# TROVA - Annotation Content Search

**The latest version can be found at: http://tla.mpi.nl/tools/tla-tools/trova/**

**This manual was last updated in May 2013**

**The Language Archive, MPI for Psycholinguistics, Nijmegen, The Netherlands**

# TROVA - Annotation Content Search: version 1.4

The latest version can be found at: http://tla.mpi.nl/tools/tla-tools/trova/

This manual was last updated in May 2013

# Table of Contents

# Chapter 1. TROVA: Contents search

If you want to perform a detailed search over multiple annotation files in different formats (including: EAF, SHOEBOX, TOOLBOX, CHAT, SUBRIP and TEXTGRID) you can use the **content search** function. This allows you to restrict the search to certain tiers (i.e. the different layers present in a hierarchy-structured annotation file), to use regular expressions, etc., while examining multiple annotation files at once.

To include multiple files in the search domain, select one or more nodes in the tree-view of the IMDI browser by clicking on the nodes while holding down the CTRL key on the keyboard. If you click on the button content search, you will include all annotation files under the selected corpus nodes in the search domain. Alternatively, you can select a parent node, in which case all the child nodes will be included in the search.



**Figure 1.1. Adding nodes to the search engine**

A new browser frame for the search engine will appear with the search domain displayed at the top of each of the three search tabs:



**Figure 1.2. Search Engine**

The voice **Domain** shows the name of the selected node(s) in the archive, whereas **Types** shows all the types of annotation files in the domain of that /those node(s). You can include or exclude types of annotation files in the search by respectively checking or unchecking them. Remember that, in order to be able to see the various annotation files, you have to log in first.

It is possible that you are not allowed to read files from certain parts of the domain due to AMS rules (http://tla.mpi.nl/tools/tla-tools/ams/). In that case all searchable files from those parts are excluded from the annotation content search. As a result, the search domain contains fewer files.

Please remember that even if you have received access to the resources through AMS, it could take up to one day for the changes to be recognised by TROVA as well.

As can be seen from figure 1.2 above, there are three tabs offering different kinds of search:

- Substring Search: it finds all annotations in which the search string occurs.

- Single Layer Search: it finds all annotations or N-grams (i.e. search strings of more than one word, either consecutive or not) in which either the search string or the regular expression occurs, both case sensitive and insensitive and possibly restricted to one (type of) tier.

- Multiple Layer Search: it finds annotations within one single tier but also in three related tiers at the same time. You can use multiple search strings or regular expressions and make constraints on duration and time slot as well as constraints on how the search strings are to be combined.

# 1.1. Substring Search Tab

This tab offers the simplest search. It just asks for a search string. After entering the search string you can click on Find (or press **Enter**) to start the search process. This will result in a page like the one below:



**Figure 1.3. Substring Search Results**

The result page above shows the annotations containing the search string plus some annotations in the context printed in italic. The default number of annotations in the context is four on both sides, but you can change it through the **Context size** menu. When the number of hits exceeds the maximum number the window can contain, you can view the rest of the hits by clicking the < and > buttons that appear above the list of hits. Besides the ones just described, there are other options you can choose from.

From the **Action** menu you can select:

**1) Show Concordance View**: it is the default mode in which results are shown. The search strings is in bold and has some context on both sides.

**2) Show Frequency View Sorted by Annotation**: the results are shown in the right bottom half of the page in alphabetical order (see the first word of the annotation).

**Figure 1.4. Frequency View Sorted by Annotation**

**3) Show Frequency View Sorted by Frequency:** the results are shown in the left, bottom half of the page, under percentage form, in a decreasing order.



**Figure 1.5. Frequency View Sorted by Frequency**

**4) Show Selected Hit in Transcription**: as the name says, first of all select one of the results, and then select the option itself (otherwise simply double click the result you need). You will be redirected to Annex interface, where you will see the annotation in the Timeline viewer [see section 2.2.1. of Annex manual; use the following link: http://www.mpi.nl/corpus/html/annex/ch02s02s01.html].

