# ORACC IN KORP USER GUIDE (Korp v.9) – July 2024

Heidi Jauhiainen

Centre of Excellence of Ancient Near Eastern Empires, University of Helsinki

Korp is an online concordance tool that contains texts in various languages. It also contains two versions of the Open Richly Annotated Cuneiform Corpus (Oracc) data: Oracc in Korp 2021 is a snapshot from June 2021 and Oracc in Korp 2019 was downloaded in May 2019.

Direct link to **Oracc in Korp 2021**: http://urn.fi/urn:nbn:fi:lb-2022031706
Direct link to **Oracc in Korp 2019**: http://urn.fi/urn:nbn:fi:lb-2019060602

Instead of using the direct links, you can also go to https://www.kielipankki.fi/corpora/oracc/ and click "Open the resource in the concordance service Korp" of either version. The page also has a link to the downloadable version of Oracc in Korp 2019. The downloadable version is a zip file containing .vrt files for each of the projects in that corpus. The files have one word per line with their tab-separated annotations.

You can also go to https://kielipankki.fi/korp and change the language from Finnish to English in the upper right corner of the window. Then select **Other languages** (referring to the languages of the text corpora) in the upper left corner. Finally, click the "Select corpora" bar next to the Korp logo and click the arrow next to the words **Cuneiform/Nuolenpääkirjoitus**. See below how to select a dataset.

The next snapshot of Oracc data for Korp is planned to be taken at the end of 2023 and to be released in the spring 2024.

Korp has

This guide replaces the one from July 2019. Korp has since then been updated to version 9 and moved to new servers. The Oracc in Korp 2021 been added with some new attributes for both Oracc corpora.

For a general guide to the Language Bank of Finland's Korp tool see https://www.kielipankki.fi/support/korp/. Note that it is for the older version of Korp.

For technical questions or complaints about Korp write to kielipankki@csc.fi.

For any questions on the Oracc data in Korp, contact Aleksi Sahala at firstname.lastname@helsinki.fi.

If you have comments/suggestions concerning this user guide, you can send them to Heidi Jauhiainen at firstname.lastname@helsinki.fi.

# Table of Contents

# 1. SELECTING THE DATASETS

You can choose the Oracc version by going to Korp using one of the urn-addresses on previous page or by going to the **Other languages** page in Korp and by ticking the box next to the name of it in the **Select corpora** dropdown. You will then be searching across all the corpora in that version. To select only one or a few corpora click the arrow left of the name of the chosen version to see the different datasets. You can select and deselect all the sets by ticking the box next to the version's name and, when unselected, choose the sets you are interested in.



Oracc in Korp 2021, contains texts from 29 different Oracc projects.
Oracc in Korp 2019 contains 18 categories, 17 of which correspond to Oracc projects. The category **Other projects** contains texts from several smaller projects:
- Idrimi: Statue of Idrimi
- akklove: Akkadian Love Literature
- Contributions Amarna
- CKST: Corpus of Kassite Sumerian Texts
- Glass: Corpus of Glass Technological Texts
- LaOCOST: Law and Order: Cuneiform Online Sustainable Tool
- OBTA: Old Babylonian Tabular Accounts
- Suhu: The Inscriptions of Suhu online.

**Note** that the *EPSD2: electronic Pennsylvania Sumerian Dictionary* is a huge corpus that may slow down your search. If you do not specifically need Sumerian texts in your search, you might want to untick both EPSD2 and ETCSRI: Electronic Text Corpus of Sumerian Royal Inscriptions. If you do not need the lexical lists, you can untick the quite large DCCLT: Digital Corpus of Cuneiform Lexical Texts.

## 2. SIMPLE SEARCH



Simple search lets you:
- search for the **transliteration** of a word
  - write the transliteration in the search box and hit Search (*Enter* does not work)
- search for a transliteration that
  - starts with the search query
    - check "✓**initial part**" before clicking Search
  - contains the search query anywhere in the word
    - check "✓**compound_middle**" before clicking Search
  - ends with the search query
    - check "✓**final part**" before clicking Search
- use case-insensitive search (e.g. get *lugal* and *LUGAL* in one search)
    - check "✓**case-insensitive**" before clicking Search

These selections can be combined!

You can add a transliteration of a second word to your search! By default, this searches for words **in order**, that is words that occur together in the order indicated. By unselecting the **in order**… box, the words can be anywhere in the document.

