/b]:	boldface	remove mark-up
[i], [/i]:	italic	remove mark-up
[u], [/u]:	underline	remove mark-up
[s], [/s]:	strike-through	remove tag
[color], [/color]:	back-ground color	remove mark-up
[center], [/center]:	centered text	remove mark-up
[quote], [/quote]:	quotation	Add quotation marks
[quote={name}], [/quote]:	quotation	<code>name</code> said: ‘ ‘ ... ’ ’
[url], [/url]:	Link	remove mark-up
[url={url}], [/url]:	Link	Leave the text.
[img ...], [/img]:	image	replace by “image”
[ul], [/ul]:	Unordered list	remove mark-up
[ol], [/ol]:	ordered list	remove mark-up
[list], [/list]:	list	remove mark-up
[li], [/li]:	list item	
[code], [/code]:	Verbatim	
[table], [/table]:	table	
[tr], [/tr]:	table row	
[th], [/th]:	table heading	
[td], [/td]:	table cell	
[youtube], [/youtube]:	URL to Youtube	remove mark-up
[gvideo], [/gvideo]:	URL to video	remove mark-up

(*methods of the main program 10a*) \equiv

```
class Contents_block:
    def __init__(self,intext):
        self.intext = intext

    def _strip_bbttag(self, intext, tagname):
        pattern = re.compile(r'\[' + tagname + r'\](.*)\[/' + tagname + ' \]\')
        return re.sub(pattern, r'\1', intext)

    def _strip_bbttagged_substring(self, intext, tagname):
        pattern = re.compile(r'\[' + tagname + r'\](.*)\[/' + tagname + r'\]\')
        return re.sub(pattern, '', intext)

    def _replace_bbttagged_substring(self, intext, tagname, repl):
        pattern = re.compile(r'\[' + tagname + r'\](.*)\[/' + tagname + ' \]\')
        return re.sub(pattern, repl, intext)

    def _unquote(self, intext):
        out = self._strip_bbttag(intext, 'quote')
        pattern = re.compile(r'\[quote=(\[^\]]*\)\](.*)\[\/quote\]\')
        out = re.sub(pattern, r'\1 said: "\2"', out)
        return out

    def _un_url(self, intext):
        pattern = re.compile(r'\[url\](.*)\[\/url\]\')
        out = re.sub(pattern, r'\1', intext)
        pattern = re.compile(r'\[url=(\[^\]]*\)\](.*)\[\/url\]\')
        out = re.sub(pattern, r'\2' + r' (' + r'\1' + r')', intext)
        return out

    def without_bbcode(self):
        out = self._strip_bbttag(self.intext, 'b')
        out = self._strip_bbttag(out, 'i')
        out = self._strip_bbttag(out, 'u')
        out = self._strip_bbttag(out, 'color')
        out = self._strip_bbttag(out, 'youtube')
        out = self._strip_bbttag(out, 'gvideo')
        out = self._strip_bbttagged_substring(out, 's')
        out = self._strip_bbttagged_substring(out, 'img')
        out = self._unquote(out)
        out = self._un_url(out)
        return out
```

◇

Fragment defined by 2, 3cd, 4ab, 5b, 6c, 7d, 10abc.

Fragment referenced in 11b.

2.9 Scrape a board

< methods of the main program 10b > ≡

```
def get_boardsoup():
    r = requests.get(boardURL)
    if r.status_code != 200:
        print("Board page {}".format(boardURL))
        print("Http request result: {}".format(r.status_code))
        print("Error exit")
        sys.exit()
    soup = BeautifulSoup(r.content, 'lxml')
    return soup
```

◇

Fragment defined by 2, 3cd, 4ab, 5b, 6c, 7d, 10abc.

Fragment referenced in 11b.

Defines: `get_boardsoup` Never used.

Uses: `BeautifulSoup` 3b, `boardURL` 3a, `print` 17a, `requests` 3b.

< methods of the main program 10c > ≡

```
def topics(soup):
    sbody = soup.body
    topictabletag = sbody.find(is_topictable)
    for topicanchor in topictabletag.find_all(is_topicanchor):
        m = re.search(topicpattern, topicanchor['href'])
        title = topicanchor['title'].strip()
        yield ['http://web.archive.org' + topicanchor['href'], m.group(1), title]
```

◇

Fragment defined by 2, 3cd, 4ab, 5b, 6c, 7d, 10abc.

