

# User Manual Araya Bilingual Term Extraction

**BiLingual  
TermExtractor**

Description how to use the Araya  
Terminology Extraction Tool

# Bilingual Extractor

- The bilingual extractor is a simple to use, but efficient tool to generate automatically term pairs from translated documents (TMX files)
  - A term pair consists of a source and target term
  - A term can consists of several words
- These term pairs can be used to create a new terminology or add new terminology to an existing terminology database.

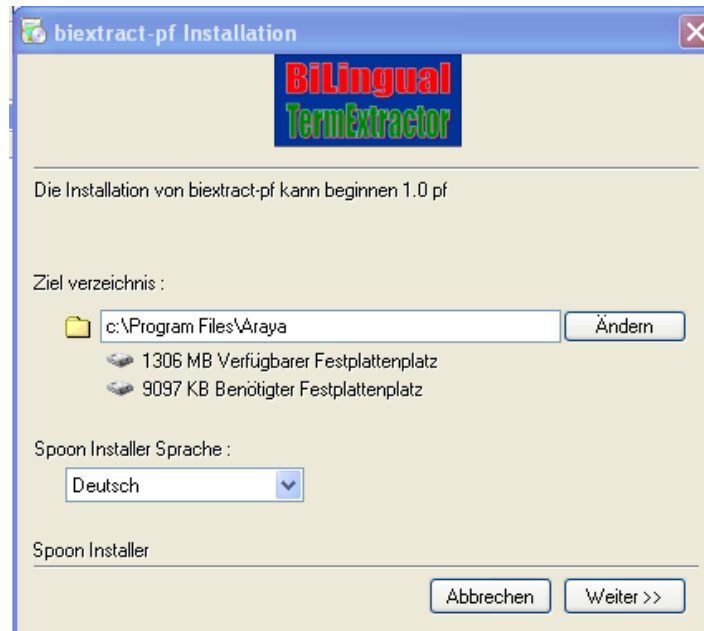
# Versions

- The Extractor was developed by Heartsome Europe GmbH and LNE International.
- It is available as
  - Single user version.
  - It is part of the Araya Server translation tools.

# Short introduction how to extract

- Extract Terms
  - File -> Extract Bilingual Terminology from File
    - (Option: Open after extraction)
- Check extracted terms
  - Mark correct translation as „validated“
- Export terms
  - Export validated terms ...

# Installation



- Installation goes into directory **c:/Program Files/Araya**. It is recommended not to change this as all initialization files map towards this directory.

# Starting Araya Extraction tool

- Go to directory:  
c:/Program Files/Araya  
Start: BiEdit.exe
- Or double click :



# The Extraction Approach

- Based on a TMX file all possible relevant term pairs are computed. This is based on a statistical approach which determines the frequency of source and target terms.
- TMX = XML Exchange format for translation memory databases

# Segment

- Extraction is based on segments which are stored in a TMX file.
- A segment can either be a sentence or a whole paragraph.
- Formats in TMX files are ignored.



# Evaluating and Validating

- Each found term pair is associated with a quality measure.
  - 2. column of the extraction table
  - Value is between 1,0 (highest probability that the term pair is a translation) and 0,5 (lowest probability that the term pair is a translation)
- Terms can be validated as correct translations.
  - Last column of table
    - Approved = checked = validated
    - Unapproved = not validated
- Validated terms can be exported.

# Validating a Term Extraction Pair

- Select line with term extraction pair
- Validate = approving using
  - Double click term pair
  - Right mouse click
- Remove validation mark using
  - a double or
  - right mouse click

0,97	30	1/1
0,97	33	1/1
0,97	41	1/1
0,97	22	1/1

0,97	113
0,97	27
0,97	27
0,97	52

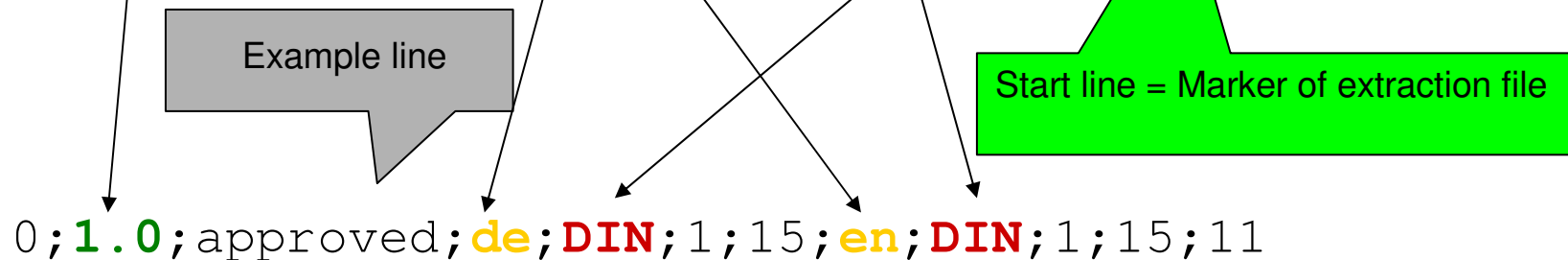
287	0,97	113	1/1	764	787	Kyoto	kyoto	unap...
288	0,97	27	1/1	32	30	Frettchen	ferrets	appr...
289	0,97	27	1/1	27	29	Frischenschlager	frischenschlager	unap...
290	0,97	27	2/2	29	30	El Salvador	el salvador	unap...
291	0,97	52	2/2	130	130	Sierra Leone	sierra leone	unap...