You might want to **deselect the box in front of the words Show statistics** (for statistics see section 6). If you use several subcorpora and statistics are calculated, the query time can get quite long.

Depending on your computer you might be able to write some of the special characters used for the transliteration and transcription of Akkadian and Sumerian words, but probably not all. **Here is a list you can copy and paste from:**

---

**Special characters:**

**Ā á ē ī í ū ú š ṭ ṣ ʾ ŋ**

**Ā Á Ē É Ī Ū Ú Ù Š Ṭ Ṣ Ŋ**

**Subscripts:**

$X_1$ $X_2$ $X_3$ $X_4$ $X_5$ $X_6$ $X_7$ $X_8$ $X_9$ $X_{10}$ $X_{11}$ $X_{12}$ $X_{13}$ $X_{14}$

---

## 3. SEARCH RESULTS

The results will show all the instances of the searched word.
- The word(s) searched for will be highlighted and located in the middle of the result list one below another (this format is called Keyword in Context = KWIC).
- You can scroll the screen sideways to see more of the context of the word.



Search results for the transliteration KU₃.BABBAR.MEŠ (see previous page for the special characters). The mouse is hovering over the gray bar revealing that the RIAo project texts have 46 instances of the word.

Above the results, there is a list of page numbers with which you can move to another page. Depending on your browser, you might also see a gray bar showing the division of the results according to the subprojects. You can move directly to the results of a subproject by clicking the bar but note that it can take a while to load.

By default, there will be 25 results per page and the results are grouped by the projects chosen. The way the results are shown can be changed from the blue bar (starting with the word KWIC) below the search box:
- Hits per page: choose 50, 75, 100, 500 or 1000 hits on the page (since changing a page can take a while, using a bigger number of hits per page is recommended)
- Sort within corpora:
  - **not sorted** (default) shows the instances in the order they appear in the corpus
  - **matched words(s)** sorts according to the matches
  - **left context** sorts according to the preceding word
  - **right context** sorts according to the following word
  - **random**

If you change any of these settings, you will have to redo your search by clicking the Search button.

The results can also be seen in a larger context by clicking **Show context** to the right of the list of pages (or under the list if your window is narrow). You will then see the whole text of each document in the search results. You can return to the default view with **Show KWIC**.

## Rīm-Anum: The House of Prisoners (Oracc 2021)

40 {GI}SA-HI.A a-na x x nam-ha-ar-ti {m}i-na-pa-li-šu ZI.GA ŠA₃ NA.KAM.TUM {ITI}x-x mu ri-im-{d}a-nu-um **lugal** unug{KI} x x x ARAD {d}NIN.SI₄.AN.NA u₃ {d}GU.LA

2(BARIG) ZI₃.DA a-na ŠUKU LU₂ DU₁₀.GAR{KI} u₃ a-hi-a-tum ZI.GA ŠA₃ E₂ a-si-ri NIG₂.ŠU {d}EN.ZU-še-mi {ITI}APIN.DU₈.A U₄ 11 mu ri-im-{d}a-nu **lugal** {d}EN.ZU-i-din-nam DUMU i-nu-x-x ARAD AN.AN.MAR.TU

1(BARIG) 4(BAN₂) a-na GEŠBUN UGULA MAR.TU-MEŠ u₃ a-hi-a-tim ZI.GA ŠA₃ E₂ A.SI.RUM NIG₂.ŠU {d}EN.ZU-še-mi {ITI}SIG₄.A U₄ 4-KAM mu ri-im-{d}a-nu-um **lugal** unug{KI} u₃ a₂-dam-bi a-pil-{d}MAR.TU DUMU {d}ŠUL.GI-x ARAD NIG₂ {d}MAR.TU na-bi-i₃-li₂-šu BISAG.DUB.BA DUMU la-ki-ta-re-me-ni ARAD ri-im-{d}a-nu-um

{m}a-wi-li-ia LU₂ {IRI}x-x-x ŠA₃ 6 {LU₂}A.SI.RUM ša i-na ta-ar-ba-şi ša GU₂ {ID₂}GIŠ.GI.NIM x {m}i₂-li₃-e-mu-qa₂-šu {m}ni-in-x-x 2 {LU₂}A.SI.RUM ša i-na {IRI}BAD₃{KI}

When you click any word (i.e. not only the searched word) in the result list, you can see information about the corpus, word, and document in the **sidebar**.