Fragment referenced in 11b.

Uses: `sbody` 6a.

< variables of the main program 11a > ≡

```
topicpattern = re.compile('.*(.)')
```

◇

Fragment defined by 8, 11a.

Fragment referenced in 11b.

2.10 The program file

"../scrape.py" 11b ≡

```

    < import modules in main program 3b, ... >
import sys
import os
import re
    < variables of the main program 8, ... >
    < methods of the main program 2, ... >

if __name__ == "__main__" :
    < get program options 3a >

    for [topic_url, topic_id, toptitle ] in next_topic(boardURL):
        topicDIR = str(boardDIR) + '/' + str(topic_id)
        print("{}: {}".format(topicDIR, toptitle))
        os.makedirs(topicDIR, exist_ok = True)
        for [postid, posttime, postnum, author, author_url, text] in next_article(topic_url):
            outpath = topicDIR + '/' + str(postid) + '.naf'
            printnaf(outpath, toptitle, author, posttime, text)

```

◇

Uses: boardURL 3a, next_topic 4a, postnum 6a, print 17a, printnaf 6c.

For now, the program just prints a mock-up of a post:

< print the testpost 12 > ≡

```

    print_post(boardnum, "CEVT", "Gallup: life got better", 1, "juddism", datetime.datetime.now(), "Come o

```

◇

Fragment never referenced.

Uses: datetime 9, print_post 2.

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 > \equiv

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 > \equiv

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_mysql.w`. Figure 1 shows pathways to translate it into

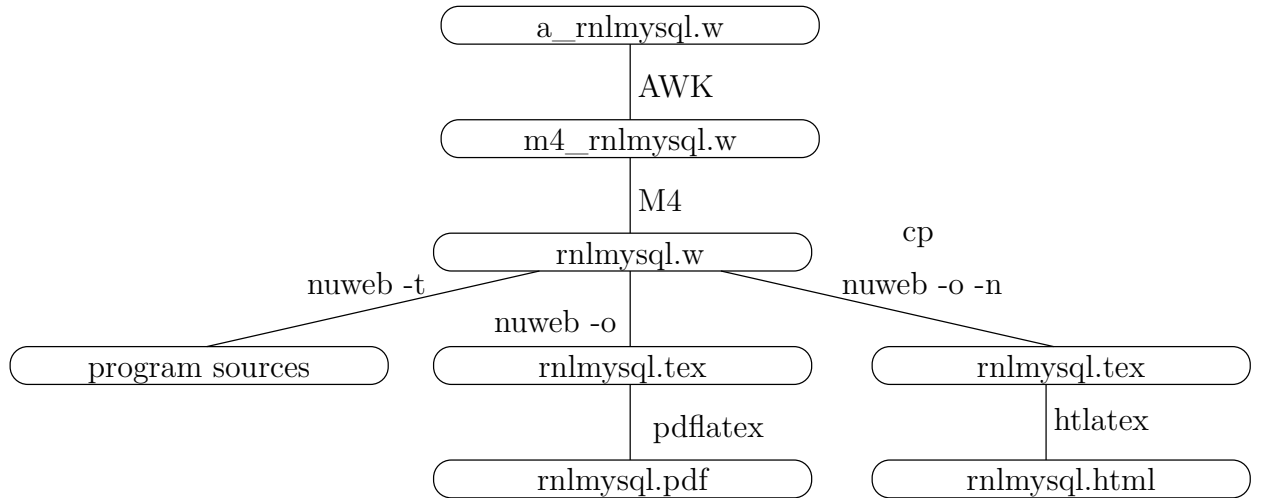


Figure 1: Translation of the raw code of this document into printable/viewable documents and into program sources. The figure shows the pathways and the main files involved.

printable/viewable documents and to extract the program sources. Table 1 lists the tools that are needed for a translation. Most of the tools (except Nuweb) are available on a well-equipped Linux system.

< parameters in Makefile 13a > \equiv

NUWEB=/usr/local/bin/nuweb

◇

Fragment defined by 13a, 14a, 15bc, 17c, 20a, 23a.

Fragment referenced in 13b.