**5) Show Selected Hit in Corpus Tree**: like with option 4), first select the result you need, and then this option. You will be redirected to the IMDI browser interface, where you will see: a) the resource to which the annotation belongs highlighted in grey in the IMDI tree; b) some information about such resource in the main content panel.

6) **Save Hits**: this option allows you to save the results you have obtained from your query. Once selected, it will open a dialogue window asking you to choose: a) the export format (either UTF-8 or UTF-16); b) the fields to export (e.g. begin and end time, tier type and name, etc.).

> **⚠ Note**
>
> If you are in doubt about which one to choose between UTF-8 and UTF-16 because you are not very familiar with character encoding, opt for UTF-8.

Moreover, below the **Action** drop-down menu, you can find a **Font** menu, which allows you to choose the style of writing and the size in which the results will be shown.

Finally, the option **Show Info Balloons**, which is displayed on the right of the **Font** menu, allows you to see some additional information about the annotations you have searched for. If you check this option, and then place the mouse cursor on one of the results, you will see a yellow balloon appearing, containing some information about that result, like for example the transcription link, the tier name and type, the begin and end time, etc. See figure 1.4. below.



**Figure 1.6. Show Info Balloons**

> **⚠ Note**
>
> All the features seen above, related to the **Substring Search**, are also present (with the same functions) in the **Single Layer Search** and in the **Multiple Layer Search**.

# 1.2. Single Layer Search Tab

The Single Layer Search tab offers a more elaborate search than the Substring Search tab. The first thing that is different from the Substring Search tab is that the Single Layer Search tab has a query history. Clicking on the **History** roll-out menu, you can have a look at the previous queries made during that ongoing session. See the figure below.



**Figure 1.7. Query History in Single Layer Search**

Furthermore, this search tab offers different modes to restrict the search. The first mode allows you to choose the way in which the results are shown. There are four options:

• Annotation: the search string is part of or exactly matches an annotation.

• N-gram over annotations: each element of the search string (the elements must be divided by spaces) is part of or exactly matches each one of several consecutive annotations.

• N-gram within annotation: the various elements of the search string (the elements must be divided by spaces) are part of or exactly match one single annotation.

• **Annotation + extra tier**: with this option you can search for annotations in two tiers at the same time. Once selected all the necessary constraints, the first tier will be shown in the results, whereas the additional tier will appear into the information balloon when hovering over such results. When you choose this option, another menu will become available, including the following options:

– **Fully aligned, any tier**: the begin and end time of both annotations are the same. In this case you do not specify the type of the additional tier.

– **Overlapping, any tier**: part of both annotations overlap on the time axis. Here again you do not specify the type of the additional tier.

– **Fully aligned, sibling tier**: the begin and end time of both annotations are the same, but the two tiers in which the annotations are contained are siblings, i.e. they belong to the same parent tier.

– **Overlapping, sibling tier**: part of both annotations overlap on the time axis. The two tiers containing the searched annotations are siblings, i.e. they belong to the same parent tier.

– **Fully aligned, child tier**: the two annotations have the same time span, but the additional one, contained in the extra tier, is child of the other one.

– **Overlapping, child tier**: part of both annotations overlap on the time axis, but the additional one, contained in the extra tier, is child of the other one.

The second mode offers the straightforward distinction between case sensitive and case insensitive search.

The third mode allows you to choose whether the annotation found should contain the search string (substring match), whether the annotation should exactly match the search string (exact match) or whether some regular expression should be used in the match (regular expression) [for further information about

regular expressions see section 1.4 below]. This mode also contains another option, which is called **variable match**, but since it has little use here in the **Single Layer Search**, it will be thoroughly explained later on, when it comes to the **Multiple Layer Search**.

Finally, you can restrict the search to one tier in particular, choosing among tier types, tier names, participants or annotators (if present), all of which can appear in a decreasing order, from the most to the least used, if you check the option **tier choice sort**.

## Wild cards and negation

When you choose an N-gram to be the form of the result, you can use two more options: a wild card and/or a negation. The wild card takes the form of a # sign. For instance, the search string `the # man` with the mode N-gram over annotations would return three annotations per hit: the first annotation contains `the` (or exactly matches it, if the mode exact match is chosen), the second annotation may contain anything due to the use of the wild card, and the third annotation contains, or exactly matches, `man`. If the mode N-gram within annotation is chosen, each hit contains one annotation. This one, in turn, contains the search string with all the possible combinations due to the # sign.

