# Starting Guide iLCM

true

14 Juni, 2021

# Contents

# 1 Before the first use

Before starting the first program, the user can decide between loading the iLCM version from a standalone or locally. The standalone version is available at https://hub.docker.com/r/ckahmann/ilcm. To use the version locally, check that docker is installed (https://docs.docker.com/install/). Having done this, the command:

```
docker run -it -d -p 3838:3838 -p 8787:8787 -p 8983:8983 ckahmann/ilcm:latest
```

does pull the image from docker hub and starts a container hosting the needed services. The user can set port mappings to different ports on their local system. Per default, the 3 services are available on this port:

- Shiny-server: http://localhost:3838
- Rstudio-server: http://localhost:8787
- Solr-dashboard: http://localhost:8983

# 2 Getting started

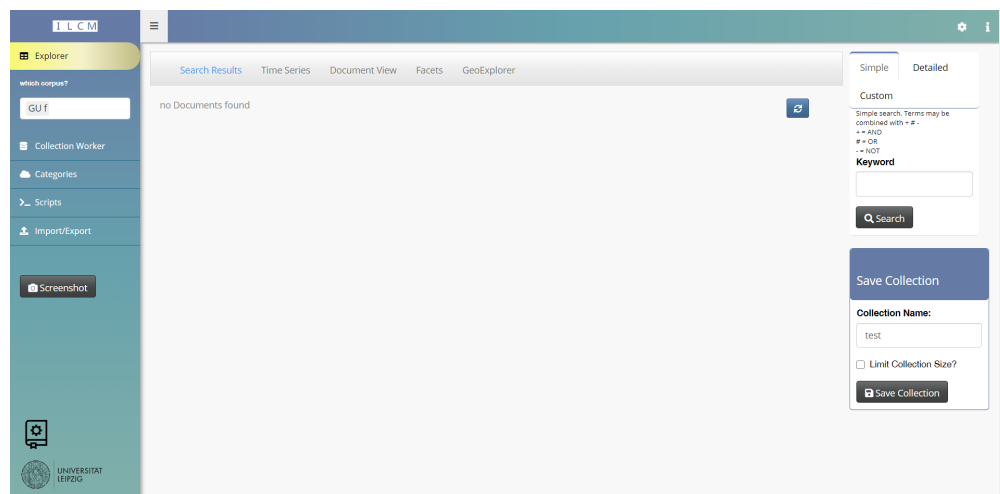After starting the iLCM-version the user should see the iLCM starting page:



Figure 1: Starting page of iLCM

# 3 Overall structure

The iLCM consists of 5 different basic tabs, which can be switched in the sidebar. These are *Explorer, Collection Worker, Categories, Scripts* and the *Import/Export.*

## 3.1 Explorer

In the *Explorer* tab, the user can search the database, get extra information about the search results, annotate text, and create a collection of documents, which serves as the starting point for further analysis.

## 3.2 Collection Worker

The *Collection Worker* is the tab where the user starts various NLP-algorithms. Also, the results of the started analysis can be accessed in the *Collection Worker*.

## 3.3 Categories

In the *Categories* tab, the user can have a closer look at the made annotations, create new annotation schemes or create further annotations.

## 3.4 Scripts

Here the user is offered five functionalities. First, if the user is familiar with R, they can change the NLP-algorithms' scripts. The user should be enabled to create different versions of the script. However,keep in mind that they have to produce the same output formats as the original script did. This is the case because the result of visualization is dependent on certain variables and variable names. The following options are given by the possibility of creating or manipulating blacklists, whitelists, dictionaries or vocabulary lists. These can be used in the preprocessing chain when applying different NLP-algorithms.

## 3.5 Import/Export

In the *Import/Export* tab, the user can upload their text data. The data will then be preprocessed and uploaded to the database and Solr. After running NLP-task or other operations, they can export collections, results, annotations or other files from the Solr database.

# 4 Import own Documents

The first thing the user might want to do is importing their data into the tool. For this purpose, there are different possibilities. The first option is importing their data as a csv-file, which can have any format. The second is importing multiple text files, like .pdf, .txt, .doc, .docx, where every file is imported as a single document. As another option they can also import REFI-projects or codebooks from other tools or upload the data directly to the database and Solr.

## 4.1   CSV-file

To import csv-files the user needs to select the csv-type in the *Importer* tab. After selecting this a *browser* bar should appear bellow. Here the user is able to search within their own directories for the csv-file they wish to upload to the tool.
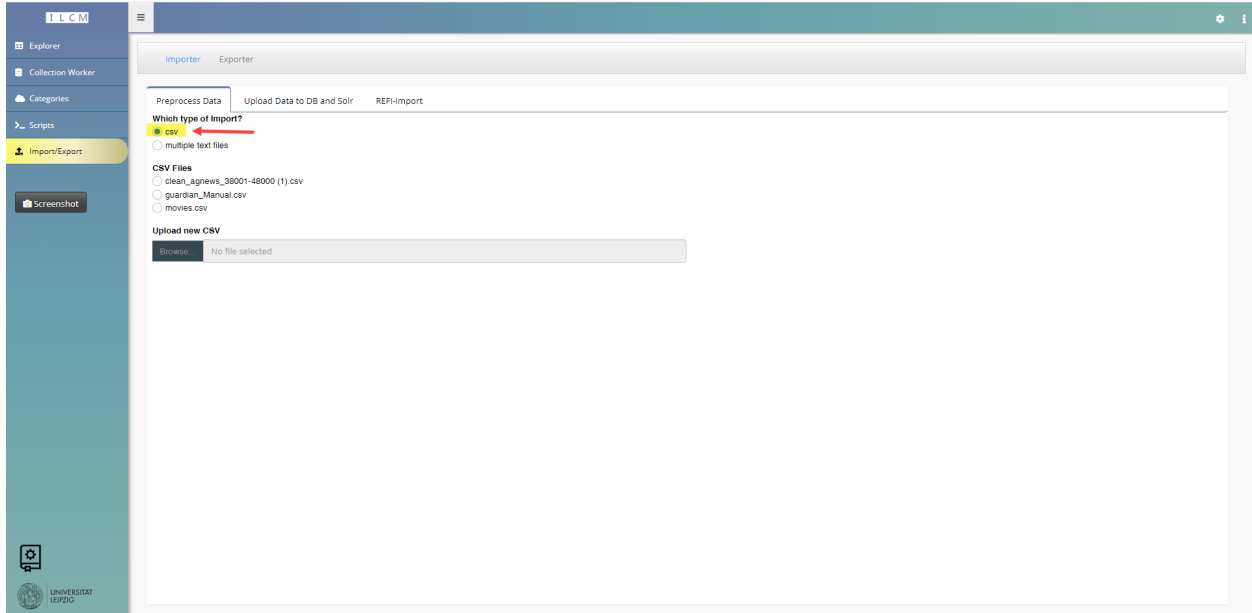


Figure 2: Import .csv-files

Once the file is loaded into the tool, the user can either check whether the data is imported correctly or directly start the mapping.
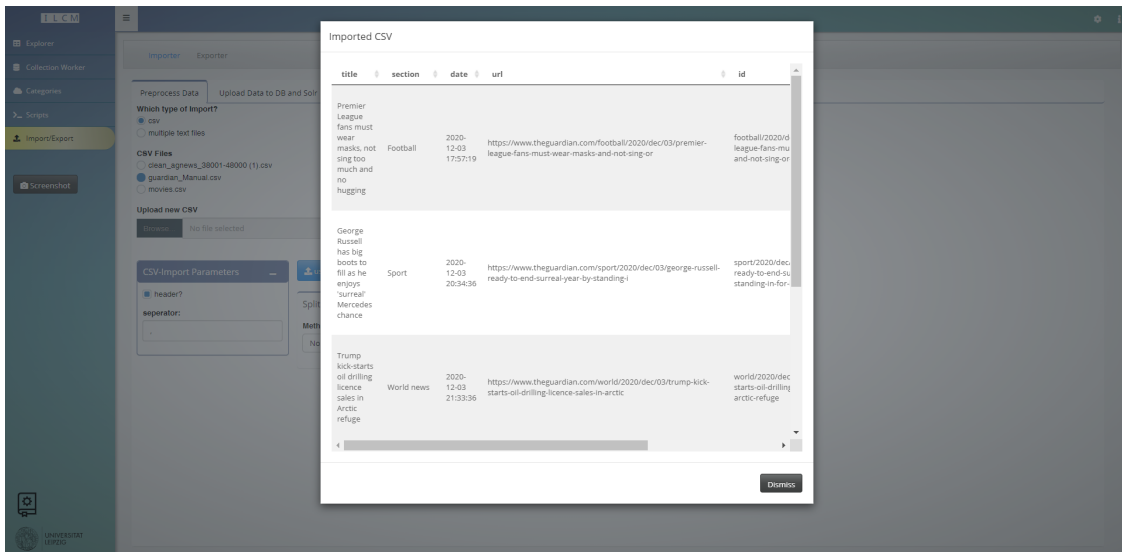


Figure 3: check if .csv-file is imported correctly

In the database, the texts are saved with certain metadata-fields. Very often, they will not perfectly match the metadata the user has. Therefore the user might want to abuse some given metadata fields for their purposes (e.g. publisher, section...). After clicking *Start Mapping*, the user can see two data tables. First, they need to set the mapping. In the header, the metadata-fields of the database are given. In the rows, the user can see the fields, which were found in the imported csv-file.



Figure 4: mapping of csv columns to iLCM database-format

Now it is the users' task to find the proper mapping. In this example, the user can see that the program can automatically guess the corresponding rows of the csv-file, but there is not guaranteed correctness. In map mde1 to mde9, the user can select how to name the columns and which of them should be included. Sometimes the user does not always have the metadata included in the csv-file. Then they have two other possibilities on how to get the data into the iLCM database: The first one is writing some R-commands in a short script for a particular database column.

```
1
2     #the vector values$Import_csv_title can be specified in this script
3     #if you want to use data from the imported csv file, you can use values$data_csv
4     #example:
5     # values$Import_csv_title<-paste0("Vortrag Nummer:",as.matrix(values$data_csv[,5]))
6     #or
7     # values$Import_csv_title<-rep("unbekannter Titel",dim(values$data_csv)[1])
8
9
10    values$Import_csv_title<-paste0("Rede Nummer:",as.matrix(values$data_csv[,"speechnumber"]))
```

save

Figure 5: use a short R-Script to input a metadata field

Here the column "title" is filled by using the word "speech number" pasted together with the values in the column "speech number" of the csv-file. For other metadata fields - like "publisher" or "type", it is sometimes the same value for the whole dataset.

If this is the case, the user can just click the "Type-button" to input the value by hand. This value will then be used for all documents.



**types (one for all)**

parlament protocoll

save

Dismiss

Figure 6: Type in a value for the certain database column

On the bottom, the user has to specify a language, the date-format and give an abbreviation for imported data. It is necessary that the abbreviation and the "id_doc" field together form a unique key. So when uploading multiple csv-files, which all belong to the same dataset (same abbreviation), the user needs to use unique id_doc values. In the second table below, the user sees the Metadata input file being created, dependent on the decisions he made in the data table above. Once the user is sure everything is set correctly, he has to click "Start Preprocessing and save csv-files" or "Start preprocessing and directly write to DB". In both options, the imported data are preprocessed using Spacy and then saved as two csv-files. After clicking the second button, the data gets imported into the database and Solr using the two created csv-files (token and meta).
Once the user restarts the app (open the link in a new tab), the data will be available and ready to work.

7

## 4.2 Multiple text files

The input for multiple text files is almost the same as for csv-files. The only big difference is the fact that for the text files only, the fields: "id_doc", "title", and "text" exist. The other missing columns have to be filled in by hand if the user wants to include them. Otherwise, he can also leave them empty ("id_doc" and "body" have to be non-empty).
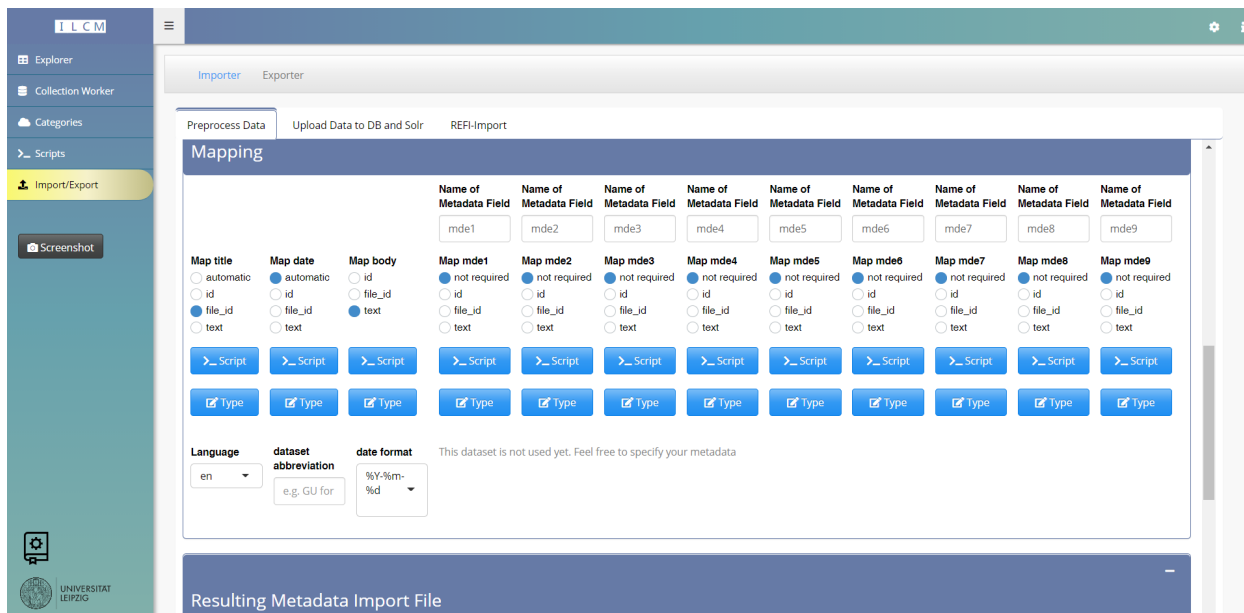


Figure 7: example input mapping for multiple text-files

Once the import is successful, the user can choose the dataset in the sidebar-panel below the explorer-button. It is possible to work on multiple datasets at the same time.

## 4.3 Upload Data to DB and Solr

If the user want to save the current imported data in the database they can go to the next subtab. Here they will see all imported corpora and are able to select one they wish to save in the current Solr-pipeline and the Maria-database. Check the options-menu on the top right of the side (gear-icon) to see if the connection to the MariaDB and SOlr are correct before proceeding to upload the data. For the upload simply click on the *Upload selected data to DB and import to Solr* button. Delete the data by selecting the corpus and clicking the *Delete* button below.

## 4.4 REFI-import

The user may want to import projects from other software. The REFI-format is a universal format to do so. Therefore the user can click on the corresponding subtab and search for a REFI-QDA Project file (.qdpx) or a REFI-QDA Codebook file (.qdc) in their personal data. On the right side the user also will see the datasets and schemes in the tool that also can be turned into a REFI-project for download. Importing a new project or codebook will update these two lists.
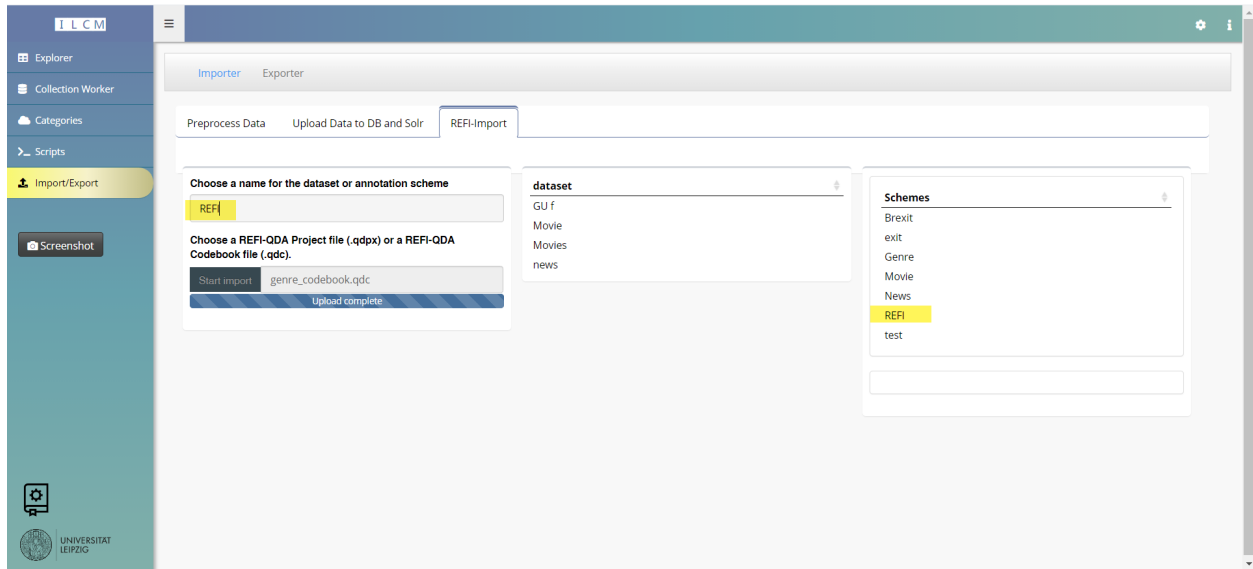
Figure 8: Importing an REFI-codebook under the name 'REFI'

# 5 Export

After finishing different tasks the user can decide to export the result-data. Different export formats allow the user to reuse the data in other NLP-programs or for visulization, documentation and other.

## 5.1 Collections

Here the user can download the created collections of documents (see section *Collections*). The user can decide whether the documents should be downloaded as Dataframes in a csv-file or as an RData-file. For the *Dataframes* selection, the user needs to select a batch size defining how many documents a csv-file should contain. Above the download button should appear a small text stating how many documents are in the collection and how many files will be generated for the selected batch size. Keep in mind that this will be a time-consuming process to generate all the files. The RData selection allows the user to download a huge file containing all the documents to the selection. This comes in handy when an experienced user wants to work with this collection in another R-program.
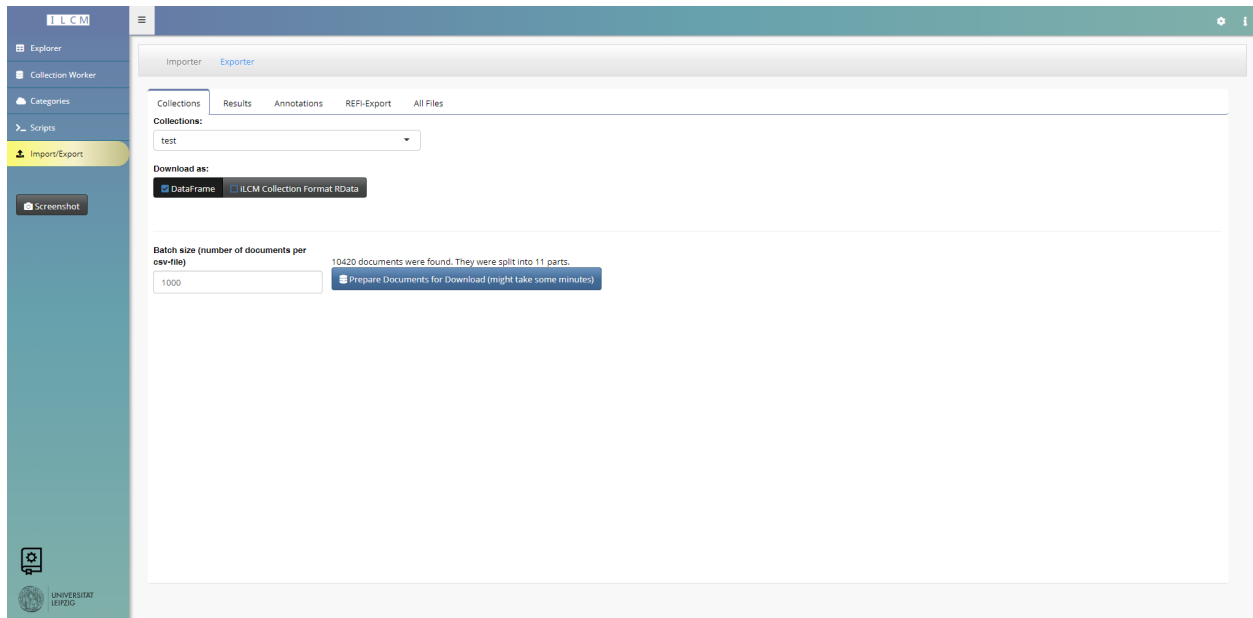
Figure 9: Datframe-export for a collection of 10420 documents and batch size of 1000 documents per csv-file

## 5.2 Results

Here the user has the option to download all the different results of NLP tasks from the tool. First, they can select the analysis task they wish to download by clicking on the corresponding field. Then they need to select a result set, where the name is generated through: *id* of the task, *name* of used collection and *time and date* when the task was started. After selecting the result set, a list with the downloadable data files should appear. Which exact files will appear depends on the selected analysis type. For example, in the *co-occurrence analysis*, there will be a file for the counts of co-occurrences, the dice parameter calculated for the co-occurrence, the log-file of the task, and so on. All these files can be downloaded by clicking on the blue download button in each file row.
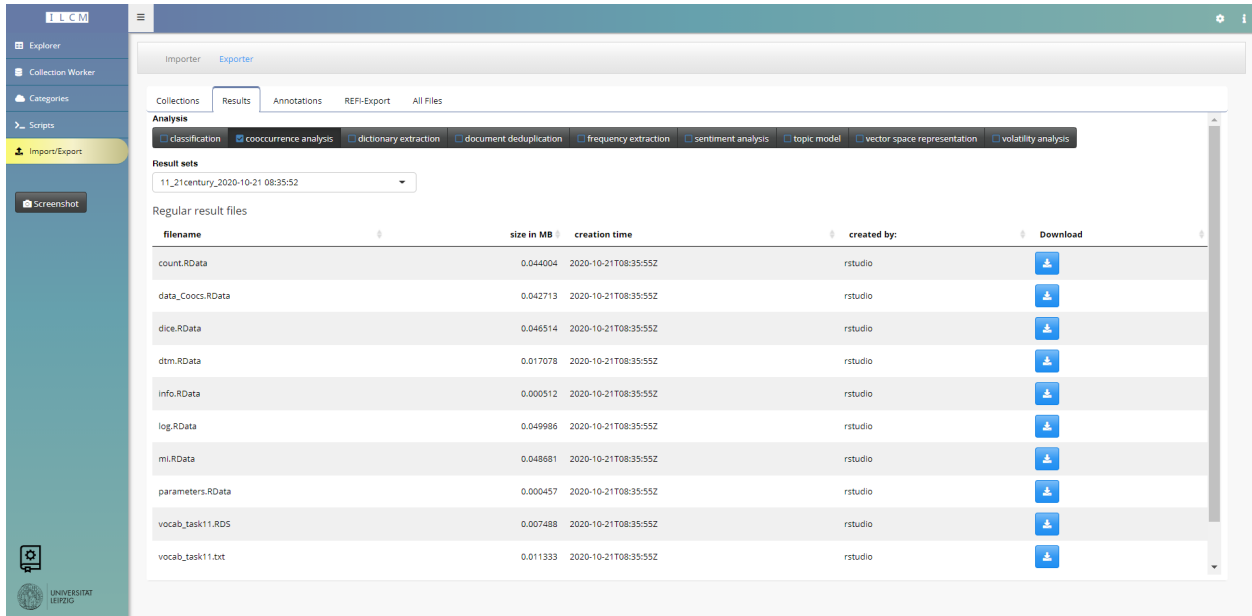
Figure 10: Result-export for a co-occurrence analysis of a movie collection