Uses: nuweb 19a.

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

A.3 Translate and run

This chapter assembles the Makefile for this project.

```
"Makefile" 13b≡
  < default target 13c>

  < parameters in Makefile 13a, ... >

  < implicate make regels 16a, ... >
  < explicite make regels 14b, ... >
  < make targets 17a, ... >
  ◇
```

The default target of make is `all`.

```
< default target 13c> ≡
  all : < all targets 13d>
  .PHONY : all

  ◇
```

Fragment referenced in 13b.
Defines: `all` Never used, `PHONY` 16b.

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

```
< all targets 13d> ≡
  myscrapexamp.pdf◇
```

Fragment referenced in 13c.
Uses: `pdf` 17a.

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 14a> ≡
  .SUFFIXES: .pdf .w .tex .html .aux .log .php

  ◇
```

Fragment defined by 13a, 14a, 15bc, 17c, 20a, 23a.
Fragment referenced in 13b.
Defines: `SUFFIXES` Never used.
Uses: `pdf` 17a.

A.4 Pre-processing

To make usable things from the raw input `a_myscrapexamp.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.4.1 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 14b > ≡
    m4_myscrapexamp.w : a_myscrapexamp.w
                        gawk '{if(match($$0, "@%")) {printf("%s", substr($$0,1,RSTART-1))} else print}' a_myscrapexamp.w
                        | gawk '{gsub(/\[\[\][\$\$]/, "$$");print}' > m4_myscrapexamp.w
```

◇

Fragment defined by 14b, 15a, 16b, 17e, 20cde, 21a.

Fragment referenced in 13b.

Uses: print 17a.

A.4.2 Run the M4 pre-processor

```
< expliciete make regels 15a > ≡
    myscrapexamp.w : m4_myscrapexamp.w
                    m4 -P m4_myscrapexamp.w > myscrapexamp.w
```

◇

Fragment defined by 14b, 15a, 16b, 17e, 20cde, 21a.

Fragment referenced in 13b.

A.5 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.5.1 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 15b > ≡
 FIGFILES=fileschema

◇

Fragment defined by 13a, 14a, 15bc, 17c, 20a, 23a.
 Fragment referenced in 13b.
 Defines: FIGFILES 15c, 20a.

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 15c > ≡
 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 13a, 14a, 15bc, 17c, 20a, 23a.
 Fragment referenced in 13b.
 Defines: FIGFILENAMES Never used, PDFT_NAMES 17b, PDF_FIG_NAMES 17b, PST_NAMES Never used,
 PS_FIG_NAMES Never used.
 Uses: FIGFILES 15b.

Create the graph files with program `fig2dev`:

< impliciete make regels 16a > ≡
 %.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 16a, 17b, 20b.
 Fragment referenced in 13b.
 Defines: `fig2dev` Never used.

A.5.2 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 `myscrapexamp.bib`. To create this file, copy the auxiliary file to another file `auxfil.aux`, but replace the argument of the command `\bibdata{myscrapexamp}` 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 16b > ≡
    bibfile : myscrapexamp.aux /home/paul/bin/mkportbib
              /home/paul/bin/mkportbib myscrapexamp litprog

    .PHONY : bibfile
◇
```

Fragment defined by 14b, 15a, 16b, 17e, 20cde, 21a.

Fragment referenced in 13b.

Uses: PHONY 13c.

A.5.3 Create a printable/viewable document

Make a PDF document for printing and viewing.

```
< make targets 17a > ≡
    pdf : myscrapexamp.pdf

    print : myscrapexamp.pdf
            lpr myscrapexamp.pdf

    view : myscrapexamp.pdf
            evince myscrapexamp.pdf
◇
```

Fragment defined by 17a, 19c, 23b, ?.

Fragment referenced in 13b.

Defines: pdf 13d, 14a, 17b, print 2, 3c, 6c, 10b, 11b, 14b, 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.


```

< implicate make regels 17b > ≡
    %.pdf : %.w $(W2PDF) $(PDF_FIG_NAMES) $(PDFT_NAMES)
          chmod 775 $(W2PDF)
          $(W2PDF) $*

```

◇

Fragment defined by 16a, 17b, 20b.

Fragment referenced in 13b.

Uses: pdf 17a, PDFT_NAMES 15c, PDF_FIG_NAMES 15c.

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.