If you want to find N-grams where a hit matches anything but one string in particular, you can use the negation operator NOT(...), entering in brackets the search string not to be matched. For instance, the search string `the NOT(strange) man` would return 3-grams in the same way as described above, but the hits where the second annotation matches, or contains, *strange* are left out.

# 1.3. Multiple Layer Search Tab

The Multiple Layer Search tab houses the most comprehensive search in Trova.

Similarly to the Single Layer Search tab, it has a query History menu, which allows you to go back to the previous queries you have made during that ongoing session. Also the mode case sensitive/case insensitive is the same as the one contained in the Single Layer Search tab.

The first new element is the Reset form button. Clicking this button will clear all the strings entered in order to make a query.

Besides, a new option has been included into the menu containing all the different types of matches (i.e. substring match, exact match, regular expression): **variable match**. As the name says, it has to do with mathematical variables, and it can be used every time you want to search for two or more annotations, contained in two or more different tiers, reporting the same text and/or the same time alignment. See the example in the figure below.



**Figure 1.8. Example Variable Match**

The buttons Minimal Duration and Maximal Duration enables you to constrict the minimal and maximal temporal duration of each result (please note in the window below that it says "Enter a Minimal/Maximal Duration for the complete pattern", which means that if you have filled in all the three strings, you will have to enter the whole duration, and not the duration of one single annotation). When you click on one of the buttons, a dialogue window appears:



**Figure 1.9. Minimal Duration**

Here you can enter the minimal and/or maximal duration as the total number of milliseconds or in the form of hours:minutes:seconds.milliseconds. A value of 0 milliseconds or 00:00:00.000 yields as undefined. Searching for annotations with a maximum duration being less then the minimum duration is impossible. Hence, entering conflicting values results in an error message saying that the combination is impossible. After entering a correct duration, this will be displayed in the corresponding button.

The buttons Begin After and End Before give a dialogue window similar to the one seen above. They allow you to restrict the annotations in the results so that they begin after a certain time and end before a certain time. Entering a Begin After time that is higher than the End Before time results in an error message saying that the query made is impossible. After entering a correct time, this will be displayed in the corresponding button.

# Search string and constraints

Beneath all the buttons and functions discussed above, you will find a table consisting of white and green fields. Search strings are entered in the white fields, while a green field between two non-empty white fields must contain a constraint.



**Figure 1.10. Multiple Layer Search**

The fields on one row have to do with the search strings and constraints that have to be matched by annotations contained in one tier. The reason for having three rows in the query table is that the search engine may find annotations contained in three tiers as one hit. Furthermore, it is possible to restrict the search to one type of tier for each row by choosing the appropriate option in the pull-down menu on the right of each row.

Let us first take a look at search strings and constraints in one row. If you enter two search strings in two white fields separated by a green field, you must fill in that green field, i.e. choose a constraint. Double clicking on the green field will open a context menu offering the following constraints:

- **= N annotations**: between the annotations containing the two search strings, there must be exactly N annotations.

- **> N annotations**: between the annotations containing the two search strings, there must be more than N annotations.

- **< N annotations**: between the annotations containing the two search strings, there must be less than N annotations.

- **= X milliseconds**: between the annotations containing the two search strings, there must be exactly X milliseconds.

- **> X milliseconds**: between the annotations containing the two search strings, there must be more than X milliseconds.

- **< X milliseconds**: between the annotations containing the two search strings, there must be less than X milliseconds.

- **No constraints**: there are no constraints between the annotations.

- **Clear**: it clears the constraint previously chosen.

When you click on Find and there is an empty constraint between two non-empty search string fields, you will get an error message. You will also get an error message if there is an empty search string field and empty constraint fields between two non-empty search string fields.

As we saw earlier, the search mechanism on this tab can construct a query within up to three tiers. Besides the constraints on the annotations present in one tier, you can also apply constraints on annotations contained in different tiers. This means that, taking into consideration the constraints entered in the green fields, the search engine will look for annotations which are contained in different tiers, and which match the search strings entered in two (or three) different rows of the query table.