Validated terms appear in green

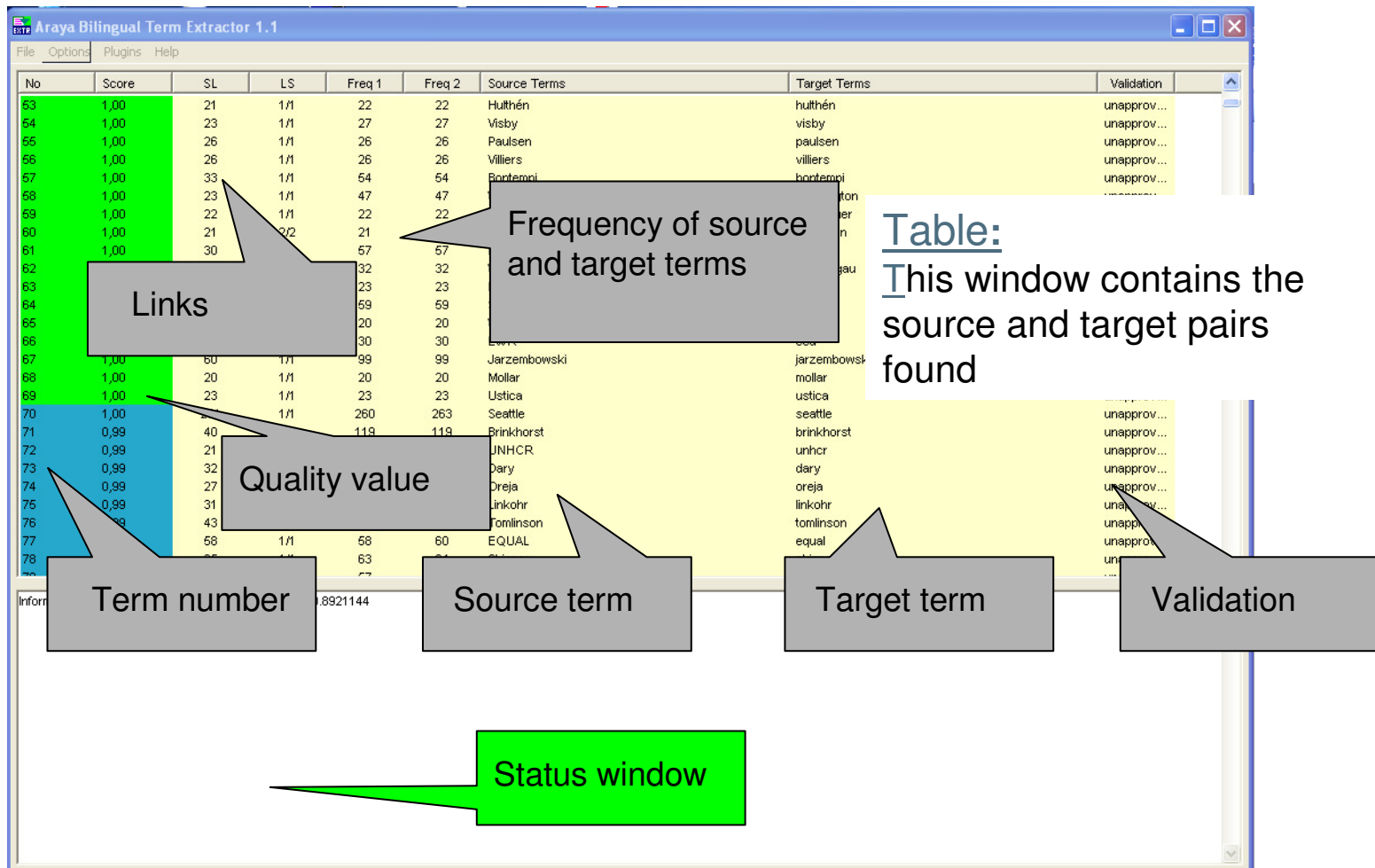
# The Extraction file

- The extraction file has the following format
- Encoding format is UTF-8!

```
nr;score;status;term1.LangCode;term1.wordGroup;term1.wordGroup
Len;term1.wFreq;term2.LangCode;term2.wordGroup;term2.wordGroup
Len;term2.wFreq;sentLinked
```



# Extraction User Interface



The screenshot shows the 'Araya Bilingual Term Extractor 1.1' window. The main area is a table with columns: No, Score, SL, LS, Freq 1, Freq 2, Source Terms, Target Terms, and Validation. The table contains several rows of extracted terms. Callouts point to specific parts of the interface:

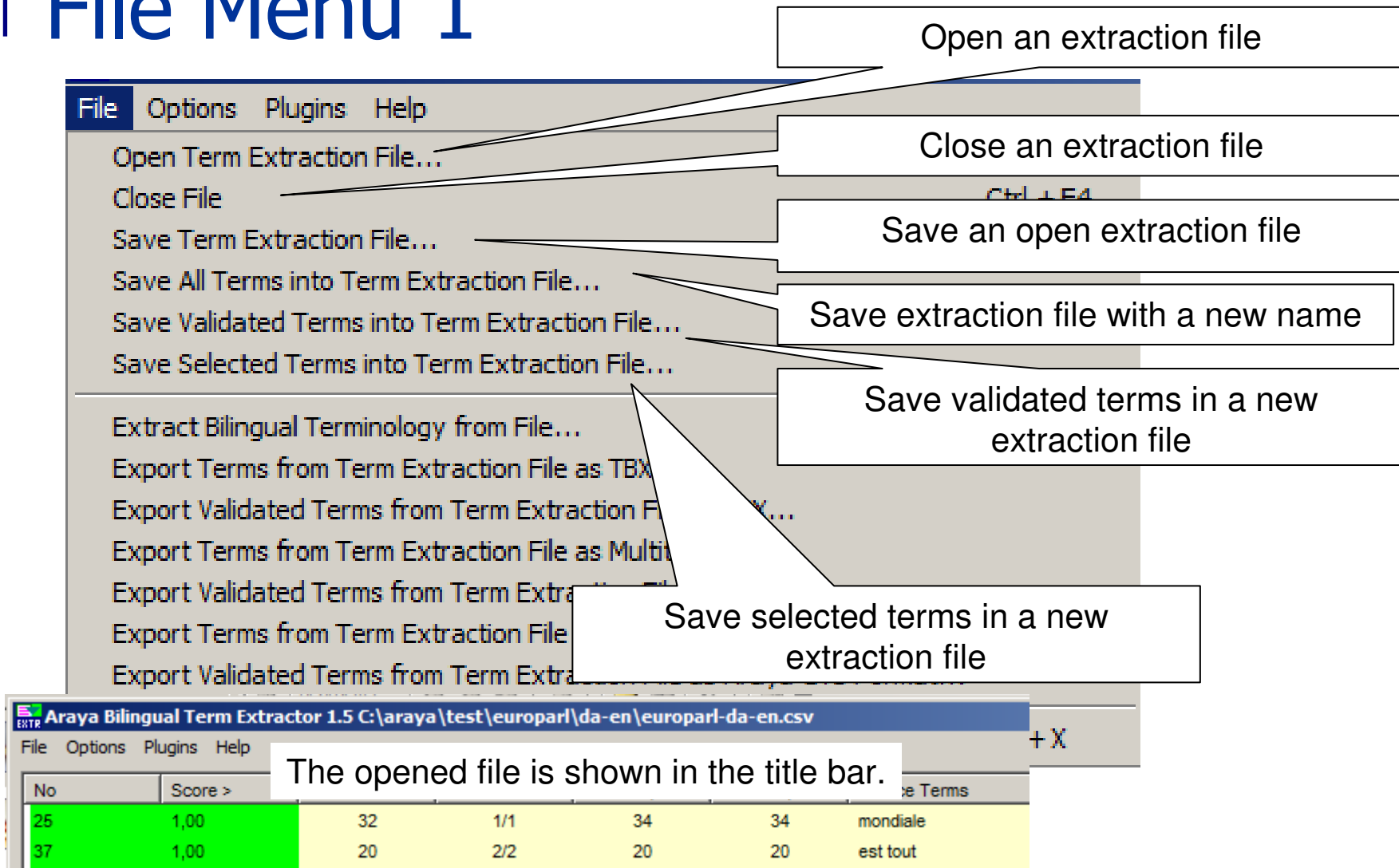
- Links:** Points to the 'No' column.
- Quality value:** Points to the 'Score' column.
- Frequency of source and target terms:** Points to the 'Freq 1' and 'Freq 2' columns.
- Table:** A text box stating 'This window contains the source and target pairs found'.
- Term number:** Points to the 'No' column.
- Source term:** Points to the 'Source Terms' column.
- Target term:** Points to the 'Target Terms' column.
- Validation:** Points to the 'Validation' column.
- Status window:** A green box at the bottom with a callout pointing to the status bar area.

No	Score	SL	LS	Freq 1	Freq 2	Source Terms	Target Terms	Validation
53	1,00	21	1/1	22	22	Huthén	huthén	unapprov...
54	1,00	23	1/1	27	27	Visby	visby	unapprov...
55	1,00	26	1/1	26	26	Paulsen	paulsen	unapprov...
56	1,00	26	1/1	26	26	Villiers	villiers	unapprov...
57	1,00	33	1/1	54	54	Bontempi	bontempi	unapprov...
58	1,00	23	1/1	47	47			unapprov...
59	1,00	22	1/1	22	22			unapprov...
60	1,00	21	2/2	21	21			unapprov...
61	1,00	30		57	57			unapprov...
62				32	32			unapprov...
63				23	23			unapprov...
64				59	59			unapprov...
65				20	20			unapprov...
66				30	30			unapprov...
67	1,00	60	1/1	99	99	Jarzebowski	jarzebowski	unapprov...
68	1,00	20	1/1	20	20	Mollar	mollar	unapprov...
69	1,00	23	1/1	23	23	Ustica	ustica	unapprov...
70	1,00		1/1	260	263	Seattle	seattle	unapprov...
71	0,99	40		119	119	Brinkhorst	brinkhorst	unapprov...
72	0,99	21				JNHCR	unhcr	unapprov...
73	0,99	32				Dary	dary	unapprov...
74	0,99	27				Dreja	oreja	unapprov...
75	0,99	31				linkohr	linkohr	unapprov...
76	0,99	43				Tomlinson	tomlinson	unapprov...
77		58	1/1	58	60	EQUAL	equal	unapprov...
78				63				unapprov...