```

< parameters in Makefile 17c > ≡
    W2PDF=../nuweb/bin/w2pdf

```

◇

Fragment defined by 13a, 14a, 15bc, 17c, 20a, 23a.

Fragment referenced in 13b.

Uses: nuweb 19a.

```

< directories to create 17d > ≡
    ../nuweb/bin ◇

```

Fragment referenced in 23b.

Uses: nuweb 19a.

```

< expliciete make regels 17e > ≡
    $(W2PDF) : myscrapexamp.w
              $(NUWEB) myscrapexamp.w

```

◇

Fragment defined by 14b, 15a, 16b, 17e, 20cde, 21a.

Fragment referenced in 13b.

```

"../nuweb/bin/w2pdf" 18a≡
    #!/bin/bash
    # w2pdf -- compile a nuweb file
    # usage: w2pdf [filename]
    # 20160920 at 1630h: Generated by nuweb from a_myscrapexamp.w
    NUWEB=/usr/local/bin/nuweb
    LATEXCOMPILER=pdflatex
    < filenames in nuweb compile script 18c >
    < compile nuweb 18b >

```

◇

Uses: nuweb 19a.

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 18b ⟩ ≡
    NUWEB=m4_nuweb
    ⟨ run the processors until the aux file remains unchanged 19b ⟩
    ⟨ remove the copy of the aux file 18d ⟩
    ◇

```

Fragment referenced in 18a.

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 18c ⟩ ≡
    nufil=$1
    trunk=${1%.*}
    texfil=${trunk}.tex
    auxfil=${trunk}.aux
    oldaux=old.${trunk}.aux
    indexfil=${trunk}.idx
    oldindexfil=old.${trunk}.idx
    ◇

```

Fragment referenced in 18a.

Defines: `auxfil 19b, 22ab, indexfil 19b, 22a, nufil 19a, 22ac, oldaux 18d, 19b, 22ab, oldindexfil 19b, 22a, texfil 19a, 22ac, trunk 19a, 22acd.`

Remove the old copy if it is no longer needed.

```

⟨ remove the copy of the aux file 18d ⟩ ≡
    rm $oldaux
    ◇

```

Fragment referenced in 18b, 21c.

Uses: `oldaux 18c, 22a.`

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 19a ⟩ ≡
    $NUWEB $nufil
    $LATEXCOMPILER $texfil
    makeindex $trunk
    bibtex $trunk
    ◇

```

Fragment referenced in 19b.

Defines: `bibtex 22cd, makeindex 22cd, nuweb 13a, 17cd, 18a, 21b.`

Uses: `nufil 18c, 22a, texfil 18c, 22a, trunk 18c, 22a.`

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 19b ⟩ ≡
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 19a ⟩
  if [ $LOOPCOUNTER -ge 10 ]
  then
    cp $auxfil $oldaux
  fi;
done
◇

```

Fragment referenced in 18b.

Uses: auxfil 18c, 22a, indexfil 18c, oldaux 18c, 22a, oldindexfil 18c.

A.5.4 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 19c ⟩ ≡
html : m4_htmltarget
◇

```

Fragment defined by 17a, 19c, 23b, ?.

Fragment referenced in 13b.

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

Make lists of the graphics files and copy them.

```

⟨ parameters in Makefile 20a ⟩ ≡
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 13a, 14a, 15bc, 17c, 20a, 23a.

Fragment referenced in 13b.

Uses: FIGFILES 15b.

```

< implicate make regels 20b > ≡
    m4_htmlldocdir/%.pstex : %.pstex
        cp $< $@

    m4_htmlldocdir/%.pstex_t : %.pstex_t
        cp $< $@

```

◇

Fragment defined by 16a, 17b, 20b.

Fragment referenced in 13b.

Copy the nuweb file into the html directory.

```

< expliciete make regels 20c > ≡
    m4_htmlsource : myscrapexamp.w
        cp myscrapexamp.w m4_htmlsource

```

◇

Fragment defined by 14b, 15a, 16b, 17e, 20cde, 21a.

Fragment referenced in 13b.

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 20d > ≡
    m4_4htfildest : m4_4htfilsource
        cp m4_4htfilsource m4_4htfildest

```

◇

Fragment defined by 14b, 15a, 16b, 17e, 20cde, 21a.