## 5.3 Annotations

Made annotations can be downloaded as well. Therefore the user can download all made annotations or filter the annotations within the list by simply searching for the annotations via one row. After the annotations where found, the user can decide whether they want to download it as a csv-file, excel-file, copy the list in the computer clipboard, or even print the annotations. Printing the annotation will allow the user to save the annotations as pdf-file as well.
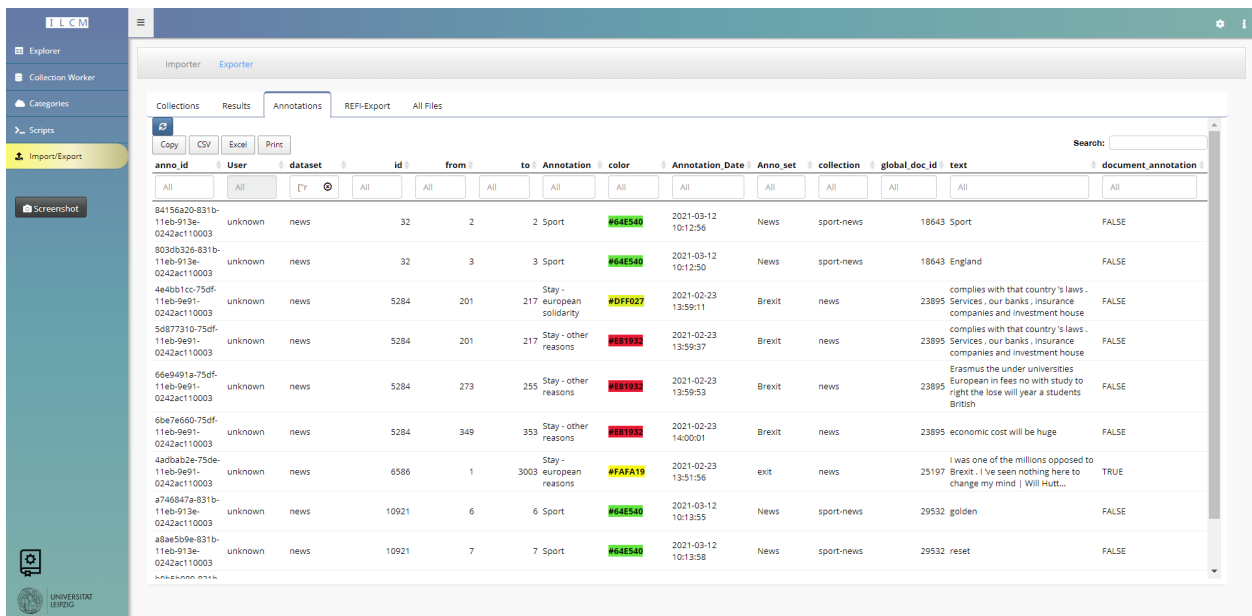


Figure 11: Annotation export for the news dataset

## 5.4 REFI-export

The REFI-export comes in handy when the user wish to transfer projects between software packages. Therefore the user can select a specific project and download it's documents or the made annotations. It's also possible to download a *Codebook* for the annotation projects in order to import it to another program. There is also the option to download the analysis results for the chosen collection.
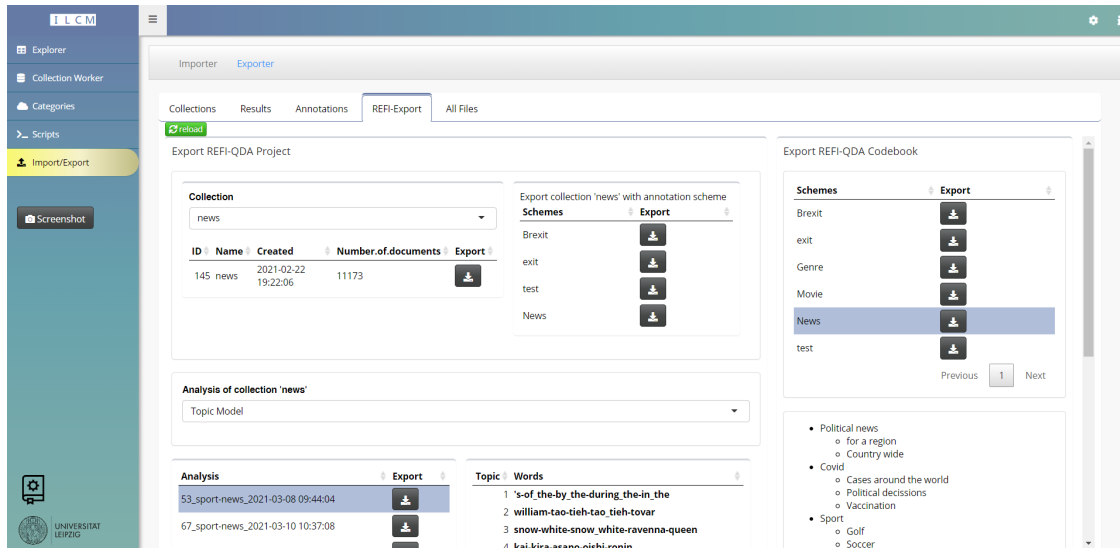


Figure 12: REFI-export for a news data collection

## 5.5 All Files

Here the user will see a *File select* button. After clicking on it, a window will open similar to an *Explorer* window on the PC. There the user can jump between different folders. The current folder direction is given on the top right of the window, below other folders will appear. In the example screenshot, the current folder is the *Results* folder containing all NLP-task results. This direction can be changed by opening the drop-down menu and selecting another folder. The folders inside the *Results* folder contain the folders of a single NLP-task of a certain kind (e.g. co-occurrence calculation). The single tasks are named with the id of the task and the name of the collection. Inside these folders, the user will find a list of files generated during the processing of the tasks. Here the user can decide which file they wish to download by clicking on the corresponding file name and the *Select* button on the bottom right of the window. If the user wants to select multiple files, they can do so by holding on to the CTRL button on their keyboard while selecting the file names with their mouse. After the file/files are selected, the small window will disappear, and on the main window, a table of all selected files is shown. Now the user needs to push the *Download* button.
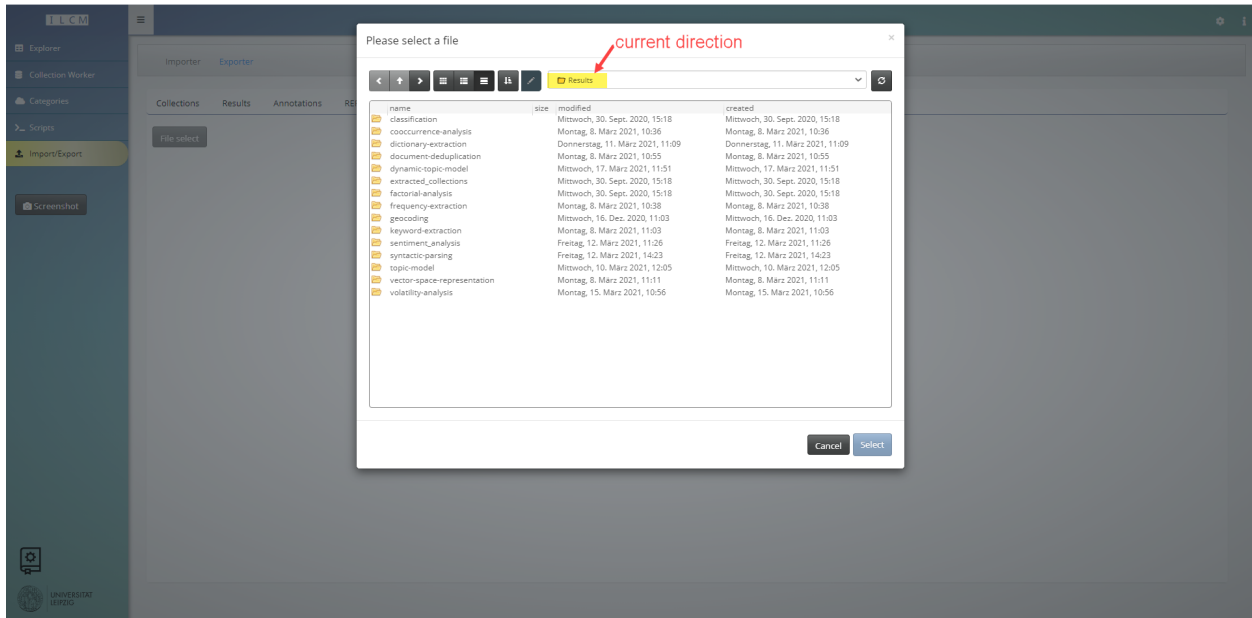
Figure 13: Explorer window to download all kinds of files from the iLCM-tool

# 6 Select imported dataset

After the user has imported a dataset, they can switch to the *Explorer* tab. There they will find under the *Explorer* tab a search field. Here the user can chose between the loaded datasets now called 'corpus'. After selecting one or more corpora they can start searching, building a collection and finally start the NLP-tasks.
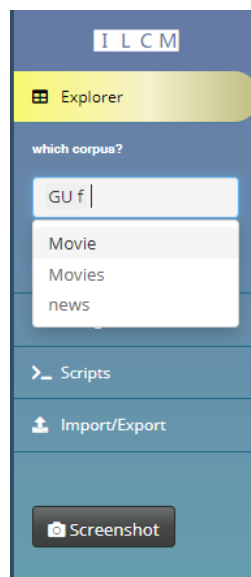


Figure 14: Selecting a dataset/ group of sets

## 6.1 Search Results

After the user has selected one or more corpora thy can see all documents in the *Explorer* tab under the *Search Results* subtab. They simply need to run an empty search request by clicking on the *search* button on the left site. Then a list will appear with all documents and their metadata information. On top of the page they user will also find the number of found documents. In order to clean the appearance of the list and focus on important information the user can change the column visibility. They therefore need to click on the *Columns visibility* button than a list of all column names will appear and by clicking on the column name the corresponding column will disappear in the documents list. To sort the documents the user can click on the column-names in the list and they will sort in ascending or descending order of the entries (e.g. clicking on the doc id will sort the documents from the smallest to the largest id or reversed).
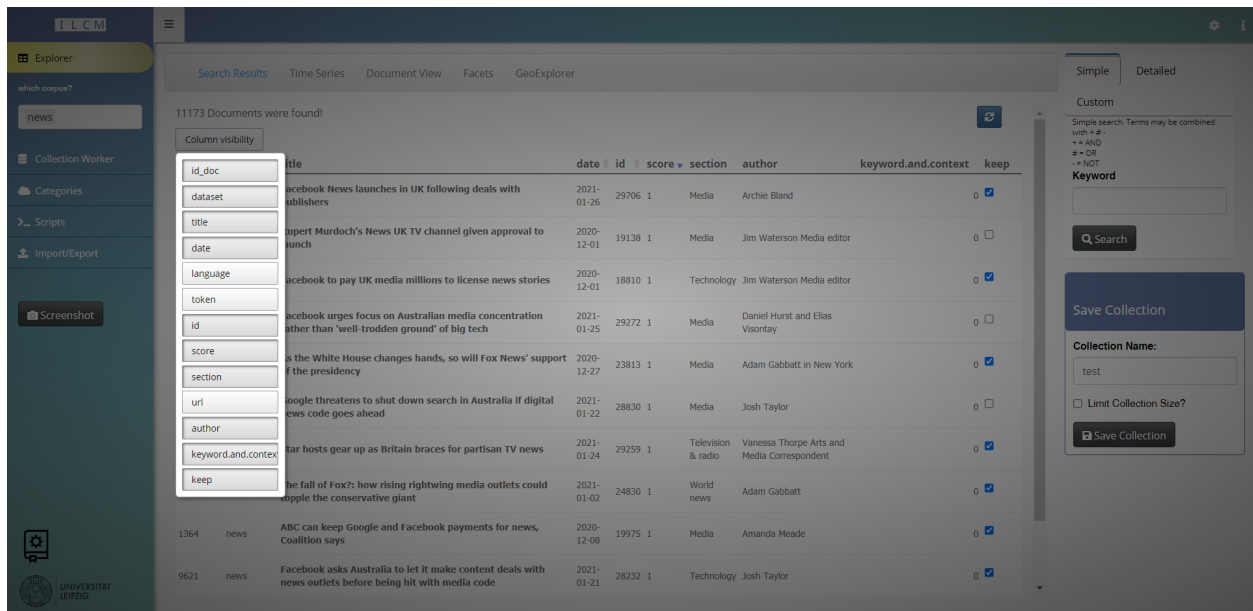


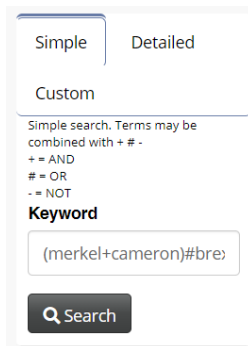Figure 15: Selecting column visibility

## 7 Search documents

One fundamental feature of the iLCM is examining vast amounts of text in a very fluent matter. The user can search for specific documents, have a closer look at some metadata for the search results and then save the documents they are interested in as a collection. When searching for documents, the user needs to make sure they are working on a suitable dataset. They can specify the dataset/datasets in the left top select-bar under the sidebar's Explorer-button. The two search options, "simple" and "detailed", the user gets suggestions while typing, based on the texts in the database.

## 7.1 Simple search

In the *Simple Search*, the user can enter search terms and combine them with the logical operators: and, or, not. The usage of brackets and quotation marks (for exact matches) for more complex queries is also allowed. The logical operators can be used by typing:

1. + = AND

2. - = NOT

3. # = OR



Figure 16: example query for simple search

The example query: *(merkel+cameron)#brexit* will return all documents in which the words *merkel* and *cameron*, or the word *brexit* is existent in the *body-*, the *title-* or in the *author* field.
The resulting documents are then listed in the *Search Results* tab. The user can see the resulting documents with metadata and some keyword and context based on the search terms. Below the search results, a slider enables the user to navigate through the other returned documents. On the complete right column of the search results table, there is a checkbox for every document, where the user can select whether he wants to use this specific document in the collection.

## 7.2 Detailed search

In the *Detailed Search*, the user has the same options as in the *Simple search*. On top of that, they can filter documents by their metadata. So the user can easily filter all the documents before a particular time or from a specific author or a certain section. The metadata filters are updated whenever the user changes the dataset/datasets they want to work on. The results are again displayed in the *Search Result* table.
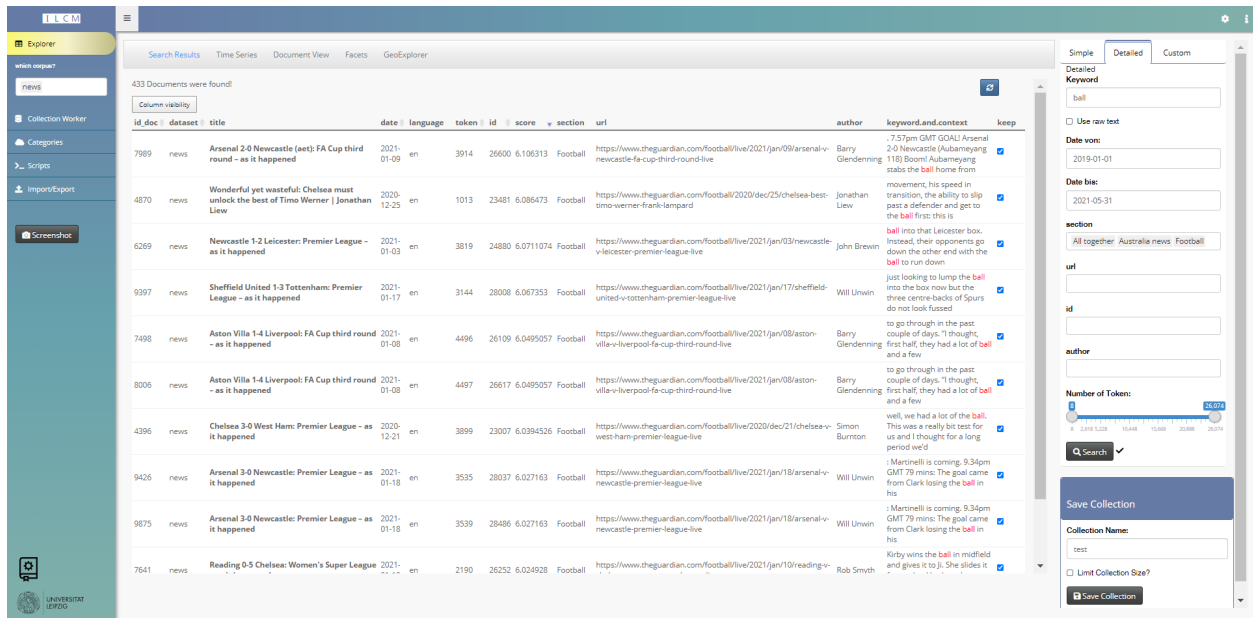
Figure 17: example query for detailed search and found results

## 7.3 Custom search

Solr is building the Search-backend. As a result, experienced users can use their own queries for Solr( Luscene query ). In the *Custom-search* tab also the sub-setting on special datasets has to be done by hand. This is not recommended for inexperienced users, however.

## 7.4 Time series

Once the users have started a search, they can look at the resulting documents. Furthermore, they can check the time-series metadata for the search query. In the *Time Series* subtab, there are two plots. In the first one, the user can see the number of documents found at a specific point in time. These plots can be compared to the last ten searches; the user has made. The time series can be plotted in days, month, and years. Also, the user can switch between relative and absolute numbers. They can download the time series data by clicking the *csv*-button. When the user clicks in the plot to a specific point in time, it will trigger a new search, whereby just those documents will be returned, which are released on the data the user has clicked on. The second plot is the calendar view. Here the user can also see the number of documents that were found for a specific day. This plot can be handy when there are some patterns in the time series data for a query (e.g. Christmas).
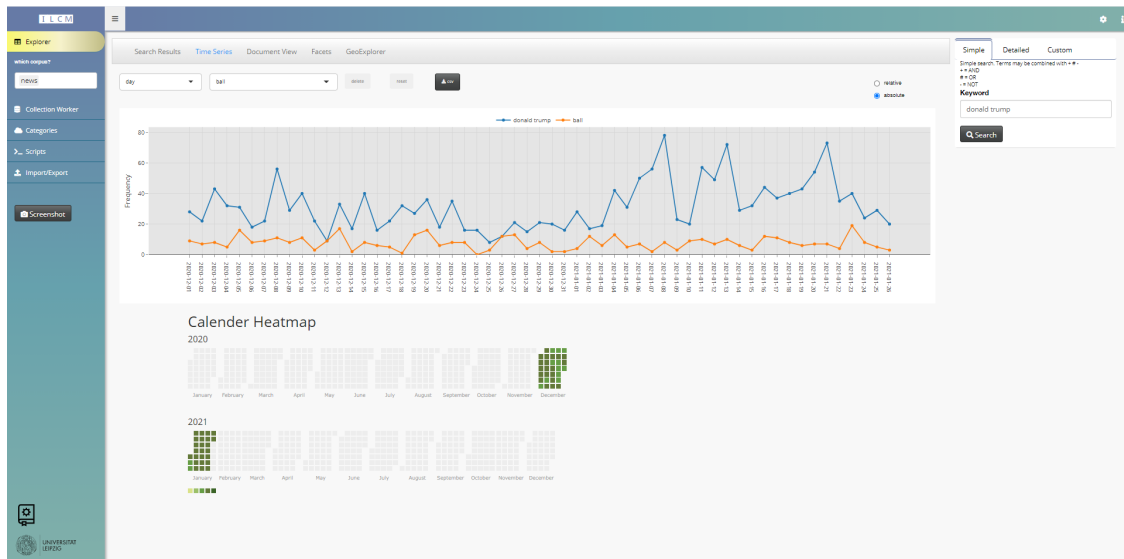
Figure 18: time series data for 2 searches

## 7.5 Facets

Despite the data of a document, other metadata might be of interest. In the *Facets* subtab, the user can see the top entries for the fields: *author, section, number of tokens,* and *entities.* For the authors, the user can set a threshold for an author's minimum occurrences to be displayed. The graph on the left side shows the absolute numbers ordered beginning from the highest values. The graph on the right side shows the relative values for the entities. In more detail, the user can see the percentage of occurrences for a particular entry, which are covered by the current search.
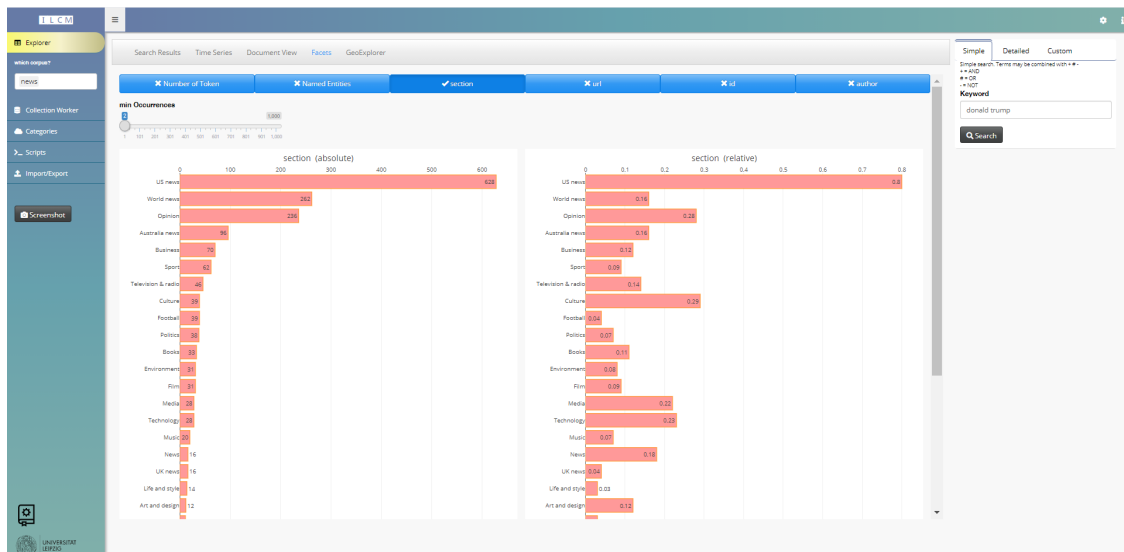


Figure 19: facet view for attribute 'section'

The number of documents is shown in a bar-chart, where the user can see the documents which exhibit a specific range of words for the current query compared with all documents in the database