# Columns

- Value
  - Statistical measure that the source and target term are translation of each other (quality measure)
- SL
  - Number of segments where both source and target term appear in.
- Freq 1
  - Number of segment where source term appears in
- Freq 2
  - Number of segment where target term appears in
- Source term
  - The source term
- Target term
  - The translation of the source term
- Validation
  - Check box, for marking correct term pairs

# File Menu 1



The screenshot shows the File menu with the following items and callouts:

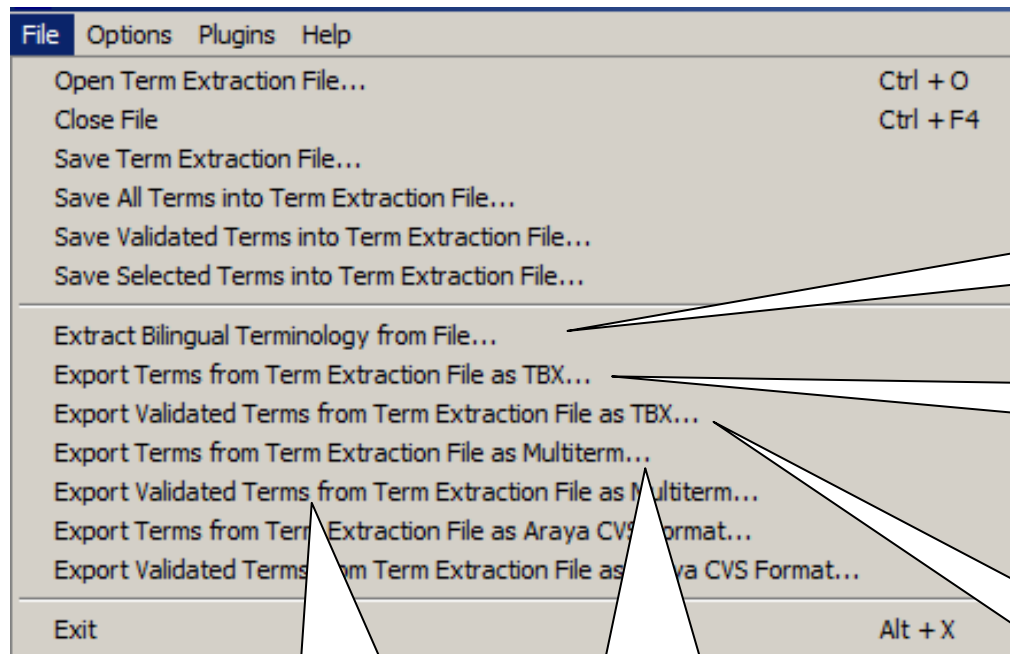
- Open Term Extraction File... (Callout: Open an extraction file)
- Close File (Callout: Close an extraction file)
- Save Term Extraction File... (Callout: Save an open extraction file)
- Save All Terms into Term Extraction File... (Callout: Save extraction file with a new name)
- Save Validated Terms into Term Extraction File... (Callout: Save validated terms in a new extraction file)
- Save Selected Terms into Term Extraction File... (Callout: Save selected terms in a new extraction file)
- Extract Bilingual Terminology from File...
- Export Terms from Term Extraction File as TBX...
- Export Validated Terms from Term Extraction File as TBX...
- Export Terms from Term Extraction File as Multi...
- Export Validated Terms from Term Extraction File as Multi...
- Export Terms from Term Extraction File as Multi...
- Export Validated Terms from Term Extraction File as Multi...

The title bar of the application window shows: Araya Bilingual Term Extractor 1.5 C:\araya\test\europarl\da-en\europarl-da-en.csv

The opened file is shown in the title bar.

No	Score >					Terms
25	1,00	32	1/1	34	34	mondiale
37	1,00	20	2/2	20	20	est tout

# File Menu 2



Extract term pairs from a TMX file

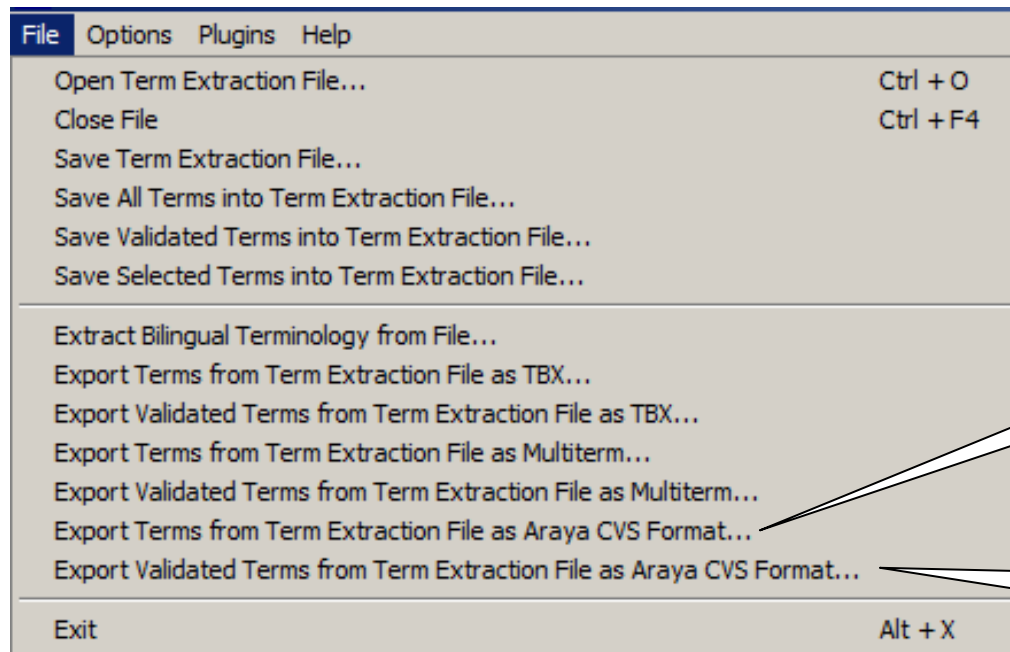
Export terms from open extraction file into TBX Format