Fragment referenced in 13b.

Copy the bibliography.

```

< expliciete make regels 20e > ≡
    m4_htmlbibfil : m4_anuwebdir/myscrapexamp.bib
        cp m4_anuwebdir/myscrapexamp.bib m4_htmlbibfil

```

◇

Fragment defined by 14b, 15a, 16b, 17e, 20cde, 21a.

Fragment referenced in 13b.

Make a dvi file with w2html and then run htlatex.

```

< expliciete make regels 21a > ≡

    m4_htmltarget : m4_htmlsource m4_4htfildest $(HTML_PS_FIG_NAMES) $(HTML_PST_NAMES) m4_htmlbibfil
        cp w2html /home/paul/projecten/cltl/emoeco/myscrapexamp/bin
        cd /home/paul/projecten/cltl/emoeco/myscrapexamp/bin && chmod 775 w2html
        cd m4_htmlldocdir && /home/paul/projecten/cltl/emoeco/myscrapexamp/bin/w2html myscrapexamp.w

```

◇

Fragment defined by 14b, 15a, 16b, 17e, 20cde, 21a.

Fragment referenced in 13b.

Create a script that performs the translation.

```
"w2html" 21b≡
#!/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 22a⟩

⟨perform the task of w2html 21c⟩

◇
```

Uses: **nuweb** 19a.

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 21c⟩ ≡
  ⟨run the html processors until the aux file remains unchanged 22b⟩
  ⟨remove the copy of the aux file 18d⟩
◇
```

Fragment referenced in 21b.

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 22a⟩ ≡
nufil=$1
trunk=${1%.*}
texfil=${trunk}.tex
auxfil=${trunk}.aux
oldaux=old.${trunk}.aux
indexfil=${trunk}.idx
oldindexfil=old.${trunk}.idx
◇
```

Fragment referenced in 21b.

Defines: **auxfil** 18c, 19b, 22b, **nufil** 18c, 19a, 22c, **oldaux** 18cd, 19b, 22b, **texfil** 18c, 19a, 22c, **trunk** 18c, 19a, 22cd.

Uses: **indexfil** 18c, **oldindexfil** 18c.

```

⟨run the html processors until the aux file remains unchanged 22b⟩ ≡
    while
        ! cmp -s $auxfil $oldaux
    do
        if [ -e $auxfil ]
        then
            cp $auxfil $oldaux
        fi
        ⟨run the html processors 22c⟩
    done
    ⟨run tex4ht 22d⟩

```

◇

Fragment referenced in 21c.

Uses: auxfil 18c, 22a, oldaux 18c, 22a.

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

```

⟨run the html processors 22c⟩ ≡
    $NUWEB -o -n $nufil
    latex $texfil
    makeindex $trunk
    bibtex $trunk
    htlatex $trunk

```

◇

Fragment referenced in 22b.

Uses: bibtex 19a, makeindex 19a, nufil 18c, 22a, texfil 18c, 22a, trunk 18c, 22a.

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 22d⟩ ≡
    tex '\def\filename{{myscrapexamp}{idx}{4dx}{ind}} \input idxmake.4ht'
    makeindex -o $trunk.ind $trunk.4dx
    bibtex $trunk
    htlatex $trunk

```

◇

Fragment referenced in 22b.

Uses: bibtex 19a, makeindex 19a, trunk 18c, 22a.

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 23a⟩ ≡
    MKDIR = mkdir -p

```

◇

Fragment defined by 13a, 14a, 15bc, 17c, 20a, 23a.

Fragment referenced in 13b.

Defines: MKDIR 23b.

$\langle \text{make targets 23b} \rangle \equiv$
 DIRS = $\langle \text{directories to create 17d} \rangle$

\$(DIRS) :
 \$(MKDIR) \$@

◇

Fragment defined by 17a, 19c, 23b, ?.

Fragment referenced in 13b.

Defines: DIRS ?.

Uses: MKDIR 23a.

$\langle \text{make targets ?} \rangle \equiv$
 sources : myscrapexamp.w \$(DIRS)
 \$(NUWEB) myscrapexamp.w

test : sources
 cd .. && python scrape.py

◇

Fragment defined by 17a, 19c, 23b, ?.