(for selected datasets). The user can toggle between relative and absolute numbers of documents. Furthermore, he can switch to the logarithmic-mode for a better overview.
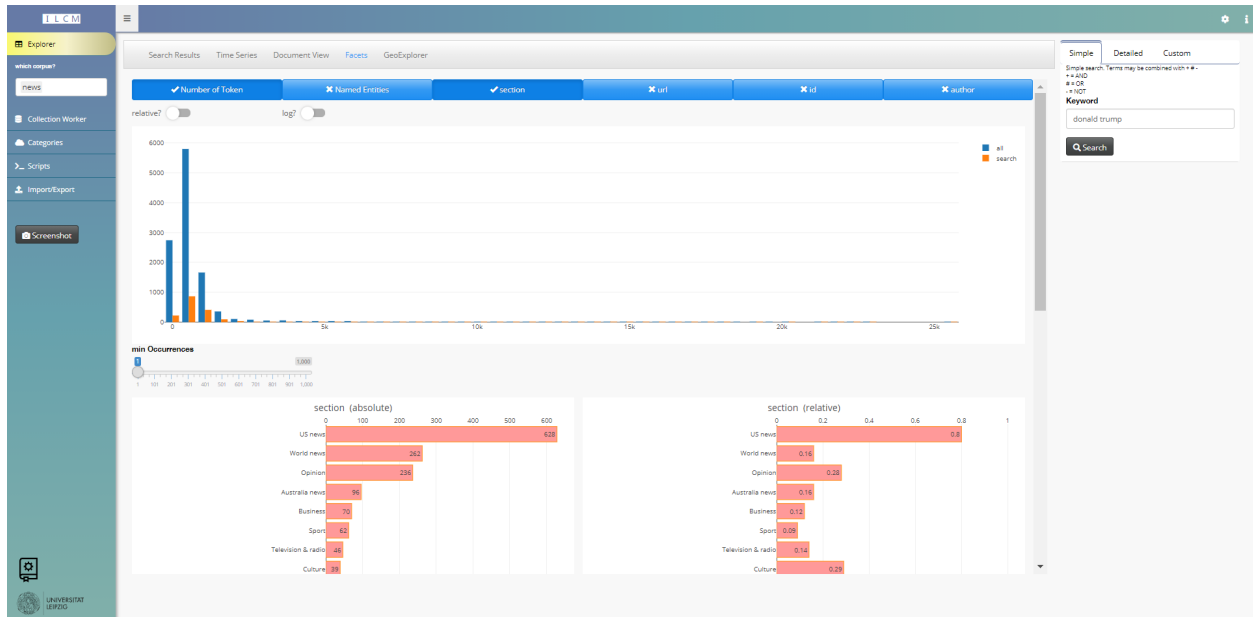


Figure 20: facet view for attribute 'number of tokens'

# 8 Options

On the top right corner of the iLCM tool there is a small gear-icon. Clicking on it will open the *Option* window. Here the user can see which languages are currently loaded in the tool and which other language models can be loaded to the tool. Below that they will see the connection points for the *MariaDB* and *Solr* database which will define where to save the result files and so on. Then the user has the option to delete some of the corpora which are imported to the system with the option to as well delete all the made annotations for this corpora. On the bottom on the page the user can sign-in to the system. This comes in handy when more than one person is working in one iLCM-tool and they want to see which person started which task or what kind of annotations a specific user has made. The last point on this page is covering the import case. There the user has the option to set boarders to the import function. When the connected database is spare on space the user can set the *max upload size* down in order to prevent the system from overloading the database. And they can also set the random seed for the different NLP-tasks in here, this is useful when it comes to comparability of tasks.
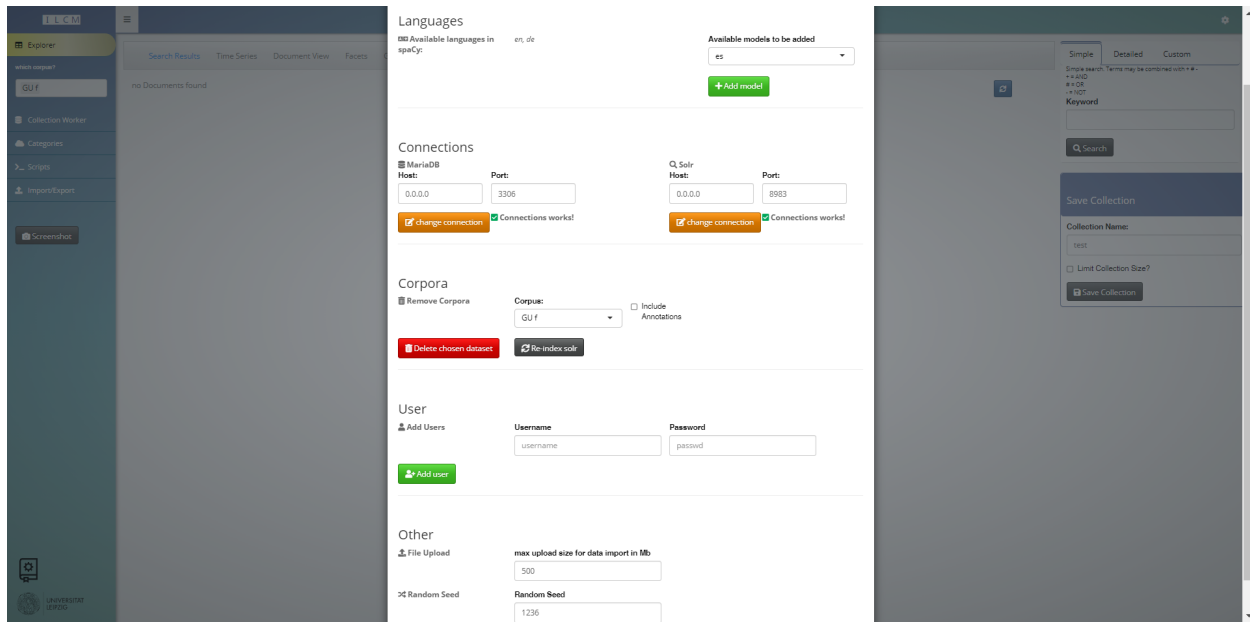
18

Figure 21: Tool-wide options

# 9 Document view

Of course, the app also allows working in a close reading approach. For this purpose, the user needs to select a document in the search-results. They can do so by clicking somewhere in the row of the document he is interested in. Next, they will automatically be forwarded to the document view. Here they can see three different boxes. In the middle, they can see the title and the whole text of the document. The metadata of the document is displayed on the top-left. Below, the metadata the user can select different POS- and NER-tags, which will be highlighted in the text. This might be helpful for some annotation tasks. The tags are calculated during the preprocessing using *Spacy*. On the top right side, the user can start a new search. If they want to go back to the results, they need to click the "Search Results" subtab in the navigation bar. Below the search-box, there is the annotation-box.
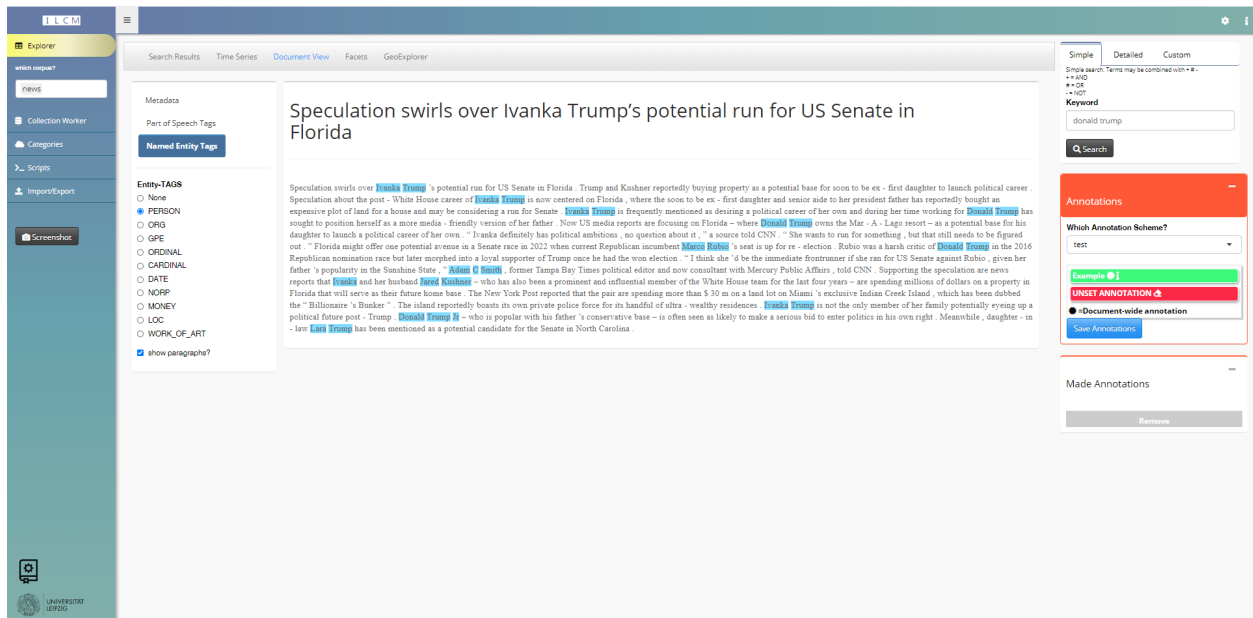
Figure 22: Document View with highlighted NER tag:'Person'

# 10 Geo-Explorer

To see the geographic information found in a certain collection at least one *Geocoding* task must be preceded for the particular collection. This task will search inside the documents and their meta data for geographic locations (like: 'Leipzig, Germany') and convert them into geographic coordinates. This allows to present the documents as points on a map corresponding which locations were found in the documents. This process is quite time consuming so its recommended to use a small collection for the *Geocoding* task as the process will search for all locations and will present them on a live map.

## 10.1 Configuration

First the user will set the configuration information right. Therefore they chose a collection and the corresponding *Geocoding* results they wish to view. After hitting the *Load data* button a couple of tables will appear on the right. These are used as parameters for filtering the results, or as stats for the calculation of numeric information/ distributions. The first table will allow the user to chose from the metadata information. Simply tick the corresponding fields in the row of a information from interest. If there are more than one value for the piece of metadata information the user needs to tick the *isMultiValue* column and select the used separator for them. The table on the bottom will allow the user to configure which *Geocoding* information should be used for the final visualization of the results. After setting all parameters up don't forget to hit the *apply settings* button on the bottom of the page.
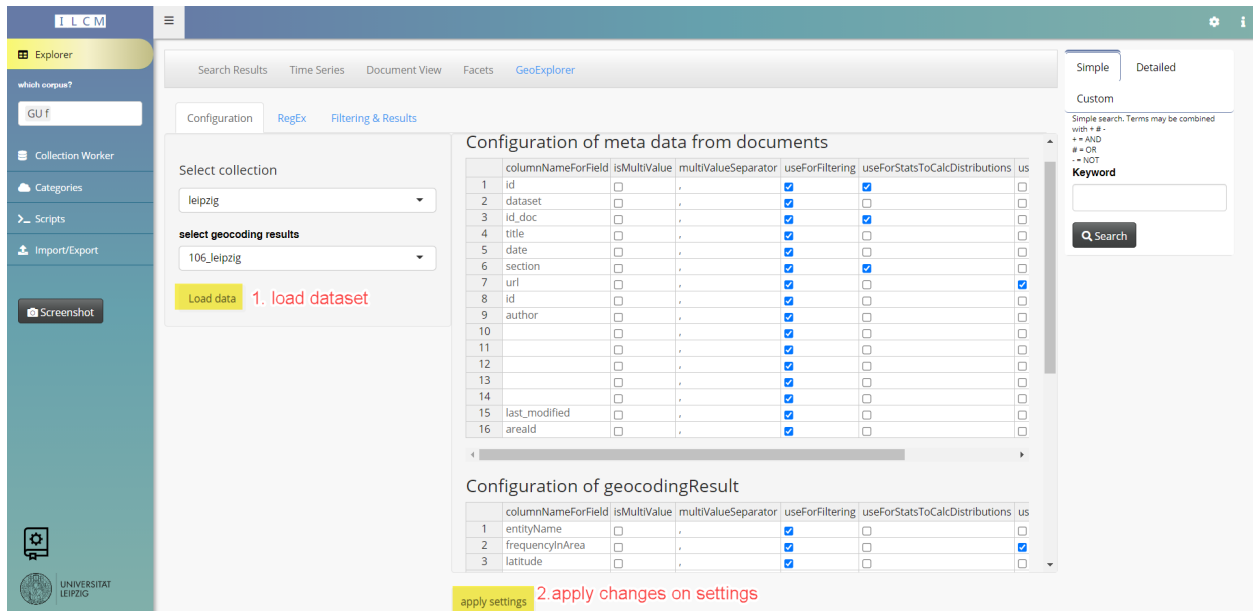
Figure 23: Setup the configuration for the Geocoding results

## 10.2   RegEx

The next subsection will allow the user to set regular expressions that should be matched to the *GeoCoding* results. The regular can help to find specific documents indicating a certain region. For example: If the user is looking for texts about *Great Britain*. Therefore documents containing words like *British, Brits, English, England* maybe indicators that the text is about *Great Britain*. So the regular expression indicating all of these words will look something like _Brit*_ and _Engl*_. To find more information on how to set up regular expression check out information on regular expressions. After choosing the regular expressions for the matching the user also needs to set the separator marks used for indicating decimal places and the separation of thousands digits. At lead the user can select parameters for the result presentation. After all parameters are set the user can click on *APPLY and perform Regex matching*. Then a counter will appear below informing about the number of found matches as well as a list with the results. If the settings should be canceled just press the *RESET* button.
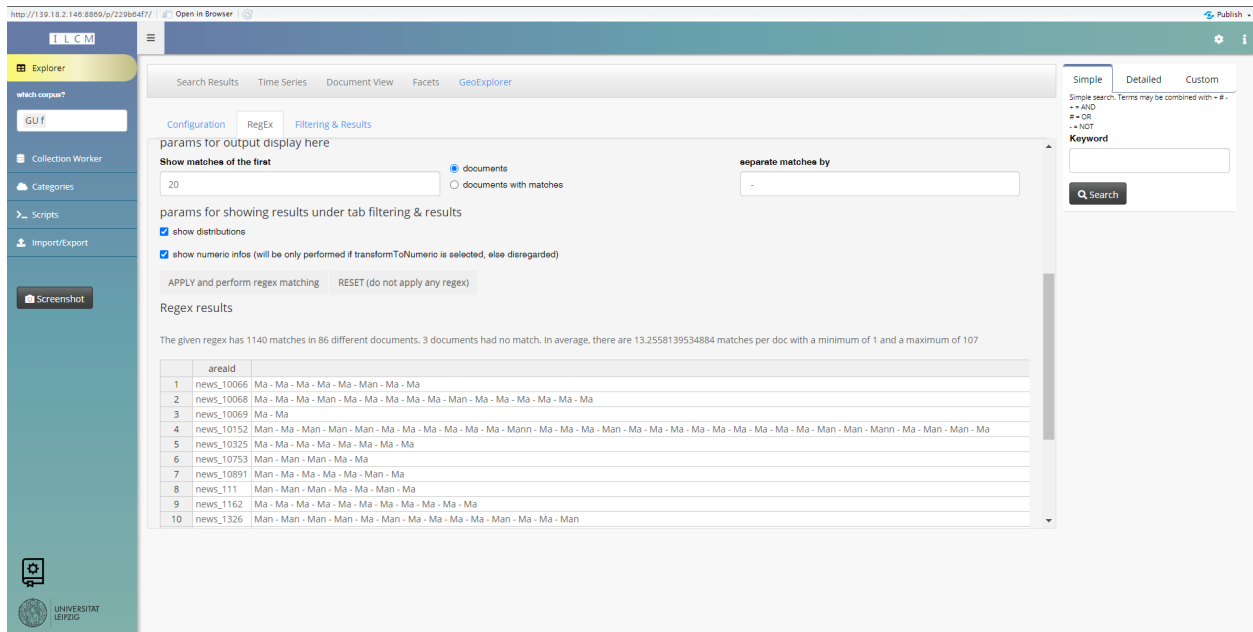
Figure 24: Found matches for the Regular expression 'Man*' in a sport themed collection

## 10.3 Filtering&Results

In this subtab all the premade settings from the *Configuration* subtabs will be visualized. This page consists of two parts: One to set filtering parameters for either metadata information, regular expression or the direct Geocoding results and one side showing the results of the filtering as an overview or on a map. In the overview part the user will first see a summary text of what was found considering the set filters. First the user will see the statistics for the metadata followed by the statistics of the Geocoding results and the ones for the regular expression task. The stats will consist of a diagramm and a table with further information. The map will show which locations were found by placing marks. Here the user can zoom in and out to see a more detailed view or a more general overview - the navigation works like every other map-tool.
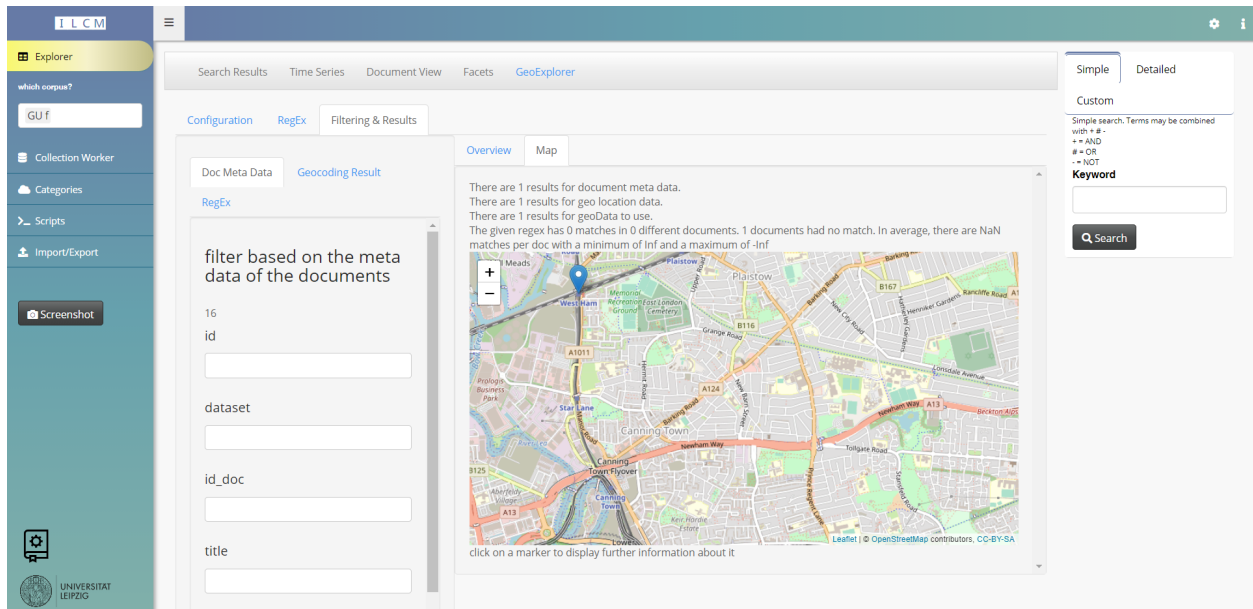
Figure 25: Map with locations of documents (closeup view)

# 11 Annotate documents

Using the iLCM allows the user to annotate their texts with predefined annotations schemes. These annotations might work as the basis for a classification-task.

## 11.1 Create annotation schemes

In the *Categories* tabs, the user finds the subtab *Create Annotation Set*. Here they can create new or manipulate existing annotation schemes. The annotation scheme can be hierarchical. An annotation class needs a name, a colour, a description, and it needs to be specified whether the annotation should target the whole document or just a specific part of the text.
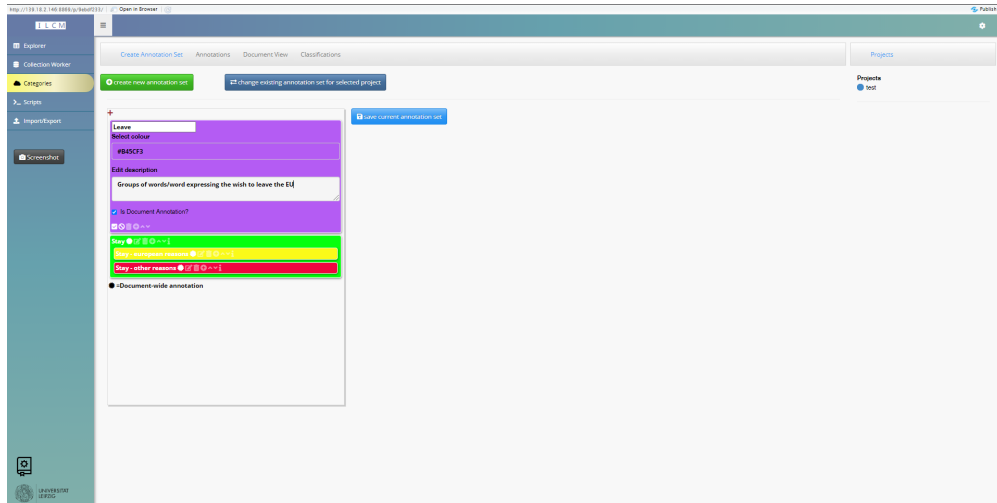
Figure 26: Subtab for annotation scheme creation

## 11.2 Annotate

Once there is an annotation scheme, the user can select it in the annotation box in the *Document View*. Afterwards, the annotation scheme will be displayed, as seen below. For annotation, the user needs to select the part of the text they want to annotate. The text will be highlighted. A click on the annotation class they want to use will then apply to the annotation; once the user is finished with annotating the document, they need to click the save annotation button. This will save the made annotations again. If the user is interested in seeing all made annotations or deleting several annotations, they can switch the *Annotations* subtab inside the tab *Categories*. Clicking on an annotation will then open the corresponding document. The user can also continue to annotate in this area.



Figure 27: Example for an annotation pro

24

## 11.3 Overview of annotations

After the user has made a couple of different annotation schemes and annotated some documents, they may want to get an overview of the annotations. Therefore the user can got to the *Categories* tab and select the *Annotations* subtab. Here they can between the 'projects' witch are build on the different annotation schemes. After selecting a project a list will appear depending of all made annotations. On top the user has the option to delete made annotations and search within the annotations or the categories.



Figure 28: List of all annotations from the project 'Genre'

## 12 Create collections

Having found the correct documents the user is interested in, they can save these documents in a collection. This is feasible by filling a name in the *save collection* box in the *Search-Results* subtab. This will save all documents in the current search result as a collection. Documents that should be excluded ("keep"-checkbox unchecked) will not be saved. Another option to limit the number of documents in the collection is to tick the *Limit collection size* in the *save collection* window on the bottom right. This will allow the user to select a limitaion method and a limitation number. To reconstruct the process of analysis, the steps of finding the documents and excluding some are saved for every collection in the *Create-collection* statement (hover collection name in subtab *Results* in *Collection Worker* tab). Each collection needs a specific name and should not include the sign: "__".