Export validated terms from open extraction file into Multiterm Format

Export terms from open extraction file into Multiterm Format

Export validated terms from open extraction file into TBX Format

# File Menu 3

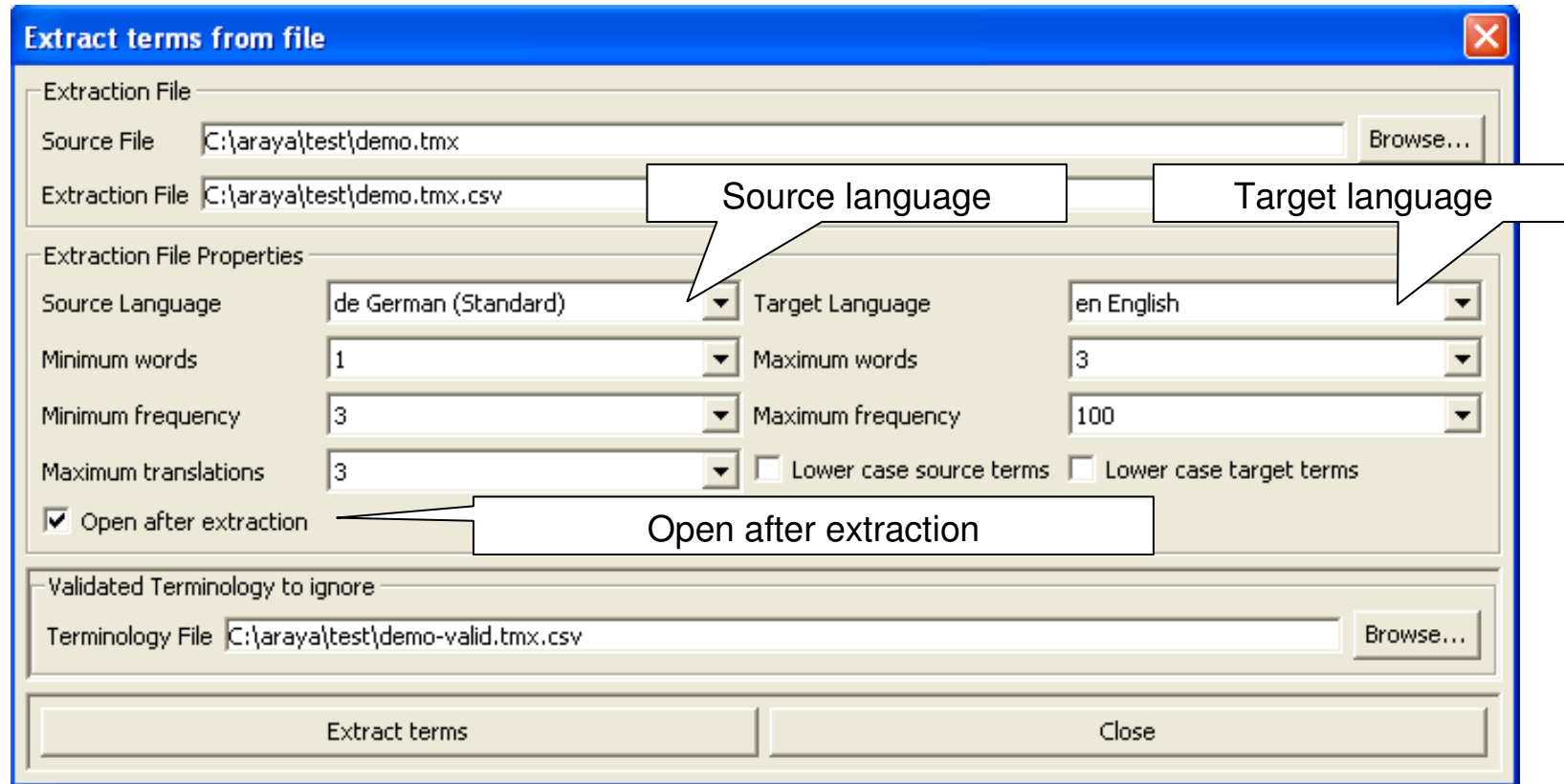


Export entries from open extraction file in Araya CSV Format

Export validated entries from open extraction file in Araya CSV Format



# Extracting Term Pairs from TMX File



# Extraction Parameters 1

- Minimum / Maximum words
  - Controls how many words should be contained min/max in the found term pair
- Minimum / Maximum frequency
  - Controls how often the term should appear min/max for the found term pair
- Maximum Translations
  - Controls how many translation should be found at maximum
- Source/Target terms in lower case
  - Controls if source and/or target terms should be converted to lower case

## Extraction Parameters 2

- Validated terminology to ignore
  - If a terminology extraction file is specified here, all terms which are marked as “validated” will be ignored.
  - Thru this know translations are ignored.

# Exporting

- Exporting can be done in different formats.
  - TBX
    - Name of extraction file + „.tbx“
  - Multiterm (™ of Trados/SDL International)
    - Name of extraction file + „.multiterm“
  - Araya CSV
    - Name of extraction file + „araya.csv“
  - Character encoding is always UTF-8
- Either all or only the validated entries can be exported
- In addition the selection filter (Options -> Export Score Filter) controls the exported terms
  - Depending on chosen value only the term pairs with a minimum score get exported (e.g. score higher than 0.6).

# Araya CSV Format

- Araya CSV Format contains the languages in the first line followed by the extracted terms

## Beispiel

**de; en**

Languages separated by ;

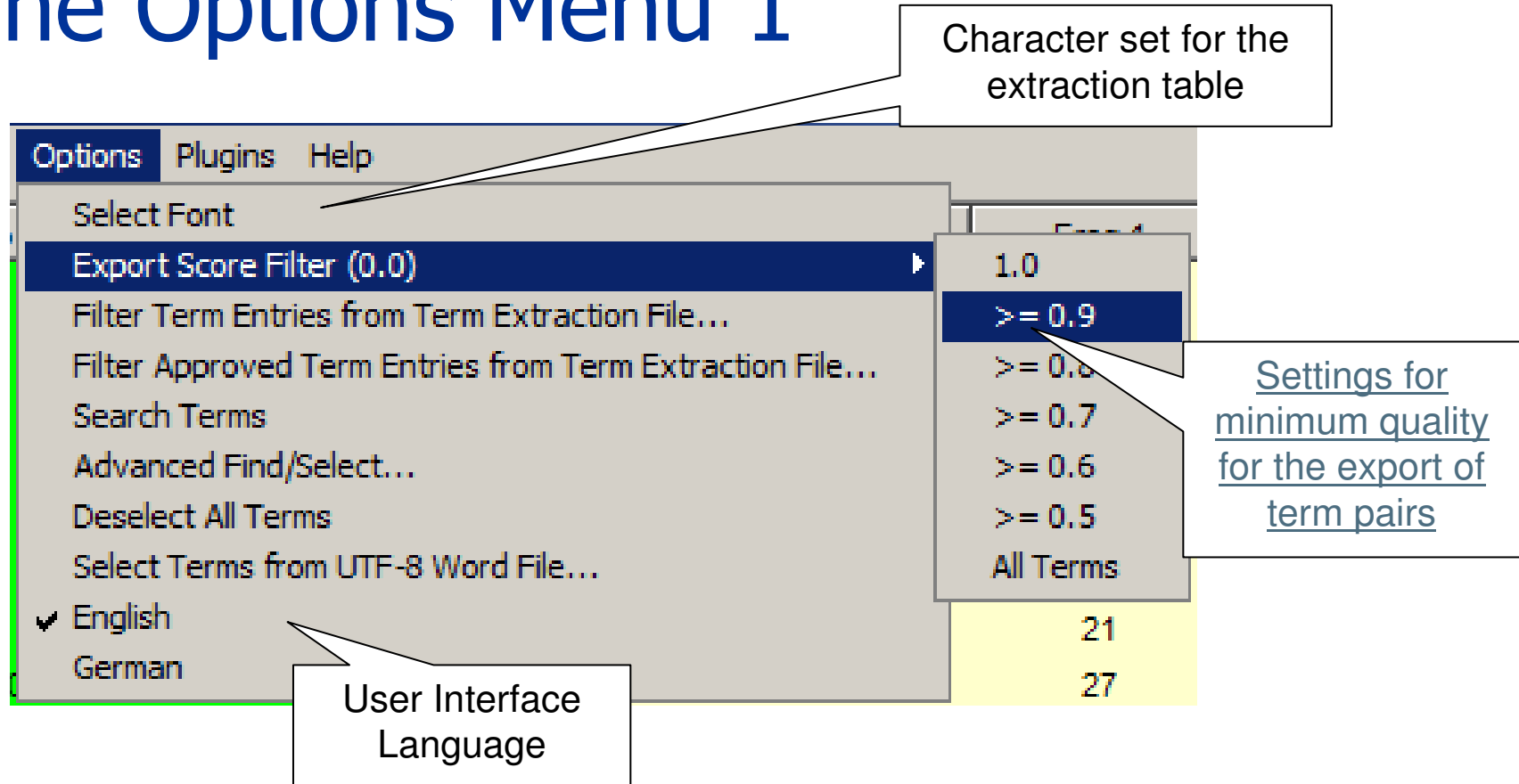
Anschlussplan; Connection diagram

DIN; DIN

Dr; Dr

Extracted terms separated by ;

# The Options Menu 1



The screenshot shows the 'Options' menu with the following items:

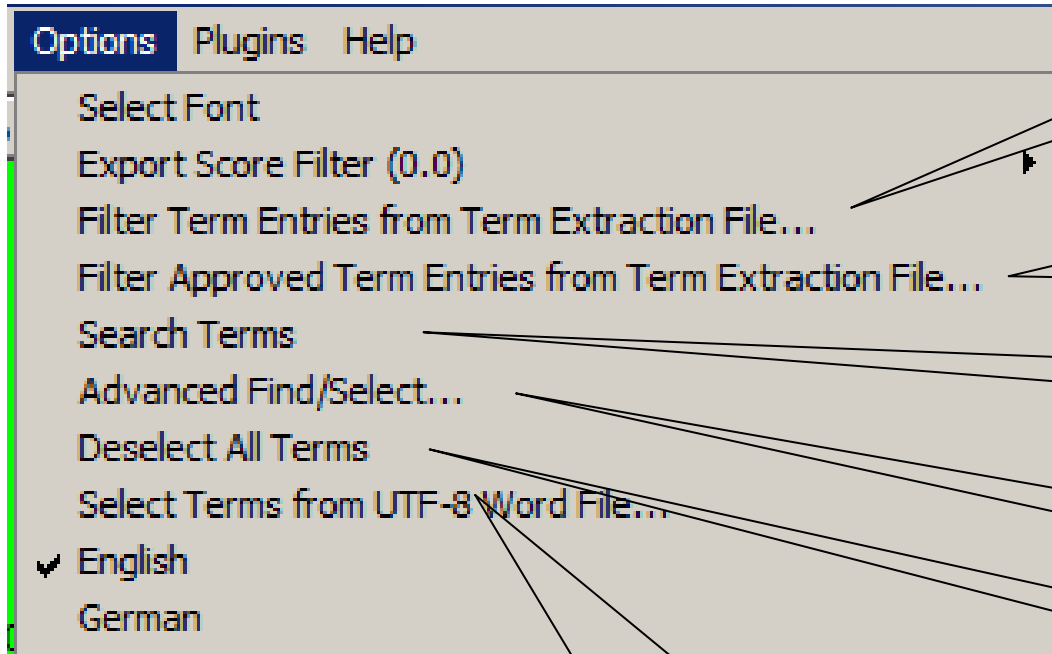
- Select Font
- Export Score Filter (0.0) ▶
- Filter Term Entries from Term Extraction File...
- Filter Approved Term Entries from Term Extraction File...
- Search Terms
- Advanced Find/Select...
- Deselect All Terms
- Select Terms from UTF-8 Word File...
- ✓ English
- German

Callouts provide additional information:

- Character set for the extraction table**: Points to the 'Select Font' option.
- Settings for minimum quality for the export of term pairs**: Points to the 'Export Score Filter' sub-menu, which lists options: 1.0, **>= 0.9**, >= 0.8, >= 0.7, >= 0.6, >= 0.5, and All Terms.
- User Interface Language**: Points to the 'English' and 'German' options.

21
27

# The Options Menu 2



Remove all term entries from the term table contained in term extraction file

Remove all term entries from the term table contained in term extraction file which are approved

Search Terms in the extraction table

Complex search for based on source and target term

Deselects all selected in the extraction table

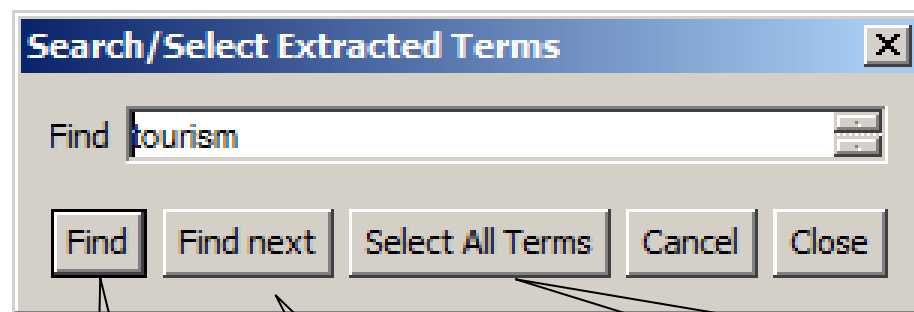
Select Terms based on a word list in a file

# Filter Term Functions

- The filter functions filters all those term entries which are contained in another term extraction file.
- The identical terms are removed from the term table.
- Depending on the chosen filter method the approved or all entries are used from the specified term extraction file.



# Search Term Functions

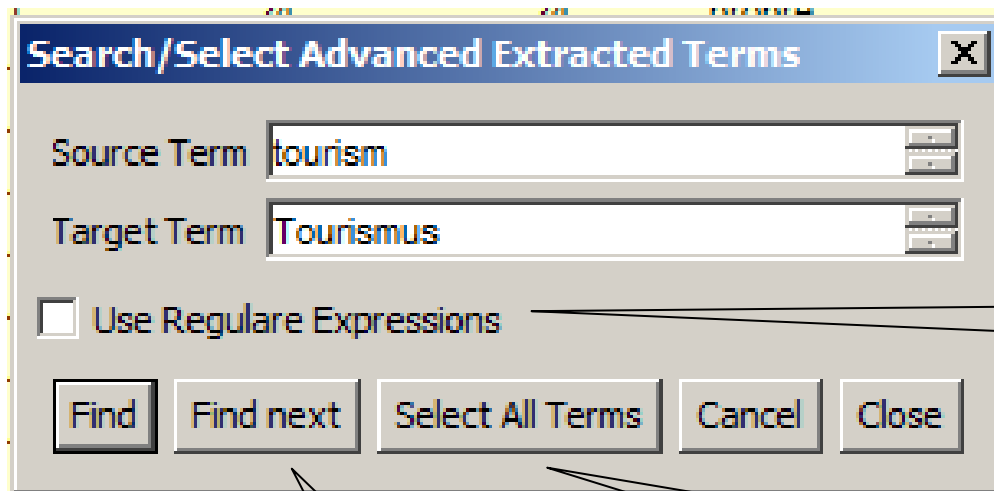


Start searching this term

Find the next matching term

Search this term and select all the matching entries in the table. The selected term candidates can then be saved with "File -> Save Selected Terms into Extraction File..."

# Advanced Search Term Functions



Use regular expressions for searching term entries

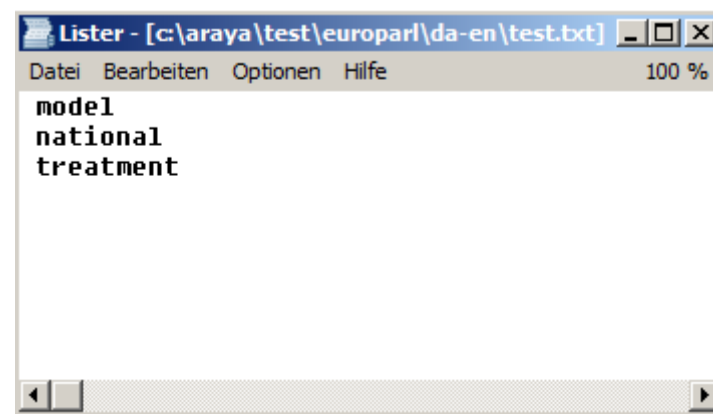
Search this term and select all the matching entries in the table. The selected term candidates can then be saved with "File -> Save Selected Terms into Extraction File..."