As you can see in the figure above, besides choosing the type of tier, you can also choose among the following options, which are contained in the green drop down menus in between the tier types menus:

- **Must be in same file**: by choosing this option Trova will search for two tiers belonging to the same annotation file, but without a specified relationship between them.

- **Must be parent and child**: by choosing this option Trova will search for two tiers, one being child of the other. Please note that if you look at the annotations in the Timeline view mode in Annex, the parent tier is generally located in the lower positions, whereas when it comes to make a query, it has to go in the top search strings.

- **Must have same parent**: by choosing this option Trova will search for two sibling tiers, which means they belong to the same parent tier.

- **Must have same participant**: by choosing this option Trova will search for two (or more) tiers whose content has been orally told by the same person. The identity (and the identifying number) of the participants can be seen on the IMDI Browser page once you have selected the required node(s).

    Multiple constraints can be selected from this menu. All you have to do is holding down the Ctrl key on the keyboard while selecting the needed items.

Double clicking the green field between two search strings opens a context menu with the following constraints:

- **Fully aligned**: the begin time and end time of both annotations are the same;

- **Overlap**: part of both annotations overlap;

- **Left overlap**: the begin time and end time of the annotation matching the lower search string lie *before* the begin time and end time of the annotation matching the upper search string;

- **Right overlap**: the begin time and end time of the annotation matching the lower search string lie *after* the begin time and end time of the annotation matching the upper search string;

- **Surrounding**: the begin time of the annotation matching the lower search string lies before the begin time of the annotation matching the upper search string; and the end time of the annotation matching the lower search string lies after the end time of the annotation matching the upper search string;

- **Within**: the begin time of the annotation matching the lower search string lies after the begin time of the annotation matching the upper search string; and end time of the annotation matching the lower search string lies before the end time of the annotation matching the upper search string;

- **No overlap**: the begin time of the annotation matching one of the search strings lies after the end time of the annotation matching the other search string;

- **begin time - begin time = X milliseconds**: the begin time of the annotation matching the upper search string must lie exactly X milliseconds before the begin time of the annotation matching the lower search string;

- **begin time - begin time < X milliseconds**: the begin time of the annotation matching the upper search string must lie less than X milliseconds before the begin time of the annotation matching the lower search string;

- **begin time - begin time > X milliseconds**: the begin time of the annotation matching the upper search string must lie more than X milliseconds before the begin time of the annotation matching the lower search string;

- **begin time - end time = X milliseconds**: the begin time of the annotation matching the upper search string must lie exactly X milliseconds before the end time of the annotation matching the lower search string;

- begin time - end time < X milliseconds: the begin time of the annotation matching the upper search string must lie less than X milliseconds before the end time of the annotation matching the lower search string;

- begin time - end time > X milliseconds: the begin time of the annotation matching the upper search string must lie more than X milliseconds before the end time of the annotation matching the lower search string;

- end time - begin time = X milliseconds: the end time of the annotation matching the upper search string must lie exactly X milliseconds before the begin time of the annotation matching the lower search string;

- end time - begin time < X milliseconds: the end time of the annotation matching the upper search string must lie less than X milliseconds before the begin time of the annotation matching the lower search string;

- end time - begin time > X milliseconds: the end time of the annotation matching the upper search string must lie more than X milliseconds before the begin time of the annotation matching the lower search string;

- end time - end time = X milliseconds: the end time of the annotation matching the upper search string must lie exactly X milliseconds before the end time of the annotation matching the lower search string;

- end time - end time < X milliseconds: the end time of the annotation matching the upper search string must lie less than X milliseconds before the end time of the annotation matching the lower search string;

- end time - end time > X milliseconds: the end time of the annotation matching the upper search string must lie more than X milliseconds before the end time of the annotation matching the lower search string;

- No constraint: there are no constraints between the annotations or the tiers.

- Clear: it clears the constraint previously chosen.