### 12.1 Create sub-collections

After the first collection was saved, the user can switch to the *Collection worker* tab. Here, they can select a collection by clicking on the collection name in the column on the right side. Further, the user can go to the *Documents* subtab to view the documents of the collection. Now they have

the chance to search inside this specific collection and create a sub-collection. After finishing the search task a list of documents will appear, which contain the searched phrases/word. The last column in the list will allow the user to decide whether they want to keep a certain document in the sub-collection or exclude it. Then the user needs to enter a name for the new collection and click on the save-icon. The new sub-collection will then appear in the *Collections* bar on the right side. Note that the made annotation will not be transferred to the new sub-collection.
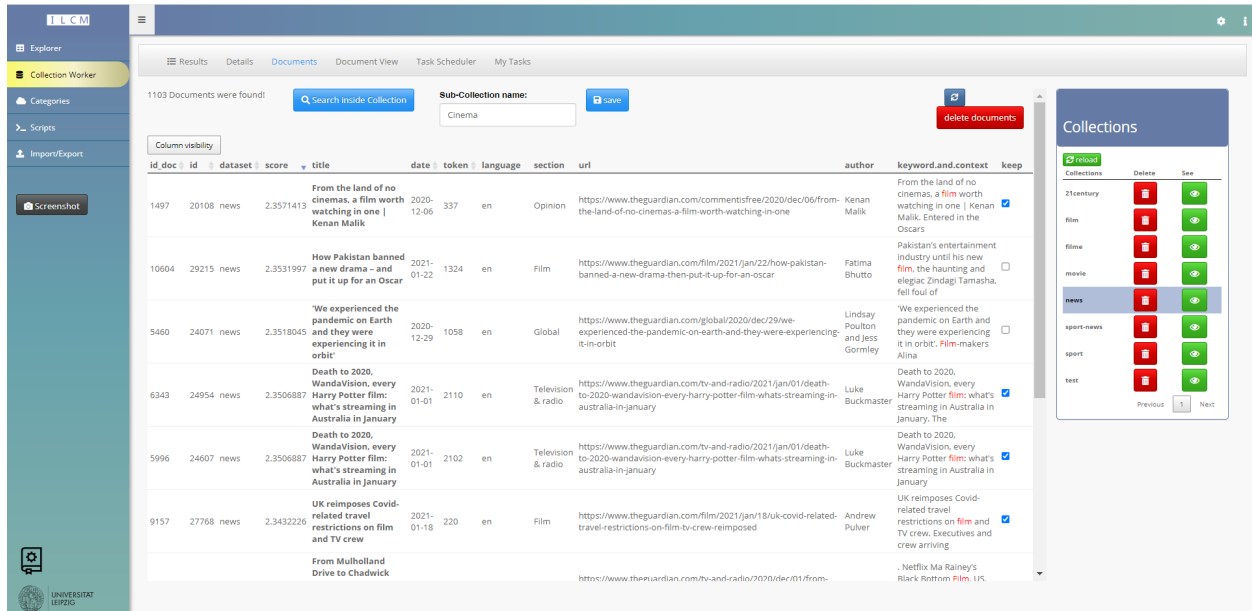


Figure 29: Subcollection of news articles after searching for the term 'film'

# 13 Start an NLP-task

Besides the search-tasks, the iLCM can also apply some of the latest NLP-algorithms to the saved collections. In the *Collection Worker* tab, there is a subtab called *Task Scheduler*. Initially, a collection and an NLP-task needs to be specified. The following tasks are implemented:
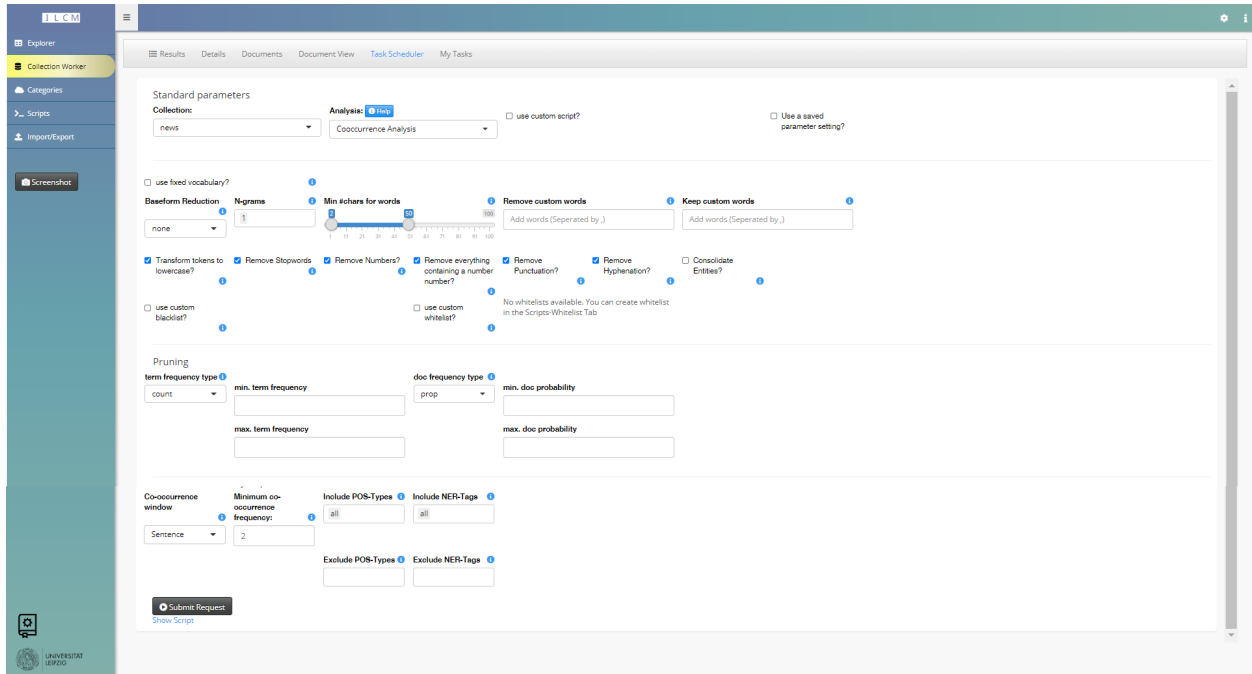
Figure 30: Task Scheduler subtab

- 1. **Co-occurrence Analysis** aims to find terms that appear significantly often together. These words can then be displayed in a network. E.g. a network of people appearing together in texts can be created.

- 2. **Frequency Extraction** counts the number of several dictionaries occurrences at specific points in time. A dictionary is composed of several words.

- 3. **Dictionary Extraction** counts the occurrence of several dictionaries at specific points in time. A dictionary is composed of several words.

- 4. **Volatility Analysis** (Context Volatility) measures the amount of context change for a specific word for a certain period.

- 5. **Topic Modeling** - A topic model is a type of statistical model for discovering the abstract "topics" that occur in a collection of documents.

- 6. **Classification** of the documents of the chosen collection depending on certain classification parameters, such as a given annotation set.

- 7. **Geocoding** find a potential location in the documents identifying where a particular document was written.

- 8. **Sentiment Analysis** finds word/phrases that lead to the reader to a certain emotion and calculates which emotion dominates a document.

- 9. **Document Deduplication** check for duplicates between the documents of a selected collection.

- 10. **Dynamic Topic Model** - dynamic modelling for "topics" that occur in a collection of documents.

27

- 11. **Keyword Extraction** - extract words that appear significant often in the documents of a specific collection.

- 12. **Syntactic Parsing** finds relationships between words in a sentence.

- 13. **Vector Space Representation** of the words/phrases of a collection for further work with word vectors.

- 14. Download Collection as CSV.

- 15. Save collection as a token-object.

- 16. Save collection as meta-object.

## 13.1   Standard parameter

The NLP-algorithms may need preprocessing before working adequately. This might be due to reasons of calculation time or achieving better results. There are some standard parameters the user can set for most algorithm. These parameters are:

- **Use fixed vocabulary**: Some tasks allow the users to decide if they want to use a fixed vocabulary in the analysis. Here the user has to enable the use of the fixed vocabulary and then chose one of the existing vocabulary lists. (Note that when the fixed vocabulary list is excluded, other parameters disappear as they are no longer relevant for the analysis).
- **Baseform reduction**: The users can specify whether they want to use lemmatization, stemming or no baseform reduction mechanism
- **Min number** of chars per word: The user can specify the minimal and maximal amount of characters a word needs to have. This helps immensely to get rid of errors in the tokenization.
- **N-grams**: The user can specify whether they want to work with uni-/bi- or trigrams. Keep in mind that using bi-/trigrams will increase the vocabulary size, leading to a more extended and more resource-intensive calculation.
- **Remove custom words**: The user can enter some words they want explicitly to be removed from the calculation. If there are more than a few words, it might be more useful to use the blacklist parameter instead.
- **Remove stopwords**: The user can enable/disable the stopword removal. The standard *quanteda* stopwords are used.
- **Transform to lower case**: The user can specify whether they want to transform all words to lowercase.
- **Remove numbers**: The user can enable/disable the removal of numbers
- **Remove punctuation**: The user can enable/disable the removal of punctuation (e.g.: ",", ";", ".", "-", "*", "#",...).
- **Remove hyphenation**: The user can enable/disable the removal of hyphenation.
- **Consolidate entities**: During the preparation of new data, named entities are extracted using Spacy. These might be composed of multiple words. E.g. the two words. *Angela* and *Merkel* from the entity "person". Consolidating the entities describes the process of binding words that form an entity together. These related words will then be used as a single word ("Angela_Merkel").
- **Use custom blacklist**: If the user wants to exclude multiple words from the analysis (e.g. bigger stopword lists, known words fragments in the data, ...), they can create a blacklist and

use it during the preprocessing. Here the user has to enable using a blacklist and then choose one of the existing blacklists.

- **Use custom whitelist**: If the user wants to include multiple words from the analysis (e.g. bigger lists of important words, known fragments in the data, ...), they can create a whitelist and use it during the preprocessing. Here the user has to enable the use of a whitelist and then choose one of the existing whitelists.
- **Pruning**: The pruning aims to remove words from the analysis whose frequency/document frequency is under or over a certain threshold. We use the pruning parameter setting of quanteda (quanteda pruning).
- **Include POS- or NER-tags**: In most of the NLP-task a POS (Part of Speech)- and NER(Named Entity Recognition)-tagging is necessary in order to complete the task. Therefore the user can choose which Tags should be considered to be important within the analysis. As an example it would be useful when the user wants to create *Topic Models* where the topics just consist out of nouns (POS-tag) with a specific NER-classification (e.g just mentioned locations).
- **Exclude POS- or NER-tags**: In contrast to including specified POS- or NER-tags the user can also exclude them. For example when the user wants to start a *Co-occurrence Analysis* the co-occurrence between words and punctuation marks (POS-tags) are not interesting therefore they can exclude them. Or they maybe are not interested in occurrence of firm-names (NER-tag) and therefore they can exclude them.

## 13.2 Create a blacklist

The blacklist can be created in the *Scripts* tab using the *Blacklists* subtab. Here the user can create new blacklists or manipulate existing ones. The words have to be entered comma separated.
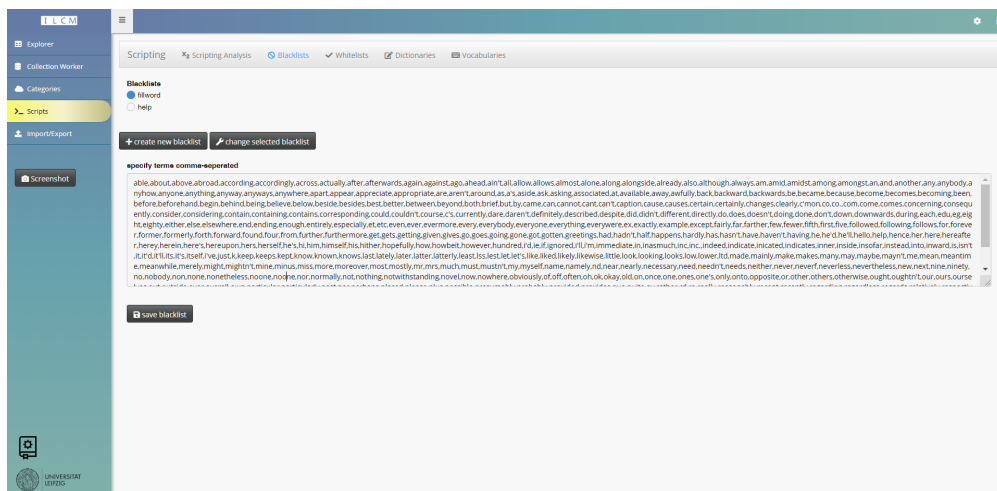


Figure 31: Blacklist Creation

## 13.3 Create whitelists

The whitelist can be created in the *Scripts* tab using the *Whitelists* subtab. Here the user can create new whitelists or manipulate existing ones. The words have to be entered comma separated.

29

## 13.4    Create dictionaries

The dictionaries can be created or changed and even deleted through the *Scripts* tab under the *Dictionaries* subtab directory. A dictionary consists of different columns. In the following example, a dictionary describing good or bad connoted words is used. To create a new column simply right-click on the table, a menu will open with different kinds of options to manipulate the current table. The concrete number of columns and rows is variable. After creating or changing dictionaries, the user needs to press the *Save dictionary* button. This will lead to a small window where the user needs to select a name for the dictionary and the columns' names.
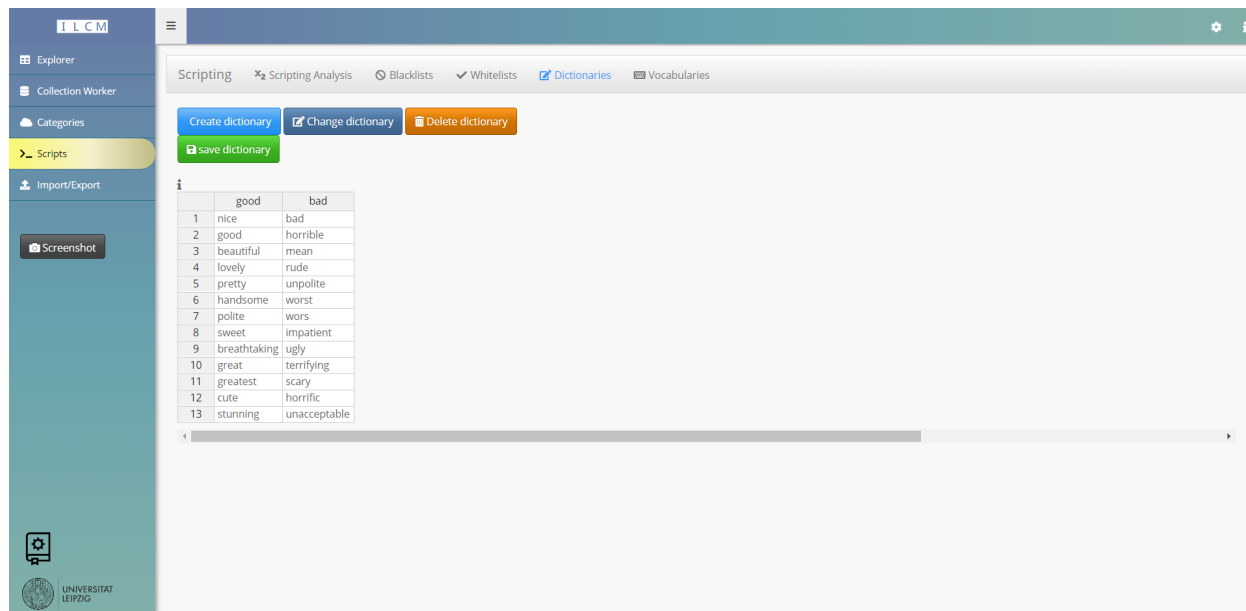


Figure 32: Dictionary Creation

## 13.5    Create vocabularies

Vocabularies can be created similar to blacklists or whitelists. Therefore the user can enter words for the vocabulary (comma-separated) or add through importing from finished analysis models. When selecting the last option, a window will open where the user can select the type of analysis and the id of the task (see *Collection Worker* under *Results*) and select the corresponding vocabulary list.

## 13.6    Special parameters

Besides the standard parameters for the preprocessing chain, each algorithm can be specified using parameters. The whole parameter setting for each started task is saved in the database. This secures the reproducibility of all calculated results.

1. **Parameters co-occurrence analysis**

- Doc/Sent co-occurrence: The calculation of co-occurrences can be done on different units of text. The iLCM offers the possibilities: document and sentence. This means a pair of words form a co-occurrence if they appear together in a sentence or a whole document.

- Minimum co-occurrence frequency: When working with vast amounts of data, the vocabulary can grow very fast. The co-occurrence matrix is a sparse term-term matrix. When the user works with a huge vocabulary, the intern object-size of the term-term matrix threshold for a minimum co-occurrence frequency, this means a combination of a term is assigned a non-0 value only if the words appear more often together than the set threshold. In most cases, this can be done without hesitation because the significance value for the word combination should be shallow anyway.
- POS-tags: For some co-occurrence analysis applications, it might be interesting to work on certain POS-tags words. This can be achieved with this parameter. Here the user can choose multiple POS-tags they want to keep. All other words that do not belong to one of the chosen POS-tags will be excluded from the analysis.
- NER-tags: Also, for NER-tags, the user can specify certain tags they want to use. The words that do not belong to the chosen tags will also be excluded. With this option, it can be pretty easy to create a network of persons. The user needs to select the NER-tag "Person", and the result will be a network of co-occurring persons.

2. **Parameter frequency extraction**

- POS-tags: Calculate the frequency counts just for words that are marked with one of the chosen POS-tags
- NER-tags: Calculate the frequency counts just for words that are marked with one of the chosen NER-tags

3. **Parameter dictionary extraction**

- Context filter list: In the *Dictionary extraction*, the dictionaries' frequency in the documents is counted. However, sometimes, the users are only interested in the counts if they appear in a particular context. In the context filter list, the user can enter some words that have to be found to be counted in order of the dictionary occurrence.
- Regular expression mode: Sometimes, the context is more complex than just a certain word requirement. Here the user can enable the regular expression mode, and as a result, the context filter list *Input* will be interpreted as a regular expression.
- Context unit: The text unit users want to control the presence of their context filter and, in the end, count the occurrence of the dictionaries. This can be done either on a sentence or document level.
- Dictionary: In the dictionary input, the user has to define his dictionaries. This has to be done in a certain format:
  Conceptname1:term1;term2;term3;...
  Conceptname2:term4;term5;term6;...
  These dictionaries can give a fundamental example for the starting point of a sentiment analysis:
  good:happy;satisfied;awesome...
  bad:sad;unsatisfied;terrible...

4. **Parameters volatility analysis**

- Minimum co-occurrence frequency: When working with vast amounts of data, the vocabulary can grow very fast. The co-occurrence matrix is a sparse term-term matrix. When the user works with a huge vocabulary, the term-term matrix's intern object-size can exceed the available hardware resources. This problem can be overcome by setting a threshold for a minimum co-occurrence frequency. This means a combination of terms is assigned a non-0

value only if the words appear more often together than the set threshold. In most cases, this can be done without hesitation because the significance value for the word combination should be meagre anyway. The co-occurrence matrices for several times slices form the basis of the context volatility calculation.

- Time interval: The aim of measuring context volatility is to calculate the context change amount in a specified time. With this parameter, we set the time unit. This can be days, month or even years.
- History: The *History* parameter determines the number of points in time compared to the specific point of time the context volatility is calculated. The number the user sets is corresponding to the set time interval.
- Number of Core: This parameter specifies the number of cores that are used for the calculation. Though at the fact, that the context volatility calculation is quite complex, it might be helpful to use more than one core.
- Method: The method parameter specifies how the difference of the expected value for the history of the significance of co-occurrence and the actual measured values is compared. This can be done by averaging the distance for every co-occurrence (standard) or calculating the cosine distance of the 2 vectors (expected values, actual measured values).
- Doc/Sent Co-Occurrence: The calculation of co-occurrence can be done on different units of text. The iLCM offers the possibilities: *document* and *sentence*. This means a pair of words form a co-occurrence if they appear together in a sentence/text.
- POS-tags: Calculate the context volatility using words marked with one of the chosen POS-tags.
- NER-tags: Calculate the context volatility just using words that are marked with one of the chosen NER-tags.

5. **Parameters topic modelling**

- Number of topics: Here, the users can specify the number of topics they want to be calculated. The number has to be at least 3 (2 does not work with the principal component for visualization with LDA-vis). The user can get an idea of the correct number of topics for a dataset by looking at the topic coherence page in the _ Result_ tab.
- Alpha: Alpha is the parameter of the Dirichlet prior to the per-document topic distributions.
- Method: There are 2 different topic model backends integrated so far. Here the user can switch between them. For more information have a closer look at the libraries: topicmodels::LDA and LDA::lda.collapsed.Gibbs.sampler
- POS-tags: Calculate the topic model using words marked with one of the chosen POS-tags.
- NER-tags: Calculate the topic model using words marked with one of the chosen NER-tags.

6. **Parameter classification**

- Context Unit: The classification of the data can be done on different units of text. The iLCM offers the possibilities: *document* and *sentence*. This means words get a specific classification depending on parameters corresponding to sentences or the whole text.
- Project: The user needs to decide to which project the current classification should belong. They can decide within a list similar to the annotation schemes.
- Mode: There are 4 different options of how the classification can be done: producing 50 new active learning examples, active learning on whole documents evaluating training set, classify the entire collection. Depending on the selected mode, there are different other parameters to specify the active learning method.

7. **Parameter sentiment analysis**

- Sentiment Dictionary: Essential for the sentiment analysis is the selected dictionary of predefined sentiments connected with the language of the data set. The user can currently decide between an English or German sentiment dictionary.
- Build document score by: The user can choose between different methods to calculate the document's sentiment (e.g. the mean of all found sentiments in the text).

8. **Parameter document deduplication**

- Similarity measurement: Here, the user can choose how the deduplication of documents will be calculated (Jaccard similarity or Jaccard bag similarity).

9. **Parameter dynamic topic modelling**

- Number of topics: Here, the users can specify the number of topics they want to be calculated. The number has to be at least 3 (2 does not work with the principal component for visualization with LDA-vis). The user can get an idea of the correct number of topics for a dataset by looking at the topic coherence page in the *Result* tab.
- Alpha: Alpha is the parameter of the Dirichlet prior to the per-document topic distributions.
- Chain variance: The user decides by which factor the dynamic topic modelling chains should variate.
- How to split dates: The user can decide how the different documents should be separated - they can choose between a separation by time or via automatic chunking.

10. **Parameter keyword extraction**

- Method: There are 4 different methods on how the keywords in the data can be extracted. Here the user can switch between them. For more information: …
- Minimal number of co-occurrence: The user can select a minimal number of occurrences per keyword (e.g. each word that appears at least 4 times in a document should be considered a keyword of this document).
- Maximal number of co-occurrence: The user can also select a maximal number of occurrences per keyword. Words that frequently appear in a text are more likely to be stopwords or no specific keywords for this document.
- Separator: Here, the user can decide what separation sign should separate the found keywords.

11. **Parameters syntactic parsing**

- Number of corse: Here, the user can decide how many cores should be used to calculate syntactic parsing. This depends on how fast the user wants results and how many cores are available in the current system.

12. **Parameters vector space representation**

- Minimum occurrence: The user needs to select a minimum number of occurrences for the words to become a vector.
- Number of vectors: Here, the user needs to specify how many vectors should be created during the task.
- Number of threads: The user needs to specify the number of threads used during the vector calculation.
- Windows: Here, the user can select a window size. This size will determine how many words will be considered to calculate the vector from a target word.

- Number of iteration: Here, the user can select several iterations to create the word vectors. This needs to be a trade-off between the accuracy of the word representation and overfitting the vectors for a specific collection of documents. Also, more iteration means the task will need longer to finish.
- Number of negative samples: The user can select several negative samples to prevent the vectors from overfitting. However, the number should be chosen in order to the size of the currently used collection.

## 13.7   Change script

Having specified all parameters, the user can start the script. R-experienced users may be unsatisfied with how the algorithms are implemented. Therefore they can also change existing scripts by clicking the *Show Script* button. This button will lead the user to the *Scripts* tab, and the users can see and edit the existing script for the chosen algorithm. The essential thing while editing a script is to keep the same output format as the original script to assure a correct visualization. The manipulated script can be saved and used for later calculation by ticking the *use custom script* field in the *Task Scheduler* subtab.



Figure 33: Manipulate Script for Co-occurrence Analysis

## 13.8   Check logs

In the *My Tasks* subtab, the users can also look at the amount of running, finished and failed tasks. They can also check the logs for each started task. Therefore they have to click on the box (running, finished or failed) and then select the task in the shown table below by clicking in the corresponding row. The log-file will then be loaded and shown:

Figure 34: My Tasks subtab

# 14 Results

The *Results* subtab is part of the *Collection Worker* tab. Here the user can find 2 columns. On the left, the user can see a box for each analysis algorithm. When expanding the boxes, the corresponding finished tasks will be displayed with their parameter setting (if possible). Clicking on a specific task will redirect the user to the *Details* subtab showing the tasks' results.



Figure 35: Results subtab with collection '21century' marked

## 14.1   Collections

In the right column, there is a box for the collections. Here the users have 3 options. They can have a closer look at the time series plot for every collection (click on *See*) or delete collections. By selecting a row in the *Collections* table, the current collection's intern variable is set. Like that, only the current collection results will be shown in the *Result* column. Furthermore, when switching to the subtab *Document* in the *Collection Worker*, the chosen collection documents are shown and can be read and annotated. When hovering over the name of a collection, its creation-statement will be displayed.