Find the next matching terms

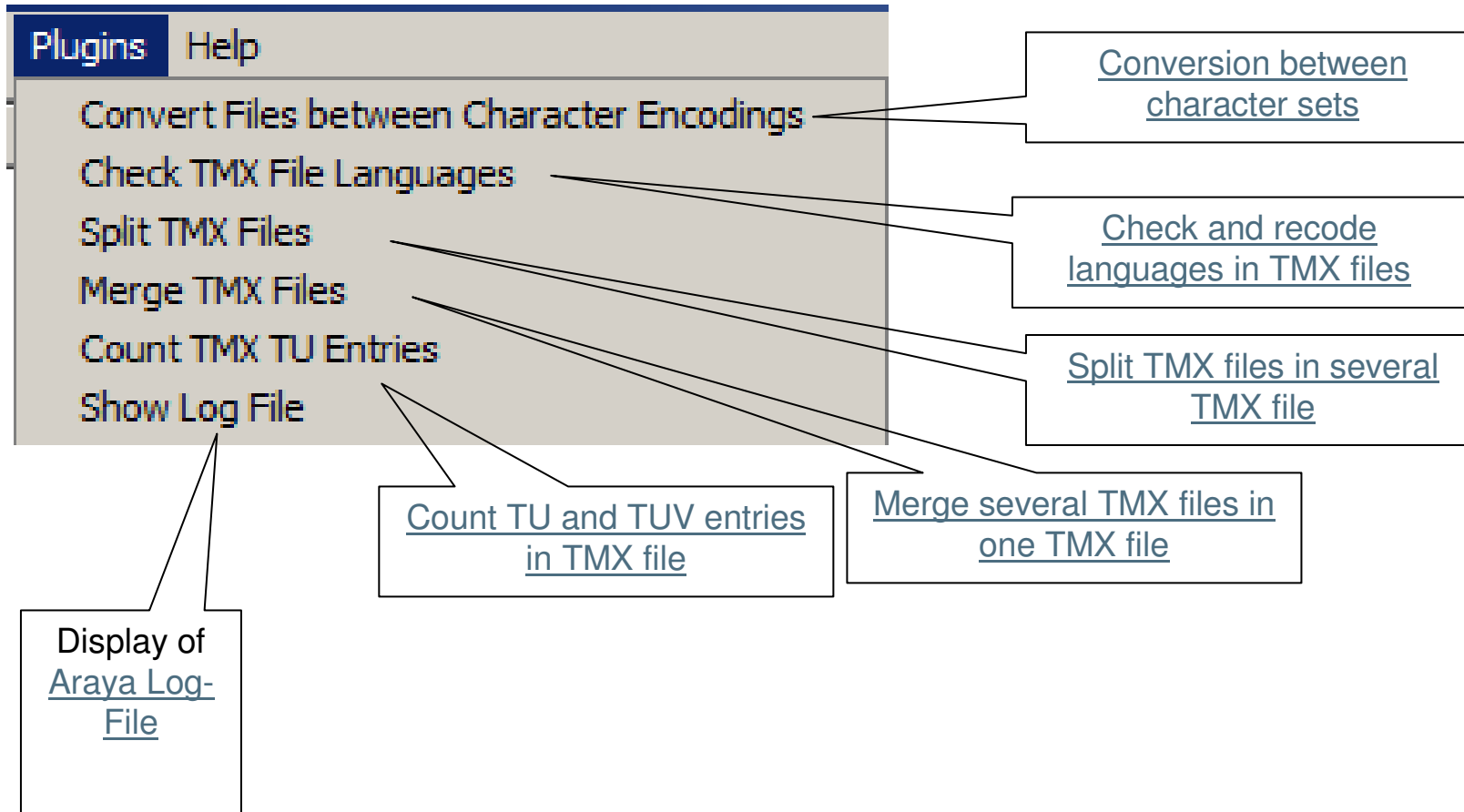
Start searching these combination of source and target terms

# Select Terms Based on a Word List

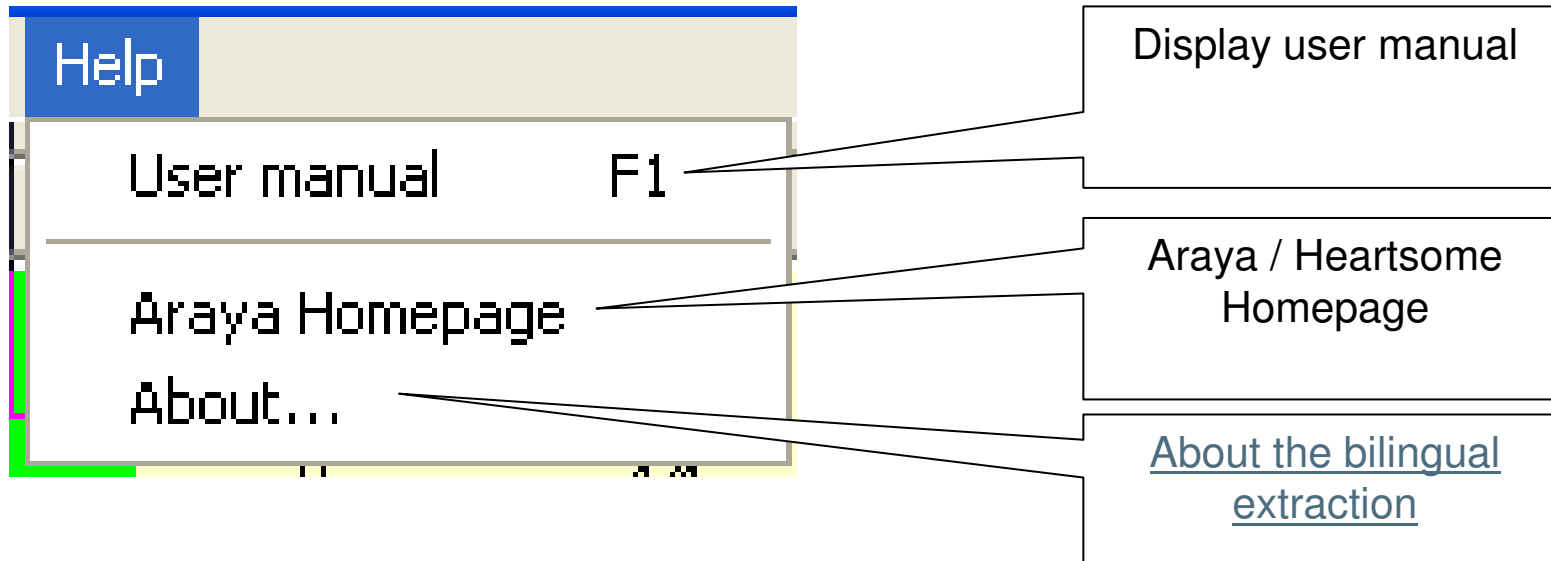
- The filter functions selects all the matching terms from an UTF-8 encoded file. The filter search the given string in both source and target language string
  - It also matches substrings
- The selected term candidates can then be saved with “File -> Save Selected Terms into Extraction File...”
- Each word (string) in the file must be written on a single line.