Because the search mechanism offers the possibility to search for patterns in three tiers and there are possibly three search strings per tier, the search results also consist of nine elements per hit. The results are presented in the form of a table. In the figure below you can see an example of a query and its results. The upper row is the parent tier, the middle row is the English gloss tier, and the lower row is the Portuguese gloss tier.



**Figure 1.11. Results of Multiple Layer Search**

Here in the **Multiple Layer Search** you can sort out the tier lists in two ways: 1) **alphabetically**; 2) **by number of matches** (the matches are presented in a decreasing order). In both cases, the tiers will be ordered as follows: tier type, tier name, participant (and annotator, if present).

In addition, on the right of the **Font** menu there is an option called **Show time alignment**, which appears ONLY when the search has been made over TWO (or more) different tiers. Once you have checked this option, you will see, above each one of the results, two blue bars (graphically representing the time span of the annotations searched) plus the begin and end time of the annotations, and their duration in milliseconds. See figure below.



**Figure 1.12. Show Time Alignment**

As soon as this option is checked, another drop down menu will open up, named **Scale**, with the item **Stretch to fit** as default option. The latter means, on the one hand, that the two blue bars will stretch as much as the page allows them to, and on the other hand that each of the result will show its own begin and end time, and its own duration time in milliseconds. If you change the scale to, for example, 10 seconds, you will see that the two bars will have changed their length accordingly (thus making the visualisation easier and quicker), that all the results will report a complete duration of ten seconds, and that the end time of the annotations will correspond to the begin time plus ten seconds. See figure below.

**Figure 1.13. Time Alignment - adjusted scale**

Finally, on the left of the option **Show Time Alignment**, you can click on **Info View Options** and decide what has to be shown or hidden in the balloon appearing when hovering over the results, in the purple column next to the results, or above the results (i.e. the time options), simply by checking or unchecking the (un)required information. See figure below.



**Figure 1.14. Info View Options**

Search hints

If you would like to use both Exact match and Substring match in one query, use the Regular expressionoption. In places where you would like to have an exact match use the ^ and $ signs to match the beginning and the end of a string (e.g. ^of$), otherwise just enter a word for the substring match.

Wild card. Instead of using the # as in the Single Layer Search, you can use the regular expression .+ to indicate any character (the dot) one or more times (the plus sign). The NOT(...) construction on the other hand can be used in the Multiple Layer Search in the same way as described in section 1.2.

One final, but not less important, remark concerns the placing of relatively restrictive search strings. As we saw earlier, the hierarchy of the rows in the query does not reflect the hierarchy in the data. While this is perfectly true, we advise you to place restrictive search strings in the left most field on in the upper most row possible and the least restrictive search string in the right most field of the lowest row possible. The reason for this is the order in which the search engine considers the search strings in the query. If it finds a restrictive search string it can filter out all the other possibilities, but if it finds a less restrictive search string it has to consider all the matches of this search string. Because of this, the search might take much more time if the non-restrictive strings are placed before the restrictive ones.

# 1.4. Regular Expressions

Regular expressions allow users to create complicated queries. Below follows a list of most commonly used regular expressions together with explanations and some potential uses.

- [abc] means "a or b or c", e.g. query "[br]ang" will match both "adbarnirrang" and "bang"

- [^abc] means "begins with any character but a,b,c", e.g. query [^aeou]ang will match "rang" but not "baang"

- [a-zA-Z] means "a charater from a/A through z/Z", e.g. b[a-zA-Z] will match "bang", "bLang" or "baang" but not "b8ng"

- . (the dot) means "any character", e.g. "b.ng" will match "bang", "b8ng", but not "baang"

- X* means "X zero or more times", e.g. "ba*ng" will match "bng", "bang", "baang", "baaang" etc.

- X+ means "X one or more time", e.g. "ba+ng" will match "bang", "baang" but not "bng"

- ^ means "the beginning of the annotation", e.g. "^ng" will match "ngabi" but not "bukung"

- $ means "the end of the annotation", e.g. "ung$" will match "bukung" but not "ngabi"

Examples

– ^[pbtd][^aeiou]

You can use this expression to search for complex onsets. It will find words that start with one of the plosives ("p","b","t","d") followed by a character that is not a vowel ("a","e","i","o","u"). An example of a matching word is "tsakeha"

– [^n]g$

You can use this expression in case you want to search for annotations ending with a "g", but not with "ng". In Dutch, you will find "snelweg" and "maandag" as the results but not words as "bang".