# 15   Details

In the *Details* subtab, the results of the chosen task are visualized. The *Details* tab depends on the chosen task and its underlying NLP-algorithm.

## 15.1   Details Topic models

- **LDA-Vis**: In the *LDA-Vis* panel the user sees the LDA-vis visualization for the calculated topic model. In the parameters-box they can change the number of words shown per selected topic. Besides there is the possibility to download the matrices for phi and theta in a csv-format.



Figure 36: LDA-Vis example for topic model calculated on guardian data for section 'sport'

- **Estimated Word Frequencies**: In the parameter box of *Estimate word frequencies*, the user needs to select at least one word to calculate its frequency. They can select between all the words used in the topic modelling documents. With clicking in the *words* text field, a list of possible words should show up. The user can choose between them. There are 3 subtabs

implemented. The "normal" *Frequency* of appearance for a word in the calculated topic will show a plot of all the topics and the number of occurrence of the word/words within the topics. Below this plot, the user will find a data table with the topic and the different frequency values of the selected words. In the *Document Occurrence* subtab, specific documents are mentioned that state the selected word/words and a counter for the occurrence for these words inside the document. In the *Document* tab, the user can look closely at the documents mentioned in the last subtab and take a closer look at the text. The selected words will be highlighted. The user also can download the estimated frequencies in the parameter box on the left side.



Figure 37: Frequency estimation over the topics



Figure 38: Occurence of the words within the documents

Figure 39: A closer look at the documents that state the selected words

- **Topic Proportion Over Time**: Here, the user can see the topic proportion of the current topic modelling over time. Therefore the user can select a lambda value for the topic labelling, the number of words per topic label and how to split the dates (by date or via chunking). They also can select a specific time interval (year, month, week, days) and the number of specified intervals perused period. This will result in a visualization that shows the proportions of the topic within a specific period.



Figure 40: Topic proportions over weekly intervals

- **Date Distribution**: The parameter box of *Document Distribution* panel shows a word

cloud for each topic. Inside the word clouds plots, the most likely words for every topic are displayed. At the top of the parameters section, the user can change the time series for the time series plot. Besides, they can switch between document count (how many documents are found), document probability (sum of all probabilities for the document belonging to the point in time for the selected topic) and the relative document count (percentage of the number of documents belonging to a specific topic, compared to all documents at the time). The belonging of documents is defined in the input below. The default setting uses the first rank mechanism (a document belongs to the topic for which it has the highest probability). If the user disables this option, they can set a probability threshold. A document then belongs to all topics; it has a higher probability than the threshold. In the top left corner of each word cloud box, the user can add the related topic (green plus-button) to the time series plot. Furthermore, they can see the overall probability of the topic and get an idea of the blue sparkline's times series plot. The time series plot will be calculated concerning the set options in the parameter box. Once the user has added at least one topic to the time series data, a data table is created with the topic number and the most likely words per topic. The user can select one or multiple topics to create a new sub-collection that contains those documents belonging to the chosen topics. The membership of the documents is calculated based on the settings at the top of the parameter box. To achieve this, the user needs to specify a new collection name. With the new sub-collection, the user can start a new NLP-task or annotate documents like the collection that forms a document search.



Figure 41: Time series graph for topic model

- **Coherence**: The coherence test enables the evaluation of the calculated topic model. There are 3 different measures implemented: Topic Coherence, Topic Intrusion and Word Intrusion. The topic coherence is in contrast to the 2 intrusion measures calculated without the users' interaction. For the topic intrusion, the user has to read a document and then tell which of the topics does not describe the document. In the word intrusion test, the user sees a group of words, whereby all except one are very likely on the same topic. The task then is identifying the word that does not fit together with the other words. The 2 parameters

*number of runs* and *set size* determine the number of examples the user has to work through before the results are saved and the number of topics/words the user has to find the intruder. The chance correction will subtract the estimated random hits from the performance measure.
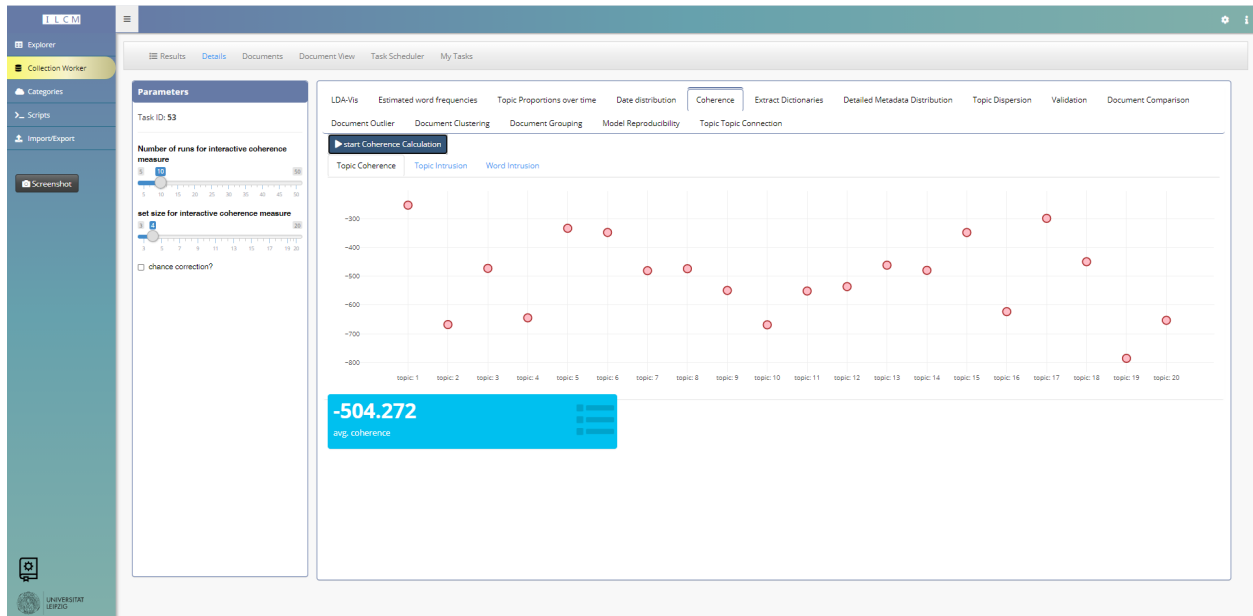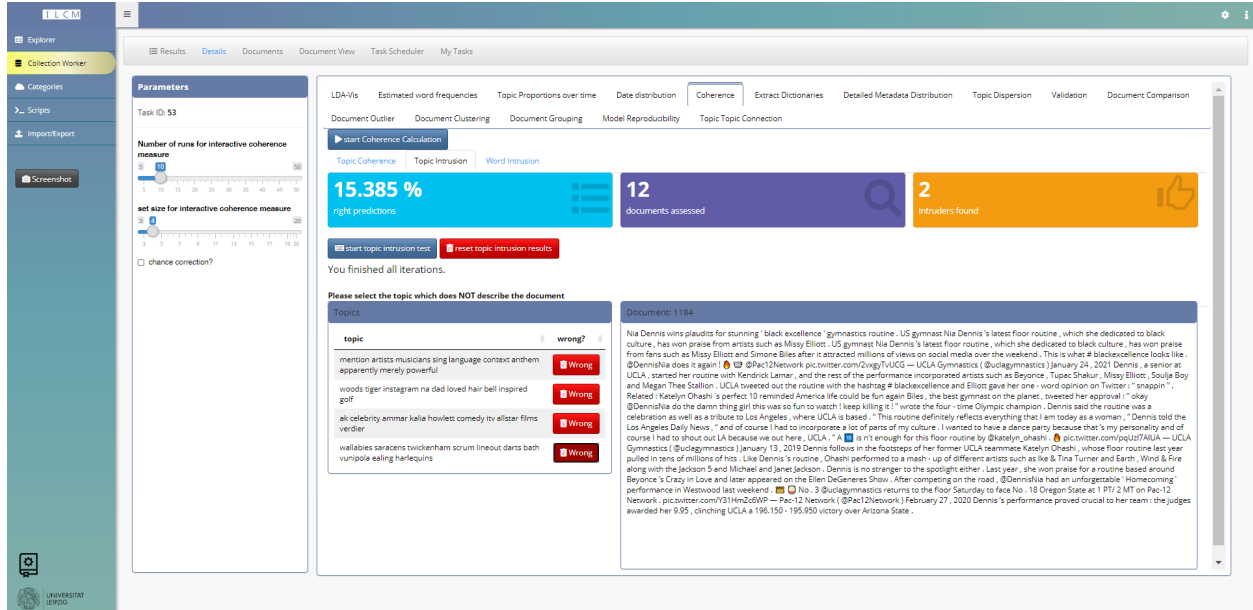


Figure 42: Example plot for a topic coherence result
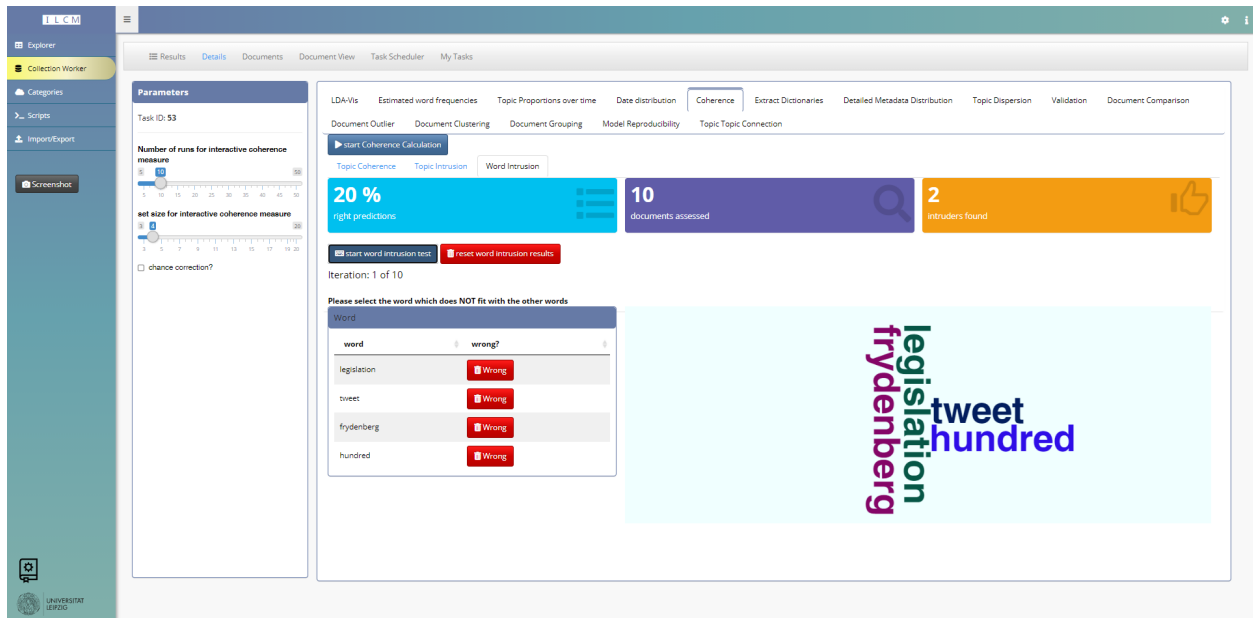


Figure 43: Example for a topic intrusion test

Figure 44: Example for a word intrusion test

- **Extract Dictionaries**: Here, the user can select different topics to extract the corresponding words into a dictionary. Each selected Topic becomes a category the user needs to make up a name for. They also need to name the dictionary in order to be able to save the dictionary.



Figure 45: An example of a potential dictionary consisting of words from three different topics

- **Detailed Metadata Distribution**: This subtab allows the user to combine the documents' meta-information with the topic model information. The user has 4 different visualization options. In the first one, the user gets an overview of selected metadata information corresponding to a selected topic. In the parameter box on the left, they can choose between the

41

shown meta-information (depending on the columns of the csv-file the collection was built on) and the number of the topic they wish to take a look at as well as a minimal occurrence of a certain value to appear in the plot. They also get the chance to chose a second metadata category. These parameters will then calculate the word cloud of the topic and a pie chart of the chosen meta category, and a data table below that the user finds the proportions of the first meta category to the second one. The next panel shows a *Membership Heatmap* of the first selected meta category above all topics. Next, the user will find a data table listing all the selected topic documents along with their metadata information. The last section is the *Topic Meta correlation*, where the user can see the relationship between all of the topics and the documents' found metadata categories.



Figure 46: Overview of metadata distributions for a selected Topic and metadata category

Figure 47: Heatmap of proportion of a selected metadata category over the topics



Figure 48: Insides in the documents of one topic and their metadata

Figure 49: Proportion of different metadata categories over different topica

- **Topic Dispersion**: This section shows the importance of different topics for all existing documents. First, the user sees a summary of all the topics and their importance for the documents. The columns show the overall importance and specific thresholds for the importance of a topic and the exact number of documents that evaluate the importance of the topic above this threshold. In the parameter box on the left, the user finds the minimal probability a topic must have to be part of a document's topic set (set when starting the topic model task). Above that, the user finds a text-box to select the topics they want a detailed analysis for. These will show in the *Details* window. Here the user finds different plots that show the topic probability distribution for the selected topics.

Figure 50: Overall view of the topic dispersion



Figure 51: Detailes on topic dispersion for selected topics

- **Validation**: In this section, the user gets an overview of essential documents for a specific topic. Therefore, they can choose whether the documents should be selected independently from the currently selected topic or by topic likelihood, leading to documents that represent the chosen topic. When selecting the independent variant, the user can switch between all possible documents. In the topic likelihood part, the user first needs to choose which topic the most relevant topics should be calculated and how many of them should be calculated. In both cases, the system will plot for each document distribution of topics, a word cloud for the most relevant words of a selected topic (see parameter box), and it will show the text of

45

the document. In this text representation, the most important and least important words for a selected topic are highlighted. The user can change the colour of the highlights as well as the method for calculating these words. These features will show when the user clicks on the gear icon in the left parameter box. In the parameter box, the user will also find metadata information about the current document. This allows the user to validate whether the current topic model has calculated useful topics or not.



Figure 52: Overview of a selected document in the validation section

- **Document Comparison**: Here, the user can select specific documents for comparison. They can choose between 3 different representations: A table view or a pie chart of the selected documents and the proportion of the topics within these documents, and the correlation of the documents. In the last part, the user can select a correlation method and a colour scheme. The euclidean and cosine similarity for the documents will also be shown in this section.

46

Figure 53: Table representation of the proportion of topics whithin selected documents



Figure 54: pie chart of the proportion of topics whithin selected documents

Figure 55: Correlation between the selected documents

- **Document Outlier**: This part allows the user to download and view comparison measures between the documents. The user can select between 3 comparison measurements: Correlation, euclidean distance and cosine similarity in the left parameter box. After selecting a measurement, the plot will show a data table consisting of the documents and their metadata and the average similarity for the documents. The user can then download the document similarity matrix for the current settings or the average similarities.



Figure 56: Table representation of document comparision

- **Document Clustering**: This section allows the user to cluster their data. In the parameter

48

box, they need to choose the number of clusters, the number of iterations, and the number of random seeds. After selecting these, a Graph will show up with different markers correlating to the number of clusters. This gives the user the possibility to see the spreading of the different clusters. The user can also change the size of the icons in the parameter box. If they want to get more information about a specific point in the graph, they can hover over the marker and get the corresponding document's title. For more information, they can also click on the marker to see the document's metadata or the whole text. In the other subtab, the user can see the clusters as data tables where each column represents one cluster, and the rows represent the documents. If the user is satisfied with the clustering, they also can download the current clustering result in the parameter box.
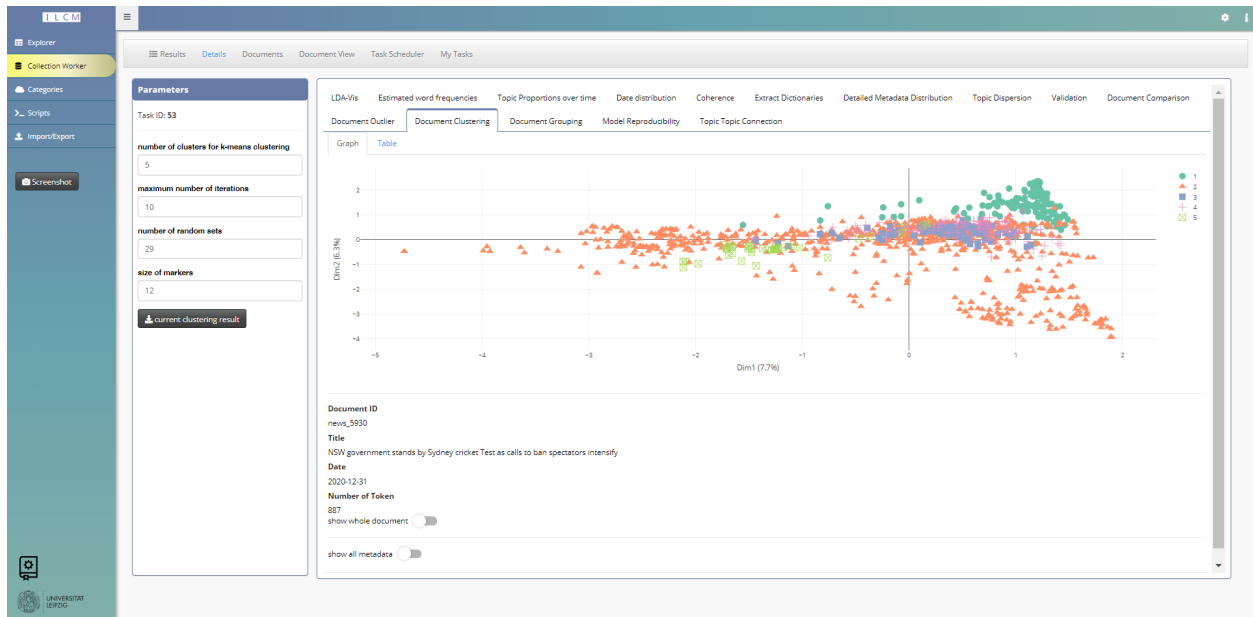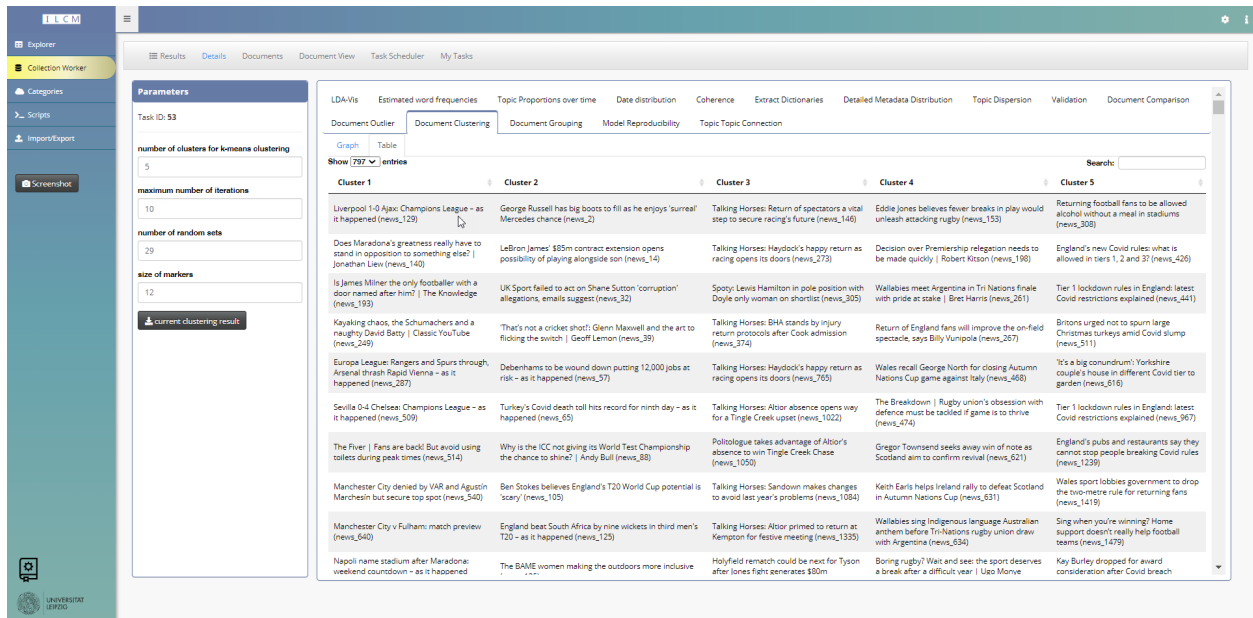


Figure 57: Graph visulization of 5 clusters

Figure 58: Table visulization of 5 clusters

- **Document Grouping**: In this section, the user can split the documents into groups depending on their metadata parameters. This comes in handy when the user wants to compare different groups of documents to the calculated topics. Let us say the user wishes to see the different impact of documents with different token sizes. Therefore the user will select in group one and two at least the token parameter. After this, a data table will be shown, and on top of it, the user will find a search bar. By clicking on the bar, the possible token lengths will appear, and the user can select a size for group one and in the next subtab for group two. In the *Comparison* subtab, the user will see the final comparison between the documents. They will find two-word clouds on top of the page summarizing the two groups' documents titles below that the user will see a data table or a heatmap of the impact the two groups will have on the calculated topics and overall correlation between the two groups.
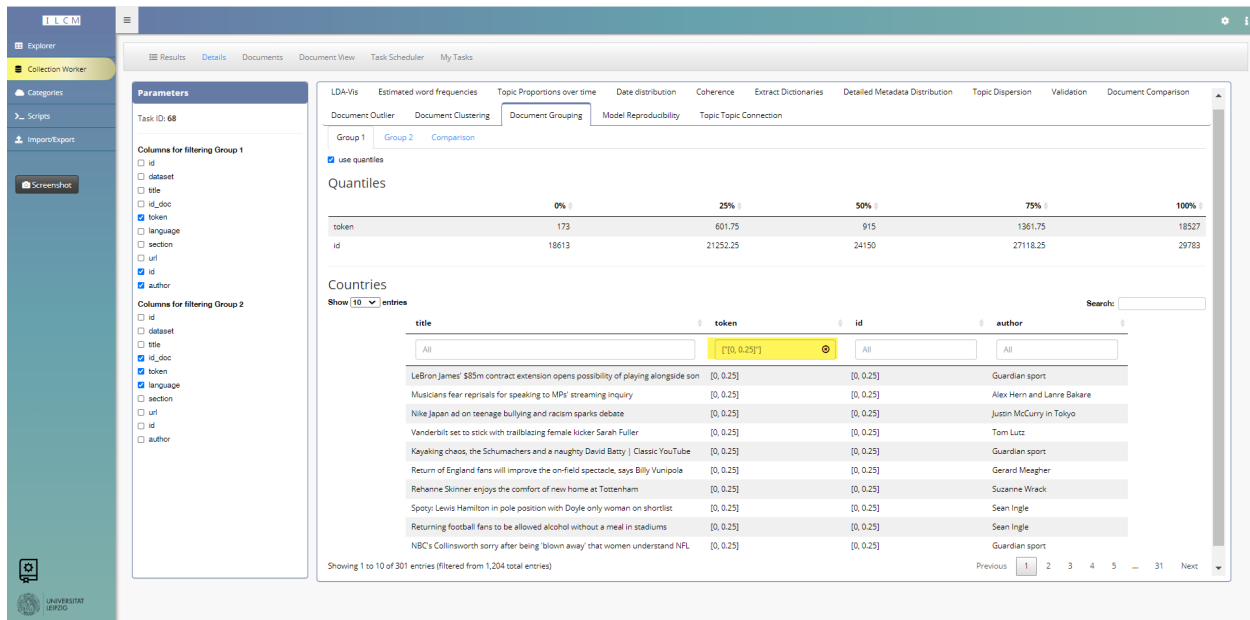
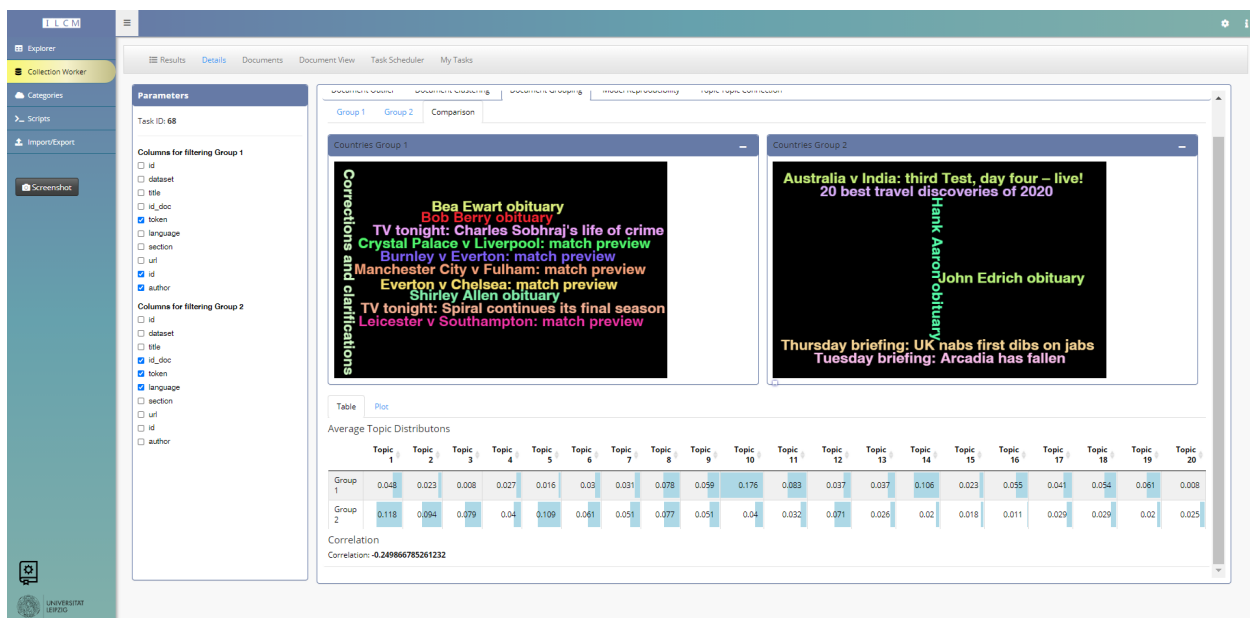Figure 59: Data table of group one after specifying the token parameter