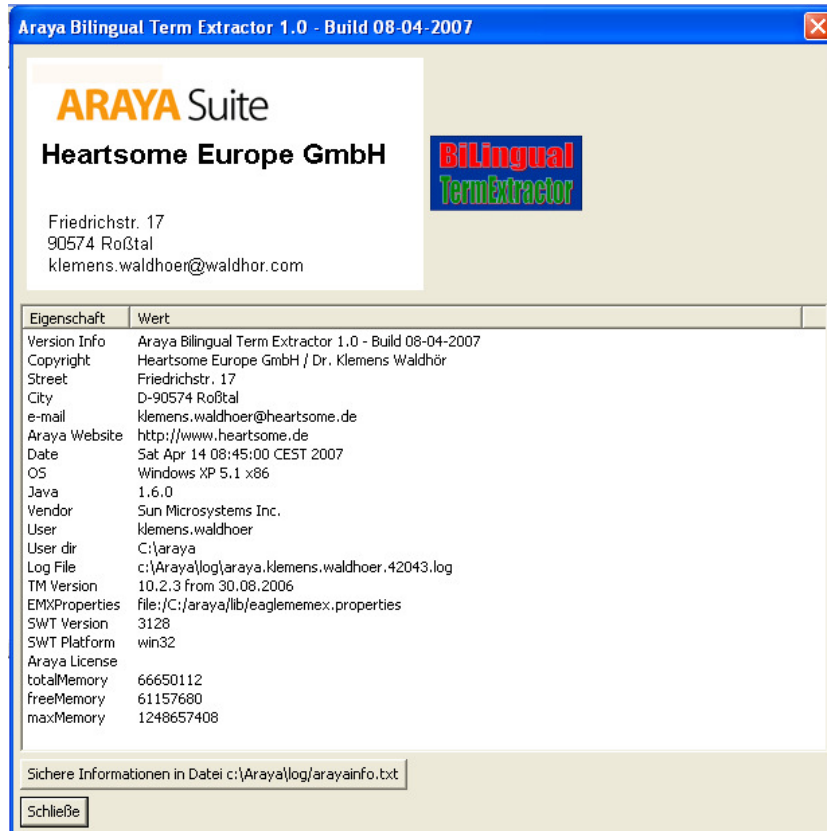
# The Plugins Menu



# The Help Menu

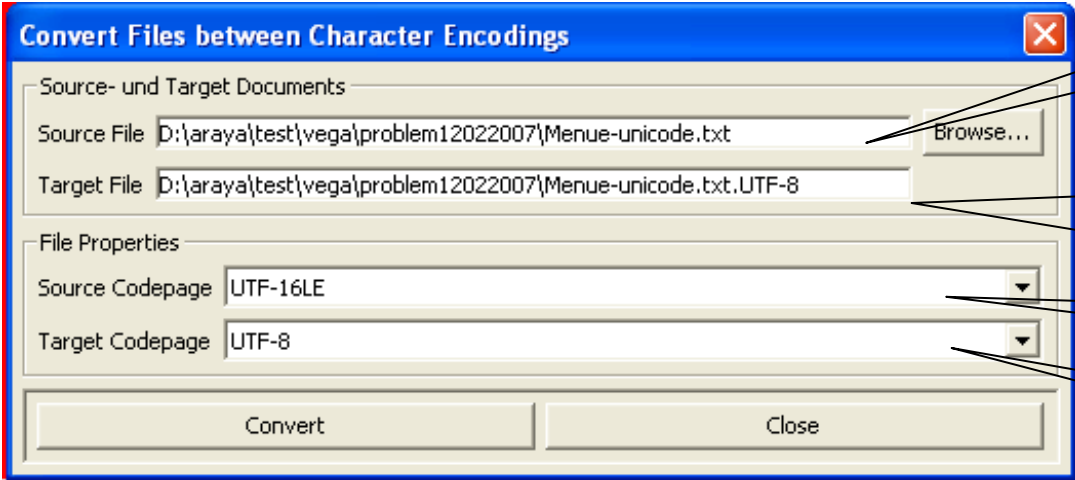


# About the Extraction



- Basic information about the Araya configuration is shown here.
- This can be important for error diagnosis
- Using „Save...“ the information is stored in the specified file.

# Converting Files between different Character Encodings



The screenshot shows a dialog box titled "Convert Files between Character Encodings". It has two main sections: "Source- und Target Documents" and "File Properties".

- Source- und Target Documents:**
  - Source File: D:\araya\test\vega\problem12022007\Menue-unicode.txt
  - Target File: D:\araya\test\vega\problem12022007\Menue-unicode.txt.UTF-8
- File Properties:**
  - Source Codepage: UTF-16LE
  - Target Codepage: UTF-8

Callouts point to the following fields:

- Source file name (points to the Source File field)
- Target file; will be created based on source file name and as extension the target encoding character set (points to the Target File field)
- Source encoding character set (points to the Source Codepage dropdown)
- Target encoding character set (points to the Target Codepage dropdown)

This function supports converting files between different character sets. Depending on the target character up to two files are written. If the target file is a UTF-8, 16, 32 or UCS file, two files are written. The second file with the extension „**nobom**“ is a copy of the first target file, the only difference is that the BOM (Byte Order Marks) are removed from this file. **This file should be used for importing, esp. when an import of an UTF-8 file is done in Araya, as the Java reading functions for UTF-8 does not over read the BOM characters. This could lead to problems when reading normal strings from those file as the BOMs are read as normal characters resulting in invalid entries.**

(This is a known bug in Java UTF-8 file reading, but will not be corrected by SUN!).

# Check and Recode TMX Files

**Check TMX File Languages**

TMX File

Source File: D:\araya\test\vega\problem15032007\translationwork-15.03.2007.tmx

Recoded File: [Empty]

Nr	TMX File Language	New Language
1	zh	zh-cn Chinese (People's Republic)
2	it	
3	zh-cn	
4	en-us	
5	cs	
6	de	
7	pt	
8	pl	
9	fr	
10	en	
11	ru	

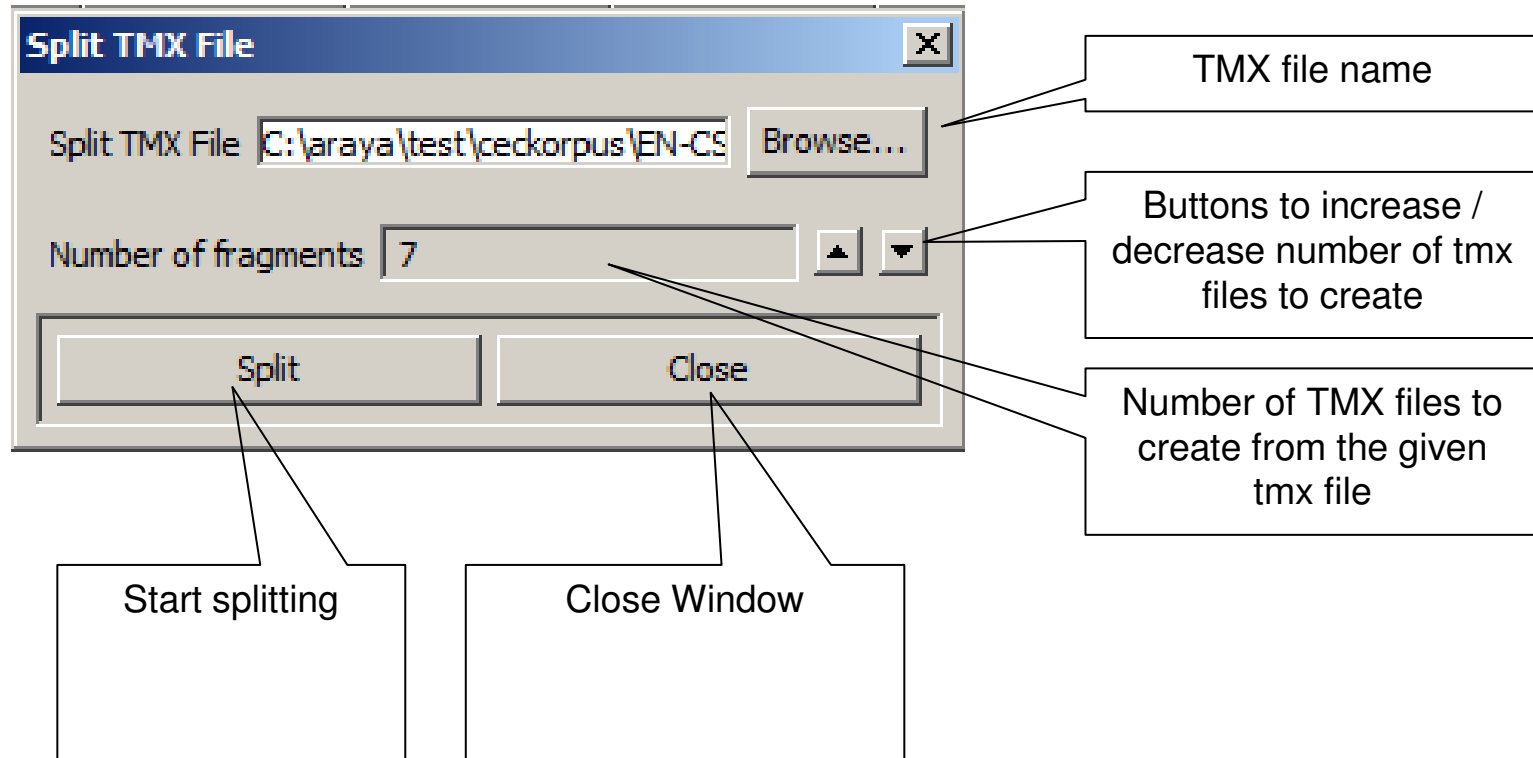
Buttons: Check Languages, Recode Languages, Close

**Callouts:**

- TMX file name
- Recoded TMX file name; Consists of tmx file name, any ".tmx" extension removed and ".recoded.tmx" added
- Language codes supported by Araya Combo box appears only in case language codes in TMX file not supported by Araya
- Light yellow lines indicate language codes not supported by Araya
- Green lines indicate language codes supported by Araya
- Start checking languages contained in TMX File
- Recode languages using new language for the specified TMX File language