– ^k.+k$

You can use this expression if you want to search for annotation starting and ending with "k" and with one or more character between them, e.g. "kitik" or "kanak-kanak"

– ^(.+)\1$

You can use this expression to search for words that are reduplicated. When you put something in bracketes, you create a variable (.+), which you can refer to as "\1". This expression then searches for an annotation that starts with one or more random characters followed by that same sequence of characters. This expression will match for instance "kulukulu".

**More about regular expressions...**

The following tables have been created by a user of ELAN (an annotation tool which has the same search mechanism as TROVA). They may result quite useful also for other users since they offer a simple and clear overview of the main symbols (partly different from the ones just seen) used in regular expressions, with a short explanation and an example for each of them. Bear in mind that the examples are taken from the language that the user is being researching, so do not pay attention to the meaning of the words but to the working mechanism of the regular expressions.

## Table 1.1. Symbols

| Symbols | Place | Meaning |
|---|---|---|
| \b | at the beginning and/or end of a string | word boundary |
| \w+ | at the end of a string | variable end of word |
| . | anywhere | any letter |
| .* | between spaces | any string of letters between spaces/any word |
| .*\ | between spaces | any string of words |
| (x\|y) | anywhere | either x or y |
| [^x] | place at the beginning | not x |
| (....)\1 | anywhere | words with four reduplicated letters |
| ? | after a letter | the preceding letter is optional |

## Table 1.2. Search for particular word forms (examples)

| Symbols | Hits | Examples |
|---|---|---|
| sa | all words containing the string *sa* | *sa, vasaku, sahata, tisa* |
| \bsa | all words starting with *sa* | *sa, sahata, sana*; NOT vasaku, tisa |
| \bsa\b | all words *sa* | *sa* |
| \bsa..\b | all words consisting of *sa* + two letters that follow *sa* | *saka, saku, sana* |
| \bsa\w+ | all words beginning with *sa*, but not the word *sa* by itself | *sahata, sana* |
| \b.*ana\b | al words ending in *ana* | *sinana, tamuana, sana, bana, maana* |
| (....)\1 | all words with four reduplicated letters | *pakupaku, vapakupaku, mahumahun, vamahumahun* |
| \b(....)\1 | all words beginning with four reduplicated letters | *pakupaku*; NOT vapakupaku |
| \b(....)\1ana\b | all words beginning with four reduplicated letters and ending in *ana* | *vasuvasuana, hunuhunuana* |
| \bva(....)\1 | all words consisting of the prefix va- + four reduplicated letters | *vapakupaku, vagunagunaha* |

| Symbols | Hits | Examples |
|---|---|---|
| \bvahaa?\b | all tokens of vahaa and vaha | *vahaa* and *vaha* |

## Table 1.3. Search for particular sequences of words (examples)

| Symbols | Hits | Examples |
|---|---|---|
| \bsaka\b .* \bhaa | string of 3 words: (1) *saka*; (2) any word; (3) the word *haa* by itself or with suffixes | *saka antee haa*; *saka abana haari*; *saka kabuu haana* |
| saka .* \bhaa\w+ | string of 3 words: (1) *saka*; (2) any word; (3) a word beginning with *haa*, but NOT the word *haa* by itself | *saka abana haari*; *saka kabuu haana* |
| (\bsaka\b\|\bsa\b) \bpaku\b | 2-word string consisting of *saka* or *sa* and *paku* | *saka paku*; *sa paku* |
| (\bsaka\b\|\bsa\b) .* \bvaha\b | strings of 3 words: (1) *saka* or *sa*; (2) any word; (3) *vaha* | *saka tii vaha*; *sa tapaku vaha* |
| (\bsaka\b\|\bsa\b) (....)\l \bhaa | strings of 3 words: (1) *saka* or *sa*; (2) any word with four reduplicated letters; (3) the word *haa* or a word beginning with *haa* | *sa natanata haa*; *saka natanata haana* |