Figure 60: Comparision of the two groups after specifying different tokens

- **Model Reproducibility**: This section will check the reproducibility of selected topic model analyses. The user needs to select the topic models he wants to compare and some other parameters, and then they need to click on *Calculate* to start the calculation process. Then a data table will show up that sums up the reproducibility of the different topics of the selected topic models. The user will also see the *Average Topic Overlap* above this table. This can help in term of the reliability of the found results.
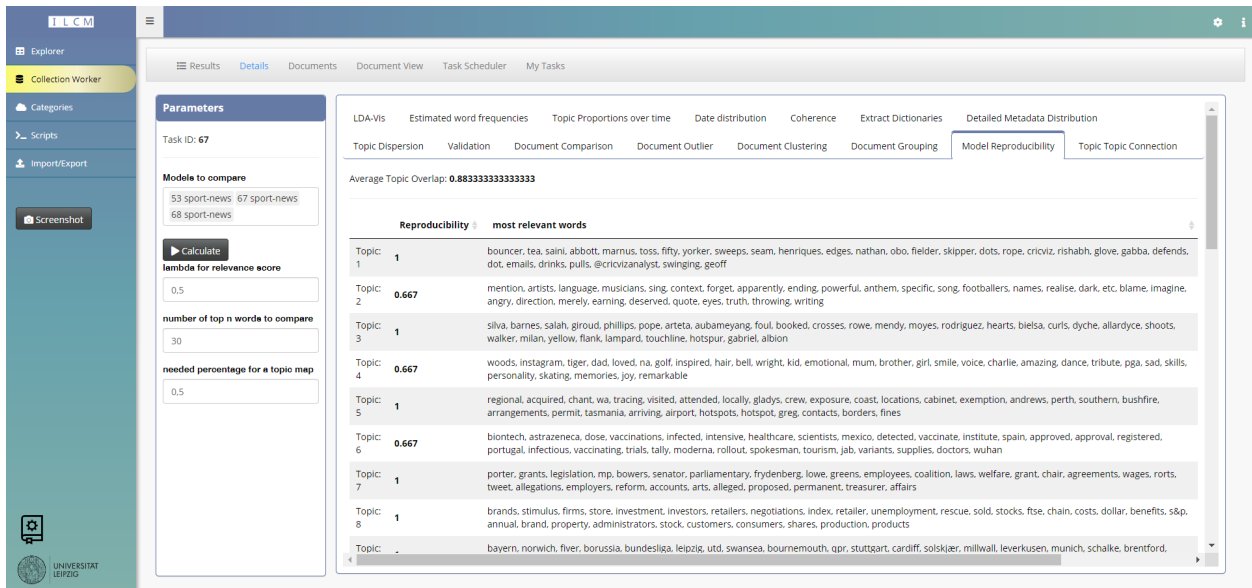
Figure 61: Reproducibility of 3 topic models over the same sport-news dataset

- **Topic-Topic Connection**: This section will show the topic-topic correlation between the current topic model analysis topics. The user can set a lambda value for the topic labels, the number of words per topic label (that get reviewed by the system to calculate the correlation) and a specific threshold. On the right, the user will see two plots: One of them is showing the topic correlation based on the topic's theta values, and the other is showing the correlation depending on whether the documents exceeds the set threshold.
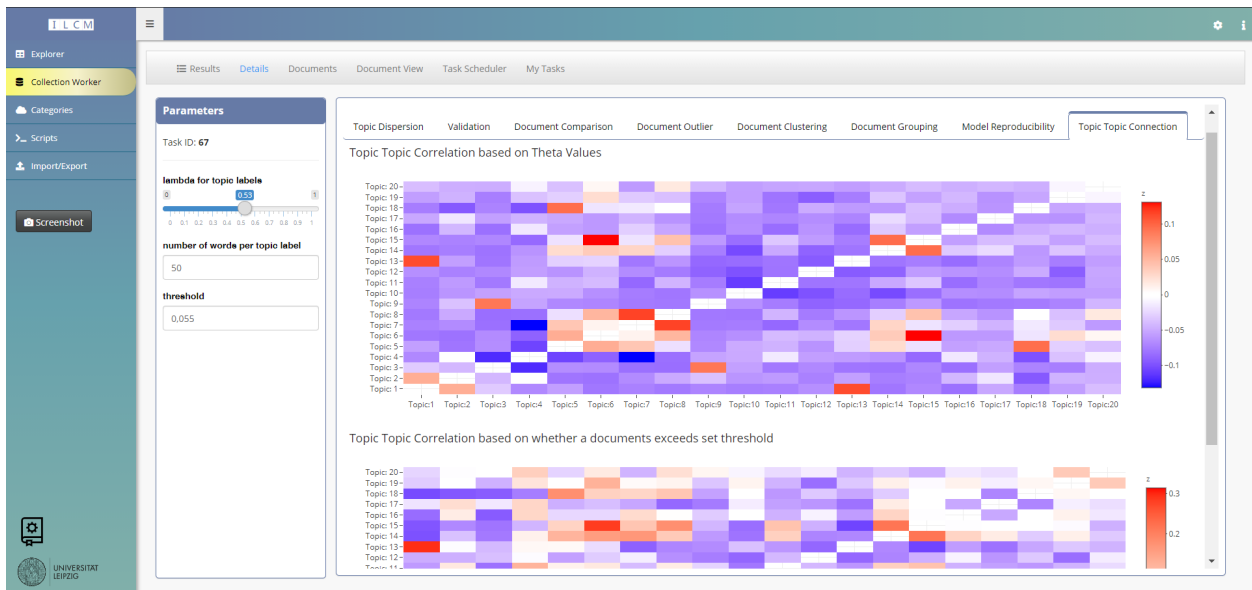


Figure 62: Heat map of the topic topic correlation

## 15.2 Details Dynamic Topic Models

- **LDA-Vis**: In the *LDA-Vis* panel, the user sees the LDA-vis visualization for the calculated dynamic topic model. The parameters box can change the time parameter to see the different topic dispersion within different time intervals.
- **Topic Dynamic over Time**: Here, the user will find a data table of all periods and the top words of a selected topic over these times. On the bottom, the user will find the option to select two time periods for comparison. After selecting the two-word clouds, one with the most and one with the least essential words for a topic within the two selected time stamps will appear. Besides this, a list of words is shown with a *change* parameter that will show the changed importance of this words over time. In the parameter box, the user can set the specific topic a number of words that should be shown and a lambda parameter for calculating which words should be shown.



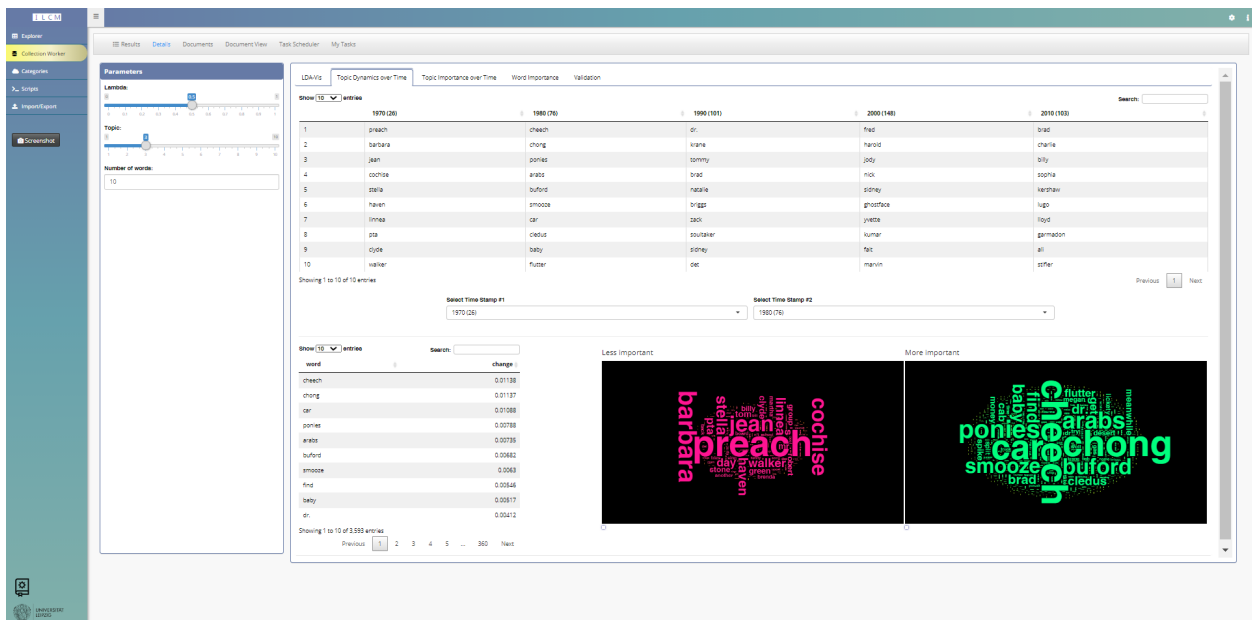Figure 63: Dynamic of topic 3 for movie data over time

- **Topic Importance over Time**: This subtab will show the correlation of the importance of the topics over the given time-stamps. At the top, the user will get an overview of the correlation of the topics within the timestamps, and on the bottom, there will be a visualization of how the importance of the topics has changed between the timestamps.
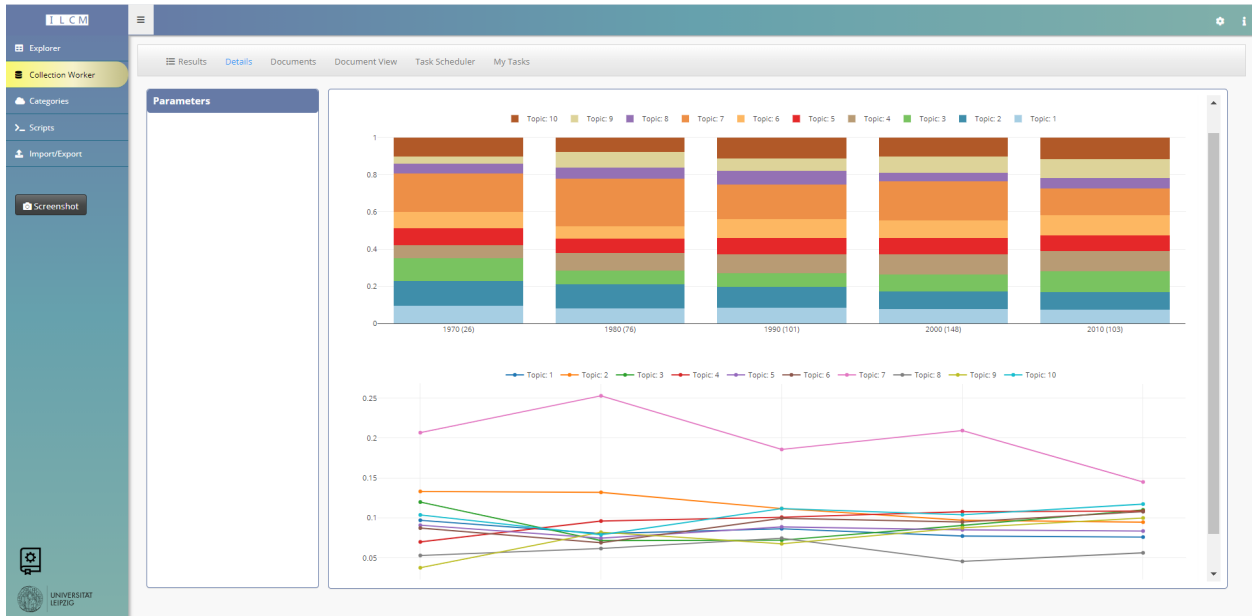
Figure 64: Importance of 10 different topics over the given time intervals

- **Word importance**: Here, the user can see the importance of selected words within a certain topic over time. Therefore the parameter box offers to select the words the user wishes to compare and a slider to select the topic the importance of the word should be calculated for.



Figure 65: Importance of the words 'bad', 'good' and 'playing' within topic 2 of a movie dataset over time

- **Validation**: This allows the user to review certain documents from the collection and check to spread topics in the document. Therefore the user needs to specify a timestamp, the number of documents in the selection, a document to view and a topic to present in the parameter box. The text of the selected document will be shown along with highlighted words that indicate

the selected topic. Additionally, the user will see a pie chart for the topic distribution for the chosen document and the most relevant words for the chosen topic.



Figure 66: Overview of the topic distribution for a chosen document in a movie dataset

## 15.3 Details Term Frequency extraction

- **Frequency plot**: In the *Frequency Plot* the time series counts for the chosen words are plotted. The time-series data can be downloaded by selecting the relevant words and then clicking on *Download Time Series*. The user can change the following settings: relative or absolute counts, time interval, counts on the document or word level and the words themselves. It is possible to plot and compare the time series data for multiple words

Figure 67: Frequency extraction results for 'sports competition' context

- **Most frequent words**: Here, the user will see the most frequent words of the collection. They have the option to enable the usage of the whole timespan. So they can choose a specific time interval and a point in time for the visualization. The system will then plot a data table with all the relevant words and their frequency counts and a word cloud with the most frequent words.



Figure 68: Most frequent words in a news related collection

## 15.4 Details on Dictionary Frequency Extraction

- **Frequency Plot**: In the *Frequency Plot*, the chosen dictionaries' time-series counts are plotted. The time-series data can be downloaded by selecting the dictionaries of interest and then clicking on *Download Time Series.* The user can change the following settings: relative or absolute counts, time interval, counts on document or word level. Its possible to plot and compare the time series data for multiple dictionaries
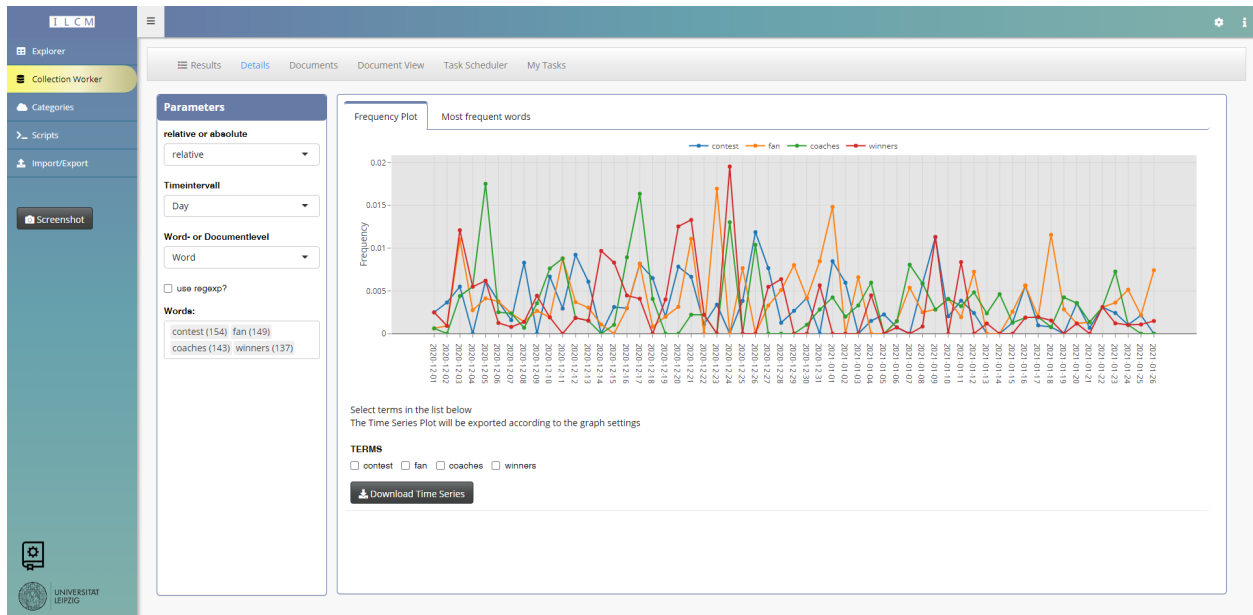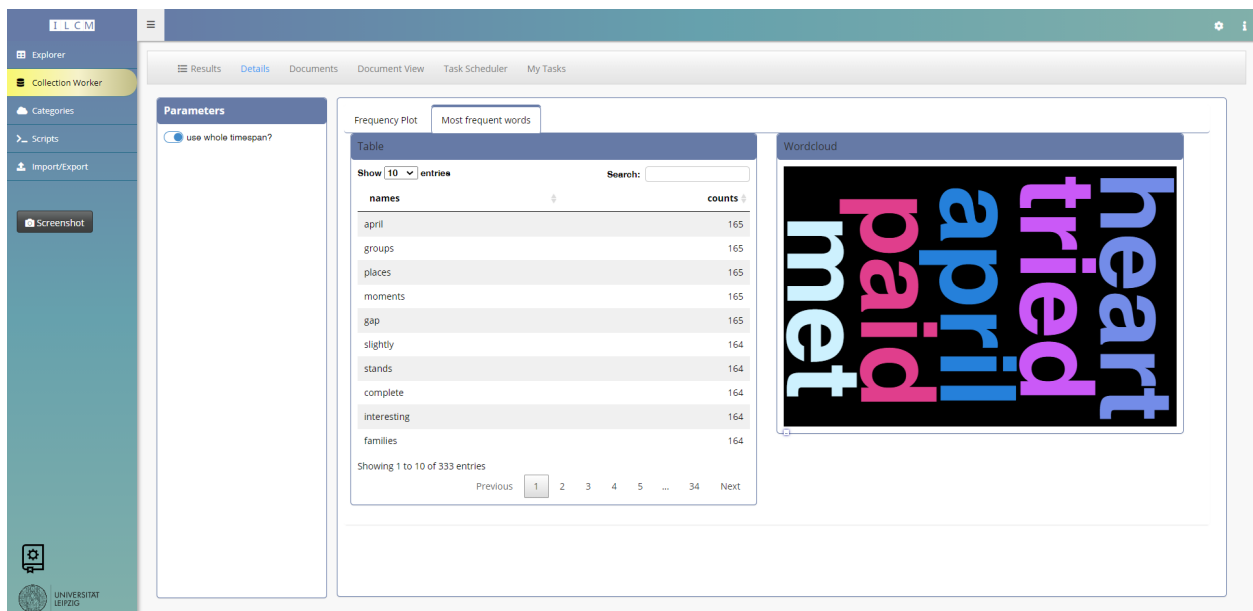


Figure 69: Dictionary frequency extraction example

## 15.5 Classification

When selecting a result from the *Classification* task, the user will be led to the *Categories* tab and the *Classifications* subtab. There they will find a list of all classifications to a selected project. They have the chance to evaluate the active learning example so and start the evaluation new in order to see a potential improvement. Additionally, they can check log and feature information from the classifications and delete a particular classification. On the other hand, the user can also see all categories on a single document example. Therefore the user needs to select training sets for the projects.