Fragment referenced in 13b.

Uses: DIRS 23b.

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_velocitylayout servletURL

Jetty: m4_jettycodehausURL

UserBase javadoc: m4_userbasejavadocURL

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

C Indexes

C.1 Filenames

"../nuweb/bin/w2pdf" Defined by 18a.

"../scrape.py" Defined by 11b.

"../template.naf" Defined by 7a.

"Makefile" Defined by 13b.

"w2html" Defined by 21b.

C.2 Macro's

<all targets 13d> Referenced in 13c.
 <compile nuweb 18b> Referenced in 18a.
 <create the naf header 7b> Referenced in 6c.
 <default target 13c> Referenced in 13b.
 <directories to create 17d> Referenced in 23b.
 <explicitete make regels 14b, 15a, 16b, 17e, 20cde, 21a> Referenced in 13b.
 <filenames in nuweb compile script 18c> Referenced in 18a.
 <filenames in w2html 22a> Referenced in 21b.
 <get program options 3a> Referenced in 11b.
 <impliciete make regels 16a, 17b, 20b> Referenced in 13b.
 <import modules in main program 3b, 6b, 7c, 9> Referenced in 11b.
 <make targets 17a, 19c, 23b, ?> Referenced in 13b.
 <methods of the main program 2, 3cd, 4ab, 5b, 6c, 7d, 10abc> Referenced in 11b.
 <parameters in Makefile 13a, 14a, 15bc, 17c, 20a, 23a> Referenced in 13b.
 <perform the task of w2html 21c> Referenced in 21b.
 <print the testpost 12> Not referenced.
 <remove the copy of the aux file 18d> Referenced in 18b, 21c.
 <run tex4ht 22d> Referenced in 22b.
 <run the html processors 22c> Referenced in 22b.
 <run the html processors until the aux file remains unchanged 22b> Referenced in 21c.
 <run the processors until the aux file remains unchanged 19b> Referenced in 18b.
 <run the three processors 19a> Referenced in 19b.
 <variables of the main program 8, 11a> Referenced in 11b.
 <yield data from articles in this soup 6a> Referenced in 5b.
 <yield topic data from soup 5a> Referenced in 4a.

C.3 Variables

all: 13c.
 auxfil: 18c, 19b, 22a, 22b.
 BeautifulSoup: 3b, 3c, 10b.
 bibtex: 19a, 22cd.
 boardURL: 3a, 10b, 11b.
 bs4: 3b.
 calendar: 8, 9.
 convert_timestring: 7b, 7d.
 datetime: 9, 12.
 dateutil.parser: 7c, 7d.
 DIRS: 23b, ?.
 fig2dev: 16a.
 FIGFILENAMES: 15c.
 FIGFILES: 15b, 15c, 20a.
 get_boardsoup: 10b.
 indexfil: 18c, 19b, 22a.
 makeindex: 19a, 22cd.
 make_soup_from_url: 3c, 4a, 5b.
 MKDIR: 23a, 23b.
 monthnums: 8.
 next_topic: 4a, 11b.
 nufil: 18c, 19a, 22a, 22c.
 nuweb: 13a, 17cd, 18a, 19a, 21b.
 oldaux: 18c, 18d, 19b, 22a, 22b.

oldindexfil: [18c](#), [19b](#), [22a](#).
pdf: [13d](#), [14a](#), [17a](#), [17b](#).
PDFT_NAMES: [15c](#), [17b](#).
PDF_FIG_NAMES: [15c](#), [17b](#).
PHONY: [13c](#), [16b](#).
postnum: [6a](#), [11b](#).
print: [2](#), [3c](#), [6c](#), [10b](#), [11b](#), [14b](#), [17a](#).
printnaf: [6c](#), [11b](#).
print_post: [2](#), [12](#).
PST_NAMES: [15c](#).
PS_FIG_NAMES: [15c](#).
requests: [3b](#), [3c](#), [10b](#).
sbody: [5a](#), [6a](#), [10c](#).
SUFFIXES: [14a](#).
texfil: [18c](#), [19a](#), [22a](#), [22c](#).
trunk: [18c](#), [19a](#), [22a](#), [22cd](#).
view: [17a](#).