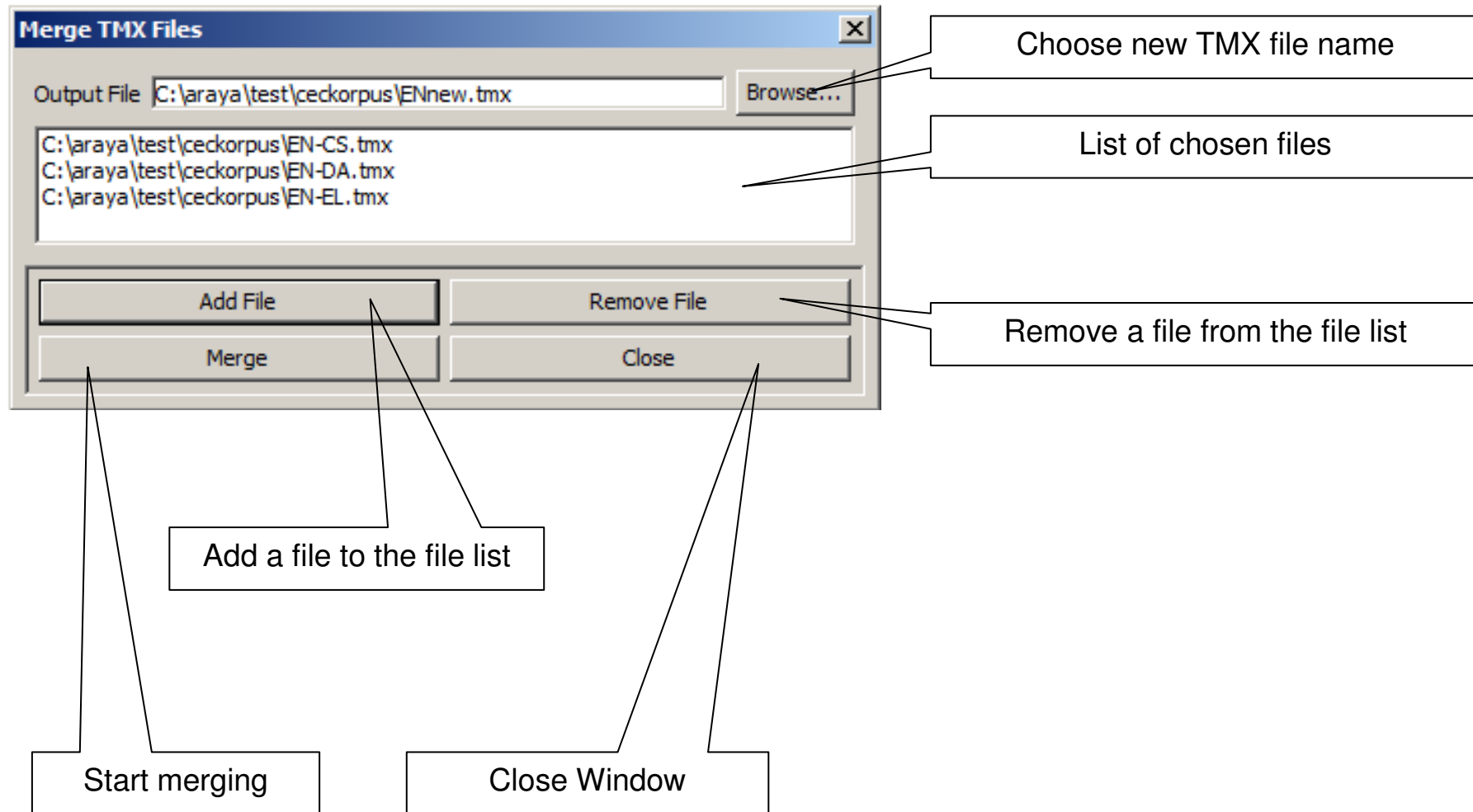


# Split TMX File



The created tmx files are composed of the original tmx file name plus a number from 1 till the number of chosen tmx file

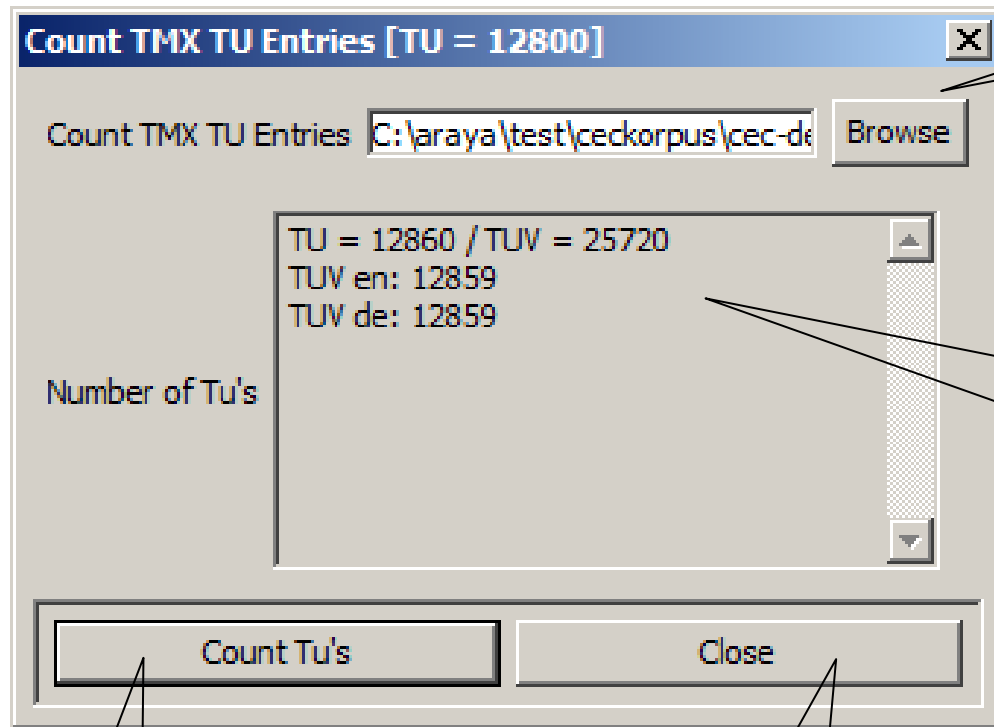
# Merge TMX Files



The screenshot shows a dialog box titled "Merge TMX Files" with the following components and callouts:

- Output File:** A text field containing "C:\araya\test\ceckorpus\ENnew.tmx" and a "Browse..." button. Callout: "Choose new TMX file name".
- File List:** A list box containing three file paths: "C:\araya\test\ceckorpus\EN-CS.tmx", "C:\araya\test\ceckorpus\EN-DA.tmx", and "C:\araya\test\ceckorpus\EN-EL.tmx". Callout: "List of chosen files".
- Buttons:** "Add File", "Remove File", "Merge", and "Close".
  - Callout for "Add File": "Add a file to the file list".
  - Callout for "Remove File": "Remove a file from the file list".
  - Callout for "Merge": "Start merging".
  - Callout for "Close": "Close Window".

# Count TUs/TUVs in a TMX File



Choose TMX file to analyse

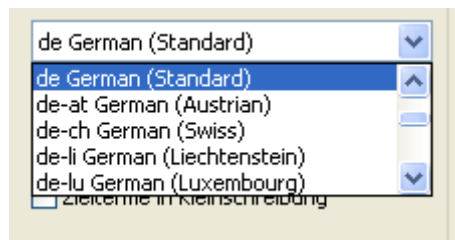
Result:  
First line: Number of TUs and TUVs  
Following lines:  
Number of language specific TUVs

Start counting

Close Window

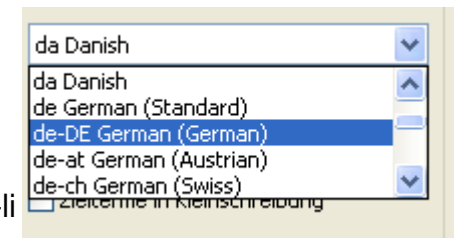
# Adding Language Codes

- A set of predefined language codes is defined in the file „ini/lancodes.txt“.
- Additional language codes can be added by modifying this file.
- Here is an example:



**Adding de-DE:** This requires to add a line like that – where = separates the name displayed on the actual language code  
German(DE)=de-DE

```
Galician=gl
German (Standard)=de
German(DE)=de-DE
German (Austrian)=de-at
German (Liechtenstein)=de-li
German (Luxembourg)=de-lu
German (Standard)=de
German (Swiss)=de-ch
```



# Impressum

- Heartsome Europe GmbH
- Friedrichstr. 17
- D-90574 Roßtal
  
- Email: [info@heartsome.de](mailto:info@heartsome.de)
- [www.heartsome.de](http://www.heartsome.de)
- © 2007, 2009 Heartsome Europe GmbH