Figure 70: Classification tasks to find movie data corresponding to the genre 'Drama'
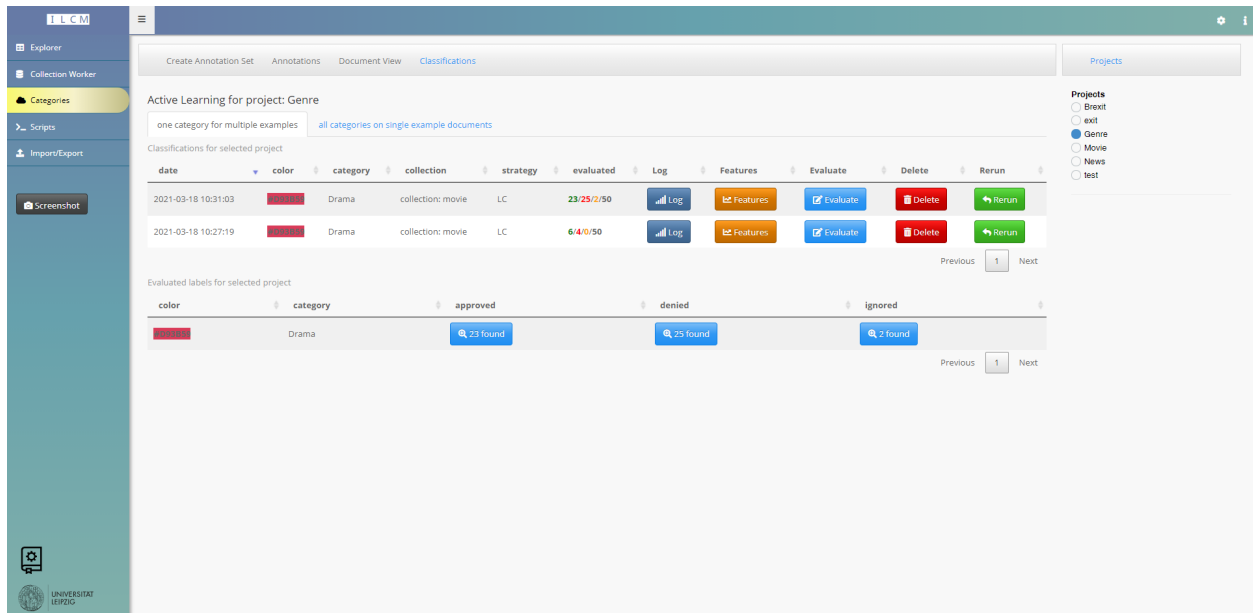
## 15.6 Details Co-occurrence

- **Co-occurrence graph**: The co-occurrence graph aims to visualize the chosen words' calculated co-occurrence behavior. For this purpose, the user has to choose a word (network root, blue-coloured dot) and a significance measurement for calculating the network (e.g. dice, mutual information, …). Then they can set the maximum number of links departing from the root. Next, they can set the maximum number of links for all other nods of the graph. With the parameter *Depth of Graph*, the number of iterations whereby new nodes will be added to the current leaves can be set. The higher this number, the more nodes the graph will have. For stability and clearness purposes, the user can also limit the number of links plotted to a certain number. With the *Charge of Graph* parameter, the user can manipulate the gravitation of the nodes (more details: visNetwork). It is possible to switch to a hierarchical layout by ticking the corresponding checkbox. Also, the user can either use straight edges or curved ones (smooth edges parameter). The node size can be determined as well by 3 measures: degree. Betweenness and centrality. The last parameter for the co-occurrence plot is the significance threshold. The values of the links reflect the calculated significance values. Only the links will be displayed with a higher value than the threshold with setting a certain threshold. When hovering over a node, the corresponding value to its size will be shown. After changing the plot parameters, the user has to hit the *Update Plot* button to see the graph's changes. When working with a depth of 4, recalculating the plot might take a few minutes. The user can also move the nodes around the canvas by clicking and moving them to the new position.
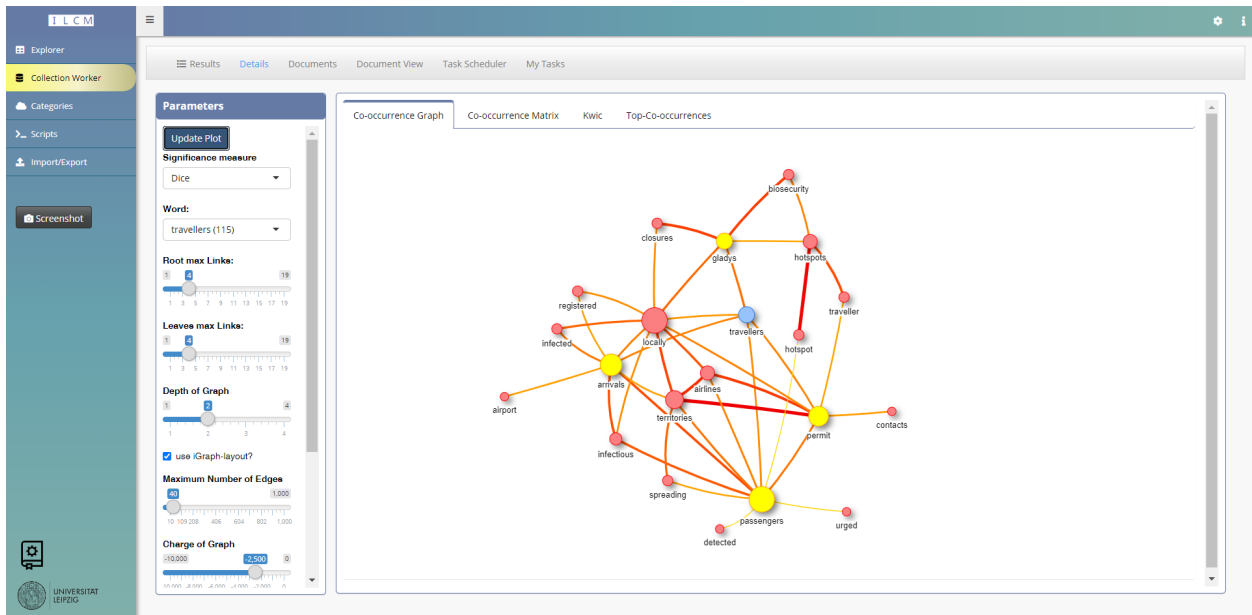
Figure 71: Co-occurrecne graph for the term 'travellers' whithin news reports from 2020 to 2021

- **Co-occurrence matrix**: In the *Co-occurrence Matrix*, the calculated significance (of a selected measurement) is displayed using a heatmap. The included words depend on the settings in the parameters section. The user can select an area in the plot to focus on a specific area.

- **Kwic**: In this subtab, the keyword and its context are displayed. The word can be chosen in the parameter setting. It is also possible to select more than one word. The number of shown examples and the size of the context needs to be defined.

- **Top-Co-Occurrences**: Here, the user will find an overview of the co-occurrence between at least two chosen words with different significance measurements. Therefore the user needs to select the first word for the calculation and an optional second one. If the second word textbox stays empty, the co-occurrence between all other words and the first selected one will be shown. In the parameter box, the user can set a maximum number of the shown occurrences, and they have the opportunity to download the results calculated with dice, mutual information or log-likelihood significance measurement.

## 15.7   Details Context Volatility

- **Single word**: In this subtab, the context volatility and the frequency for a single word are plotted. The user has to select a word in the parameters section and click *update plot.* In the time series plot, the orange line indicates the frequency of the selected word. The blue line shows the calculated context volatility values. Below the times series plot, there are 2 boxes. In the upper one, the most significant co-occurrences for the selected word and a specific point in time are displays in a table and a corresponding wordcloud. The user can specify the point in time by clicking on the time series plot to correspond to the times they are interested in. The two points in time are the last points the user has clicked on in the time series plot. Most divergent means those words, who have changed their significance value the most. Green words in the word cloud indicate words with a higher significance value on the

59

last clicked point in time than in the one before. Pink words show those words that reveal more minor significances.
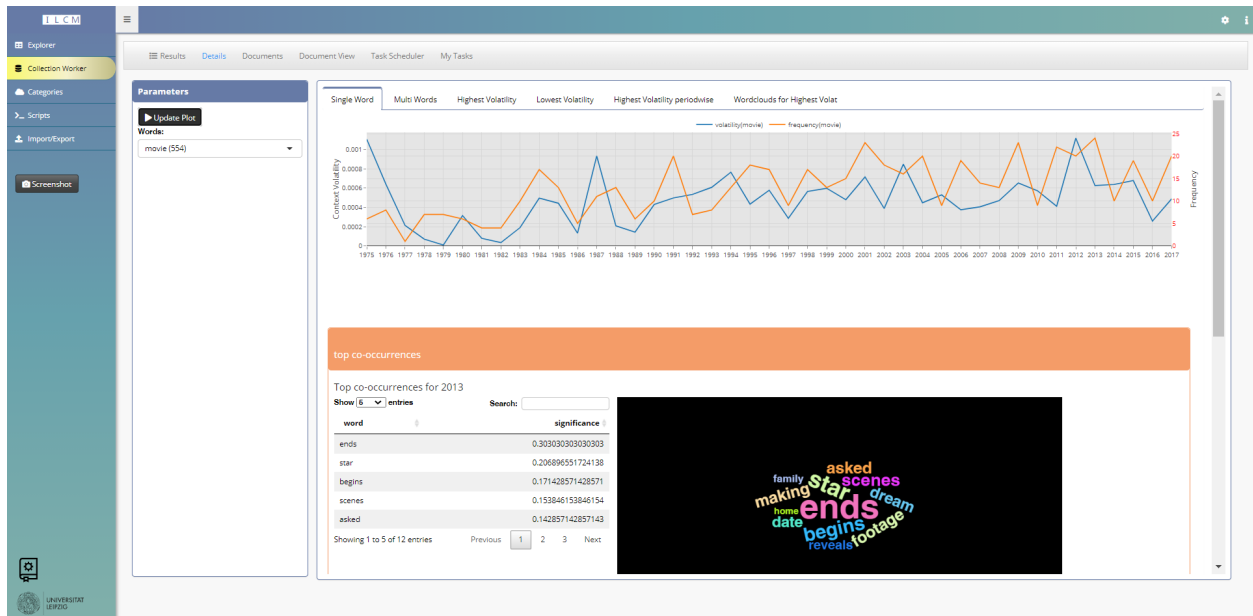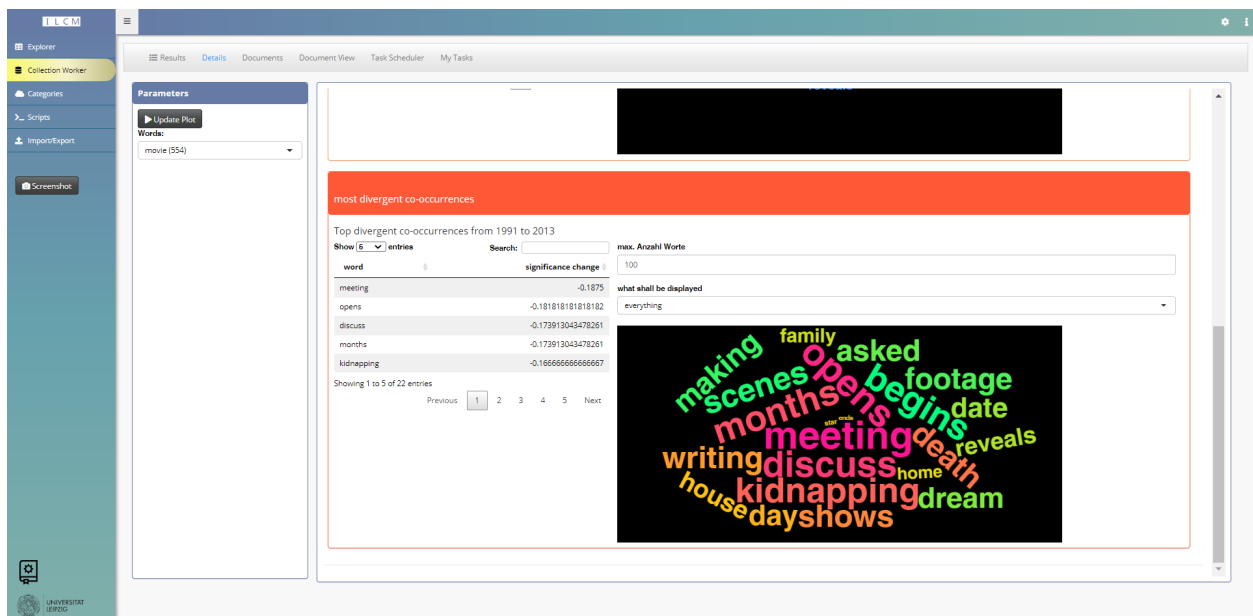


Figure 72: Context volatility for the word 'movie'



Figure 73: most divergent co-occurrences for the word 'movie'

- **Multi words**: In this subtab, multiple word's context volatility can be prepared. Therefore the user can select multiple words from the select input on the parameters settings.
- **Highest volatility**: In the *Highest Volatility* subtab, a table with the average word frequencies over time and the overall average context volatility is shown. In the parameters section, a maximum and a minimum frequency can be specified. The results are ordered to decrease

by their context volatility.

- **Lowest volatility**: In the *Lowest Volatility* subtab, a table with the average word frequencies over time and the overall average context volatility is shown. In the parameters section, a maximum and a minimum frequency can be specified. The results are ordered increasing by their context volatility.

- **Highest volatility period wise**: The subtab *Highest Volatility Periodwise* enables the user to identify the words with the highest context volatility values for a specific point in time. This point in time can be specified in the parameters section. Also, in the parameters section, the user can limit the shown words to a specific NER- or POS-tag. In the graph, the words are plotted concerning their frequency (x-axis) and their context volatility (y-axis). When hovering a point in the plot, the corresponding word is shown with the calculated frequency and context volatility.
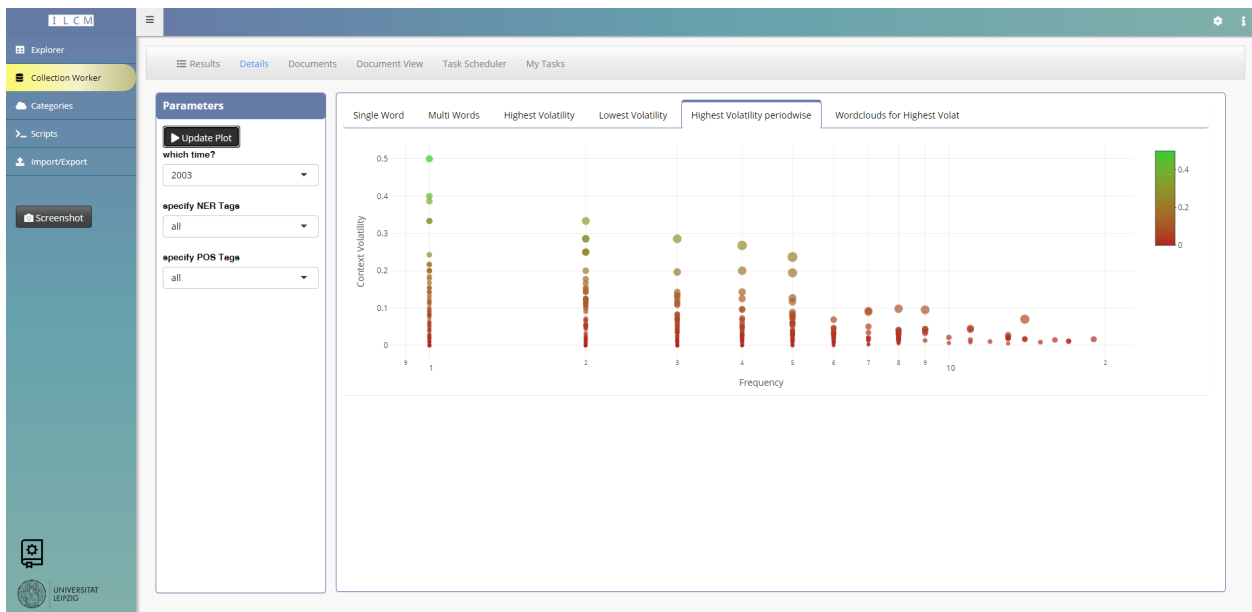


Figure 74: Most volatile words periodwise

- **Wordclouds for highest volatility**: In the subtab *Wordclouds* for Highest Volatility, the most volatile words are displayed at a specific point in time. Nevertheless, this time, the results are shown in the wordcloud for each point in time available. On the right to each wordcloud, there is a table with the exact values calculated. In the parameters section, the user can change the time intervals depending on whether the chosen time interval in the parameters of the *Task Scheduler* was set today, month or year. If "day" was the chosen time interval, the periods could be aggregated to month and years in the visualization. If the choice was "month", The periods can be month or years. When "year" was the chosen option, it was impossible to split the data back into data for months or days—furthermore, it is again possible to restrict the words to a specific NER- or POS-tag. In addition to that also a minimum and a maximum frequency can be defined. The last options declare the maximum number of words included in each wordcloud.
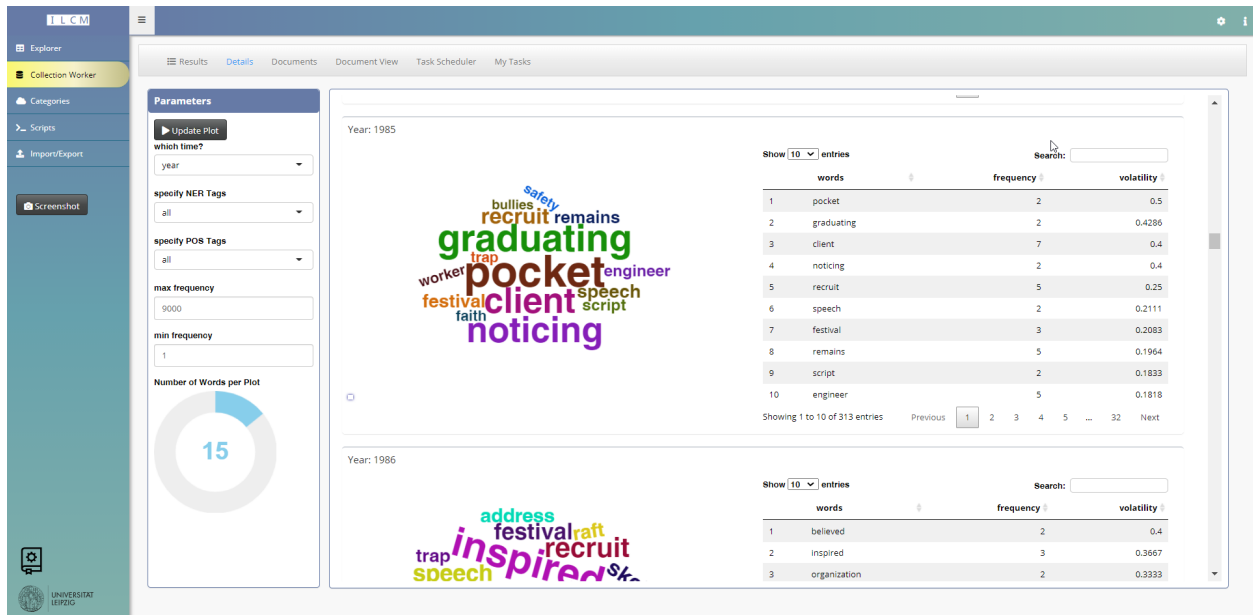
Figure 75: Most volatile words periodwise - wordclouds

## 15.8 Details Sentiment Analysis

- **Time Series**: In this part, the user will see the average sentiment scares along with the number of found documents over time. Therefore the user can select a specific time interval. The user can also decide if the plot should only show points or points connected with a line. On top of that, the user has the chance to download the current plotted data.
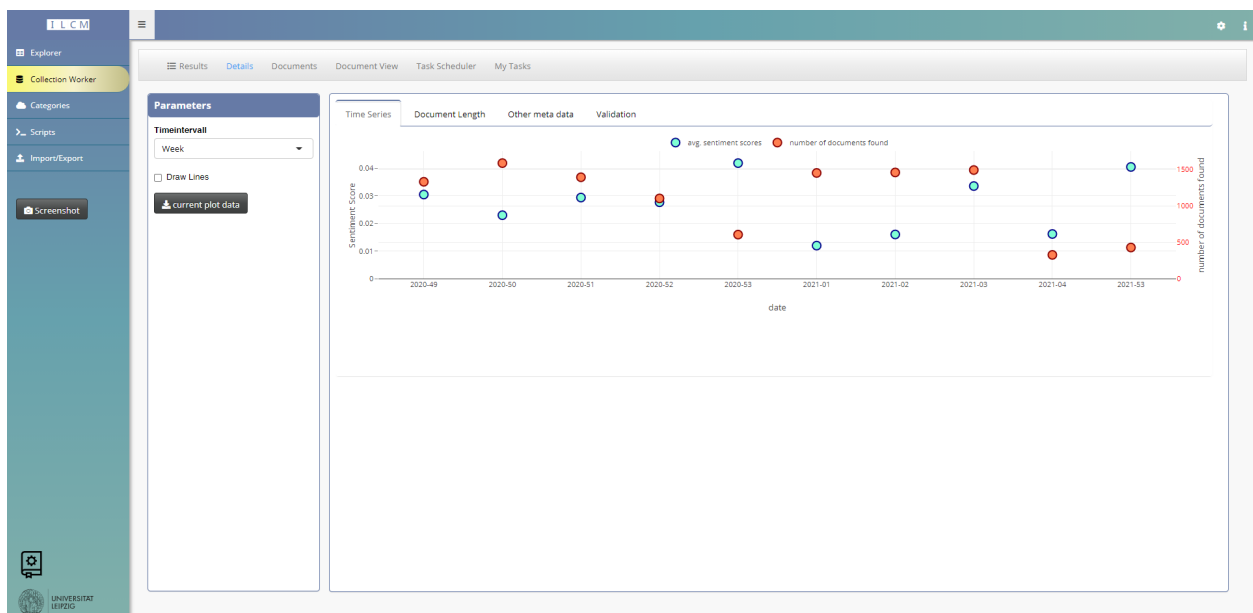


Figure 76: sentiment scores and number of found documents whithin a weekly time intervall

- **Document length**: In this subsection, the user will see a sentiment score in correlation with the number of found documents and the number of words per document. The user can

change the maximal number of breaks which will make the plot more accurate. They can, again, select to draw lines between the shown dots or download the current plot.

- **Other metadata**: Here, the user gets the chance to see the different sentiment score corresponding to the used collection's metadata. Corresponding to the metadata, several subtabs will be shown each for one specific piece of metadata. For each part, the user can download the plot and specify the minimal number of occurrence to be shown in the plot.
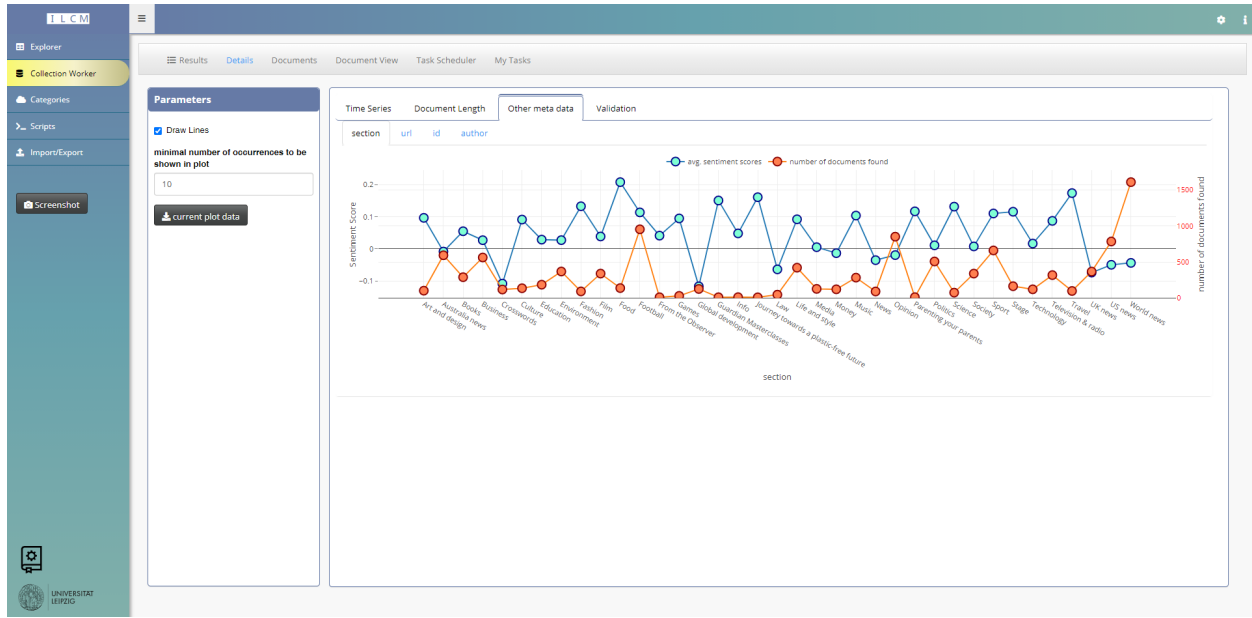


Figure 77: Sentiment scores corresponding to the categories of news articles

- **Validation**: In this subsection the user has the chance to validate the results of sentiment analysis. Therefore they can either show single documents with the found sentiments highlighted (red = negative, green = positive) or see a list of the top documents with the most positive and negative documents.
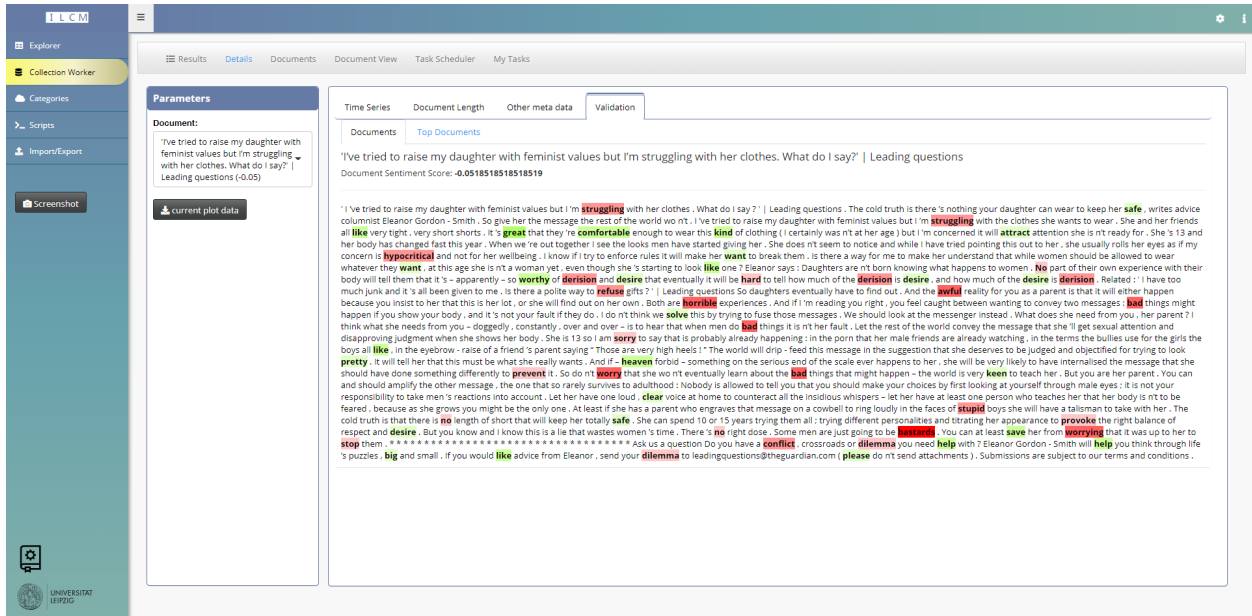
Figure 78: Sentiment analysis of a news aerticle

## 15.9  Details on Vector Space Representation

- **Similarity**: In this section, the user can compare the selected words' vector representation. Therefore the plot will show a coordinate system and presents the selected words as red dots and other not selected words within this space as blue dots. In the parameter section, the user can select the number of shown points within the coordinate system and a reduction method. On the left, the user will also see a list of the distance between the shown points.
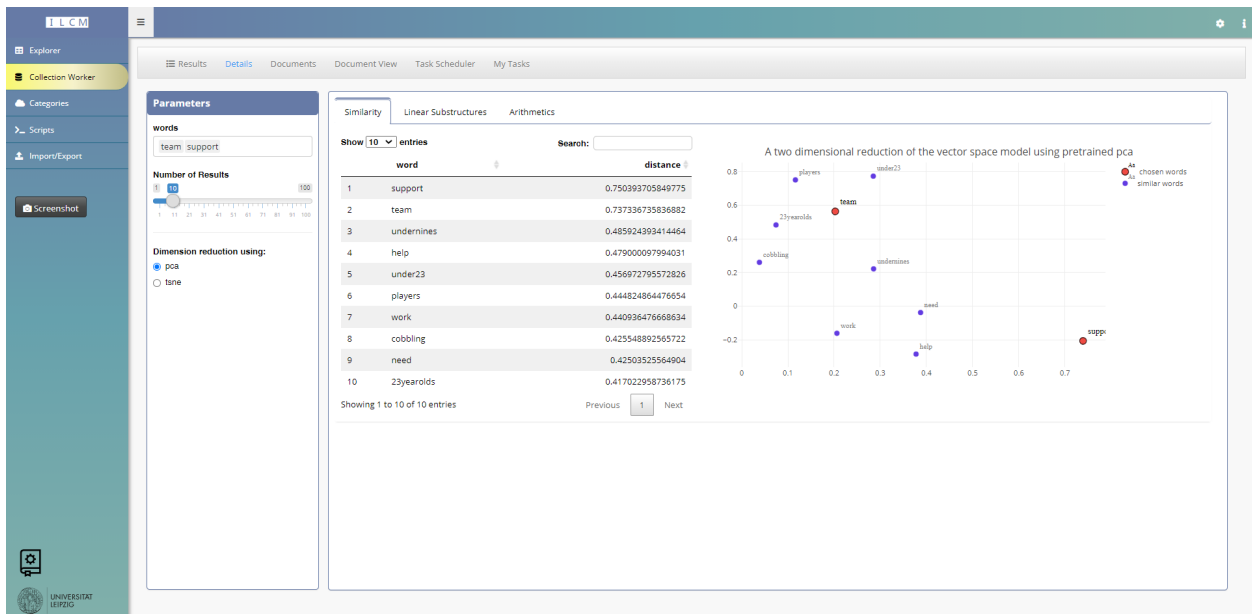


Figure 79: Vizulisation of the vetor space representation

- **Linear Substructures**: Here, the user can select a certain number of word pairs or triples

plotted in the coordinate system. This will give the user a vague idea of how close the selected words are related when calculating the vector space representation with PCA or tsne.

- **Arithmetics**: In this section, the user can select specific words (at least one) for which the cosine similarity to other words will be calculated. This calculation is shown in a data table. The user can also set the number of relevant words (the ones with the nearest similarity to the selected words) they wish to see along with a list of words for which the cosine similarity should not be calculated (e.g. for stopwords).
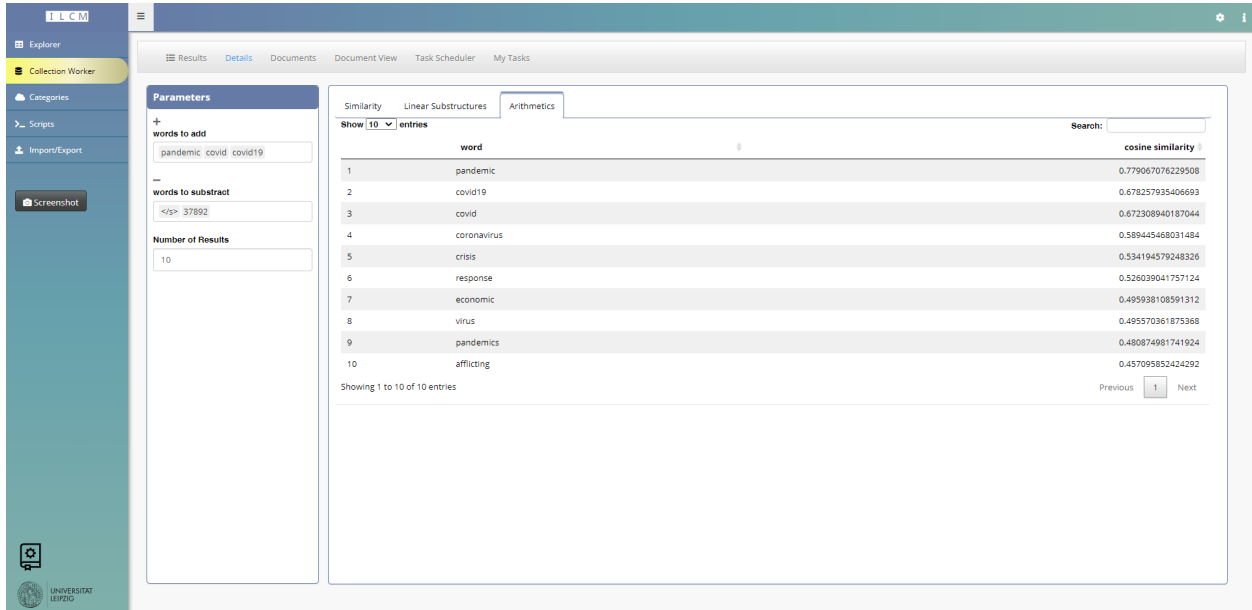


Figure 80: Cosine simularity of chosen words in a news copus

## 15.10   Details Document Deduplication

- **Graph**: This section will show the user a graph consisting of all documents as nodes. There will be an edge between two nodes if they are in some way identical. In the parameter box, the user can set the default strategy of placing the nodes in the graph and a threshold defining when two documents are similar. The user has the chance to download either a list of duplicates or duplicate free data. They can also reset their input or save the current collection of words. Within the graph, the user can zoom in in order to see the names of the nodes.
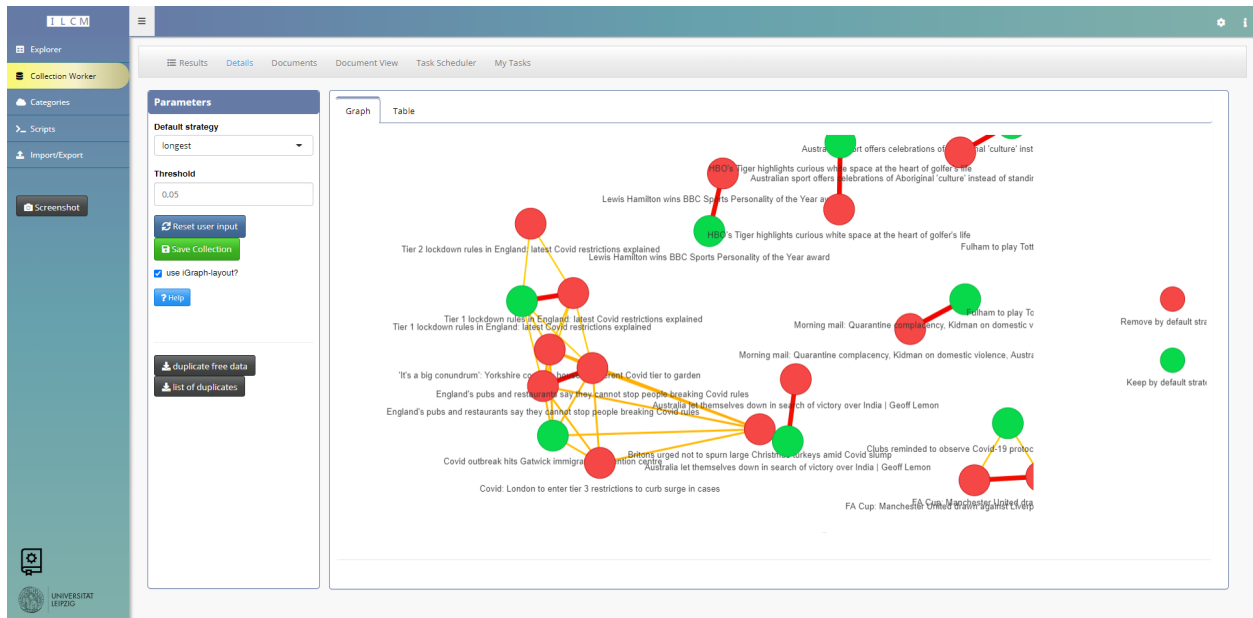
Figure 81: Part of a news document graph, zoomed in to an example of simular documents

- **Table**: Here, the user will find two data tables. The first one will show a list of documents the system would keep (to prevent duplicates), and in the second one, the user will find documents the system would delete from the collection. The user can now scroll through the lists and either remove or keep a document in his collection by clicking on the corresponding button in the document row. They can also check manually if the document is similar to another one. Therefore they need to click on the button in the *Diff* column of a document row and select via the given document-id-list a document they want to compare the selected one with.
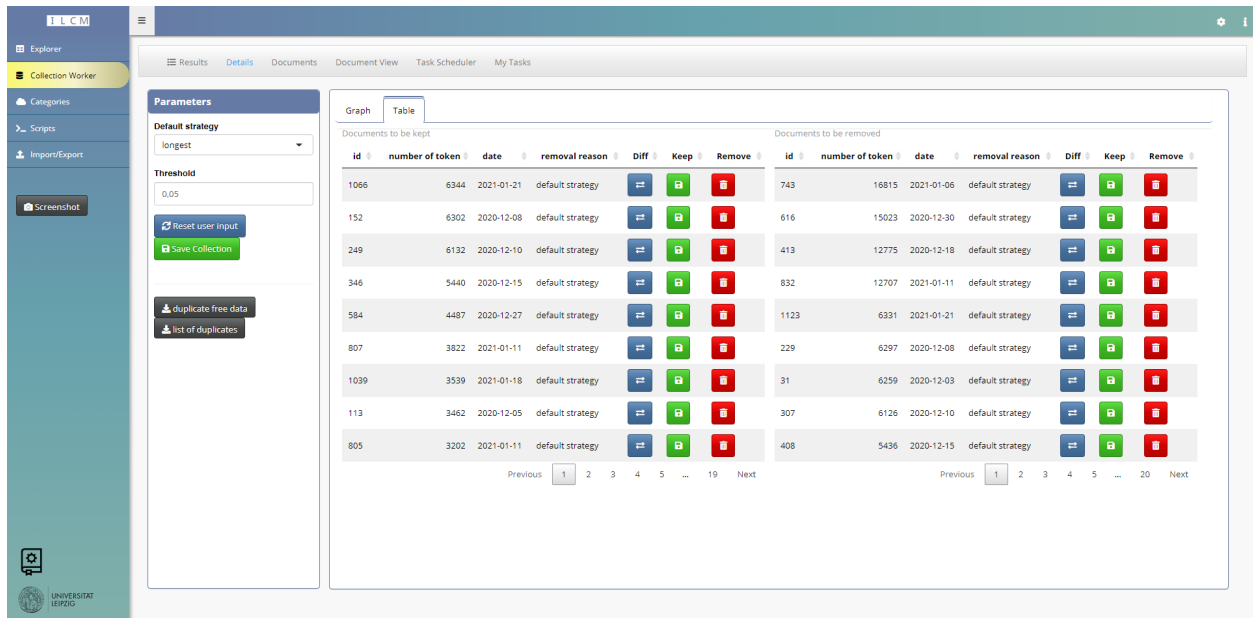
Figure 82: Data tables of documents to keep and delete from a collection according to theyre similarity to each other

## 15.11 Details Keyword Extraction

The parameter in this section will depend on the corresponding method but in some sort they will always allow the user to specify the frequency each word needs to fit in, in order to appear in the plot along with a parameter to specify the number of words/phrases the user wants to see.

- **Plot**: Here the user will see a plot of selected keywords or key-phrases and their frequency. The plot can be changed by setting the parameters for a frequency interval or a limit for the number of worda within key-phrases, etc.
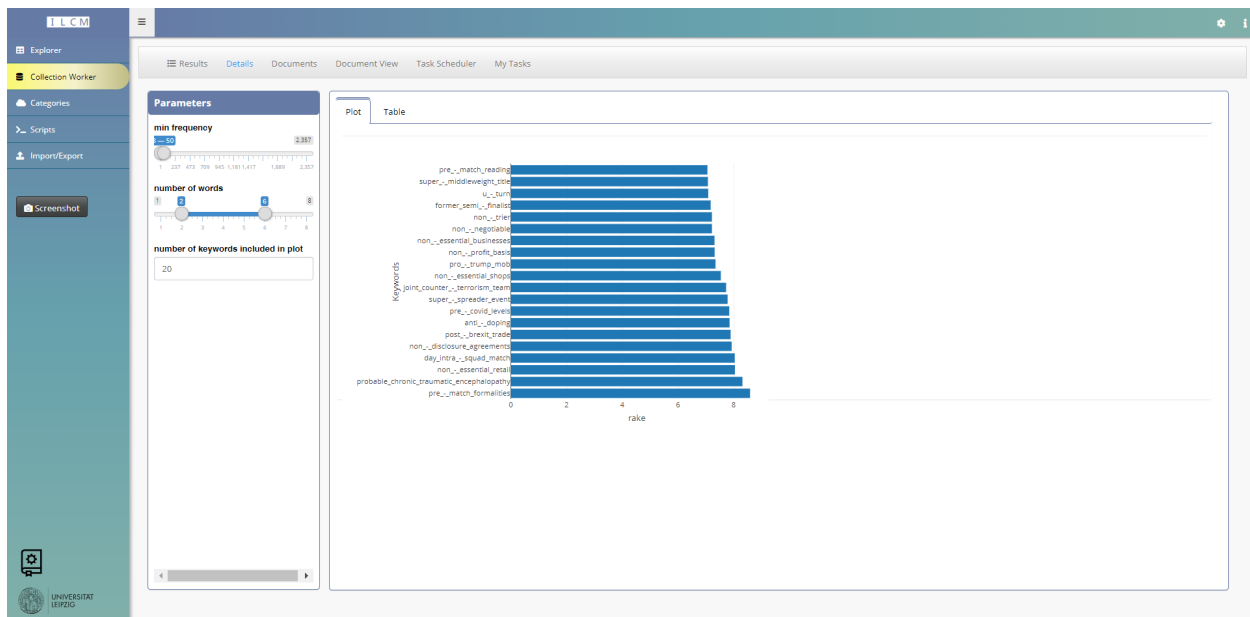
Figure 83: Plot for the keyword extraction with the rake method

- **Table**: In this subsection, the user will see the extracted keywords in a data table. With the parameter, they can again specify which specific kind of keywords they wish to see. They also can copy this list of keywords or download it as a csv-file or in excel-format.

## 15.12   Details Syntactic Parsing

- **Search**: Here, the user can get an overview from syntactic parsing by selecting a specific subject, predicates and objects he can specify the sentence results they wish to see. The changing of these parameters will result in a change of the list of documents. The user can select a specific sentence in the data table and see the calculated syntax tree. This sentence will be extracted from the documents in the collection.
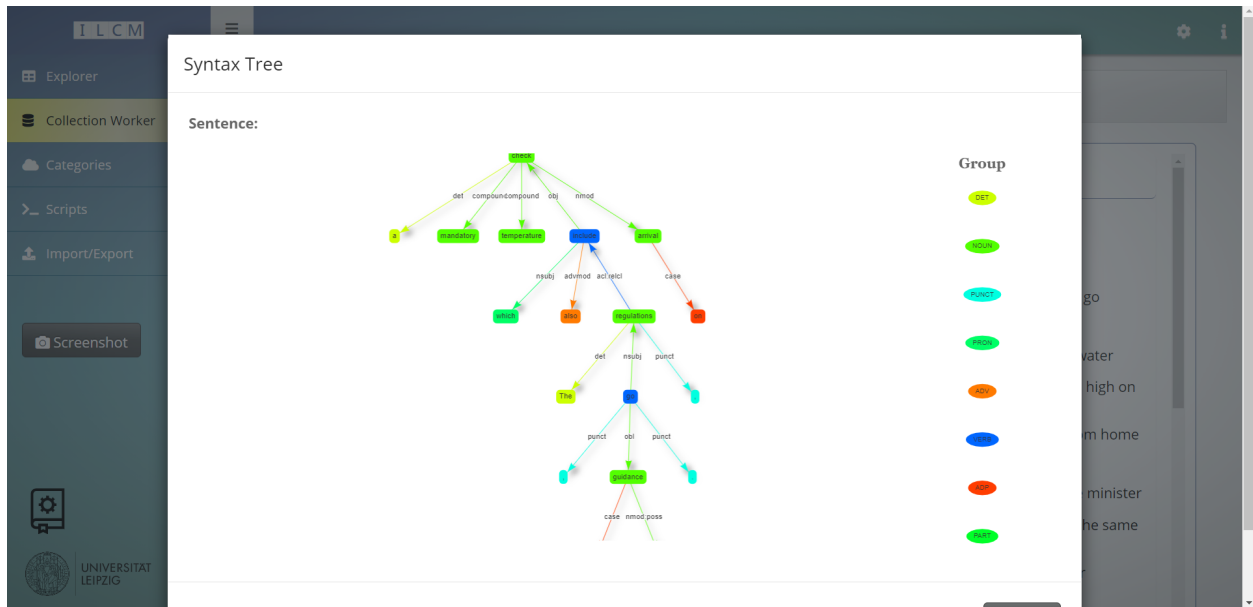
Figure 84: Syntax tree for a sentence containing the selected words 'regulation' and 'include'

# 16 Closing iLCM and restarting

The iLCM can be stopped using the starndard docker command for managing running containers (docker stop). All results, collections, annotations, blacklists, and essential data to the database and Solr will also be available when restarting the images. This is secured with the use of docker volumes. The stop and the start commands are using the name of the running container. The name can be seen using 'docker ps'.

# 17 Working with RStudio

While using the local version of iLCM the user can open the whole iLCM-tool in RStudio. This allows the experienced user to get a quick overlook of all the data used for performing the NLP-Tasks. After importing the iLCM project the user can see a couple of folders. In the folder *collections* the user will find all personalized data from the iLCM, such as *annotation-schemes*, *blacklists*, *whitlists*, and so on. Clicking through these folders the user will find *Rdata* or *CSV* files labeled like in the iLCM program. In the folder *logs* the user can additionally find the the log-entries of all started task separated into *failed*, *finished*, and *running*. This is similar to the lists the user will see under the *Collection Worker* tab in the *My Task* subtab. To change existing files in RStudio just open them by clicking on the name in the directory list then the selected file will open and the user can change the content of the file respecting the format of the file, don't forget to save the changes after editing the files. Keep in mind that *.Rdata* files will not open as a file but be loaded into the global environment of the RStudio distribution. To change these kind of files use the console and write the corresponding R-commands. After the changes are made the user can search for the *app.R* file in the iLCM directory and open the file and click on the play button that will appear on the top right in order to start the iLCM tool.
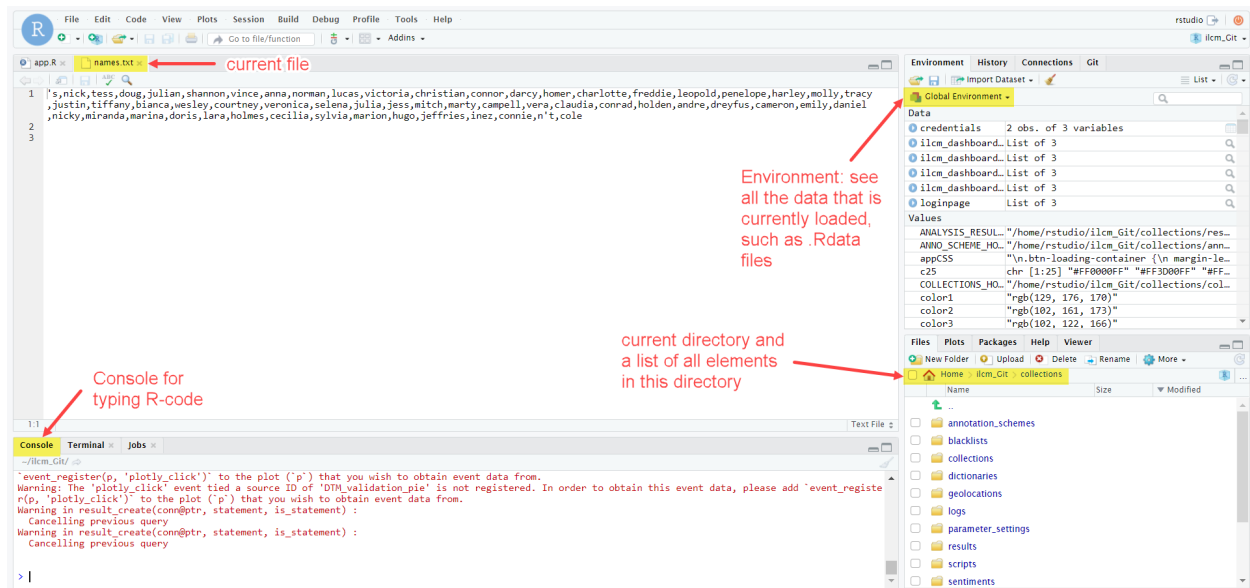
69

Figure 85: Overview of Rstudio

# 18 For Questions

If you are facing any questions, feel free to contact kahmann@informatik.uni-leipzig.de. For more information please visit the iLCM website.

# 19 Problens and Feature Requests

Any problems and feature requests can be reported via the Github-Page.