**CORPUS** section has several links:
   Credits > the corpus in Oracc
   Metadata > MetaShare with information on the whole Oracc in Korp data
   information page > Oracc in Language Bank of Finland
   Link to corpus in Korp
   Persistent identifier of the MetaShare page
   Licence
**TEXT ATTRIBUTES** section has information on the document according to the information available in ORACC (see Section *4. Extended Search* for the attributes and their explanation).
**WORD ATTRIBUTES** section has information on the word selected (see Section *4. Extended Search* for the attributes and their explanation).

## 4. EXTENDED SEARCH

In addition to transliteration (word), the extended search tab lets you search by:
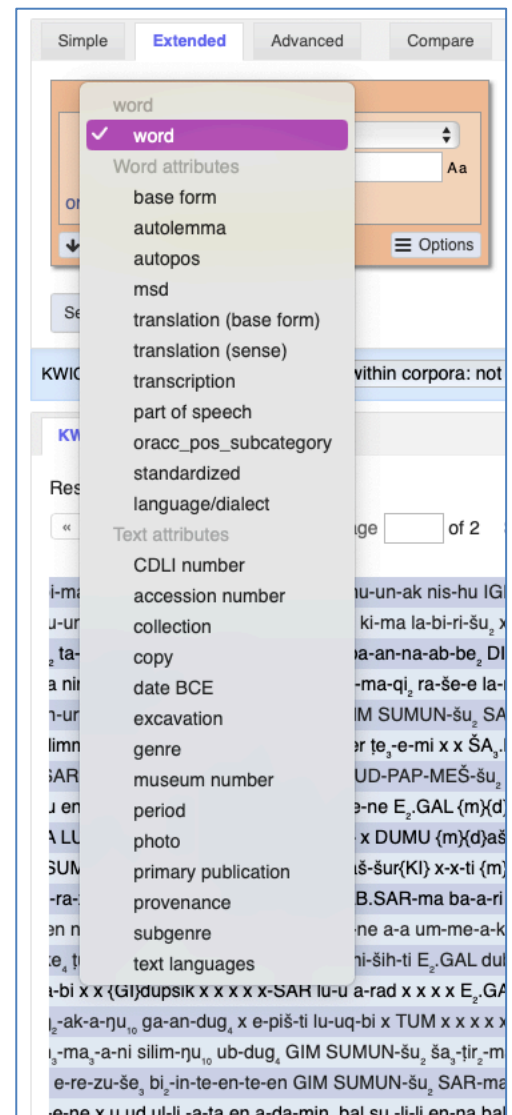- word attributes
- text attributes

In the orange box, click the field that says "word" and you get a dropdown list of attributes. For the Oracc explanations of the attributes see here.
The word attributes are:
- word = transliteration
- base form = dictionary headword in *A Concise Dictionary of Akkadian* (The Citation Form in Oracc)
  - if you want to search for a compound word (e.g. *mār bārê*), use "&&" to combine the base forms of the words (māru&&bārû)[1]
- autolemma = base form produced with a tool called BabyLemmatizer at the University of Helsinki (https://github.com/asahala/BabyLemmatizer)
- autopos = part-of-speech tag produced with BabyLemmatizer
- msd = morphosyntactic descriptors for Sumerian
- translation (base form) = first translation of the dictionary form in CDA (The Guide Word in Oracc)
- translation (sense) = an optional meaning of the word in the context (Sense in Oracc)
- transcription = the normalization for the form (Normalization in Oracc)
- part of speech = the basic grammar of words, i.e. word classes, pronouns, and other tags used in Oracc.
- oracc_pos_subcategory = An optional part-of-speech tied to the current syntactic context (The Effective Part-of-Speech tag in Oracc)
- standardized = normalized form of names of gods and places. See https://github.com/anee-helsinki/OraccInKorp/tree/master/VersionMay2019 for the lists used (the same lists have been used for the Oracc in Korp 2021 version). Not in Oracc.
- language/dialect

The text attributes are:
- CDLI number: P-numbers and their matching designations have been assigned by the CDLI. Q-numbers refer to Oracc's central Q-catalogue of composite texts.
- accession number
- collection
- copy (of)
- date BCE
- excavation
- genre
- museum number
- period
- photo
- primary publication
- provenance
- subgenre (as defined in Oracc)
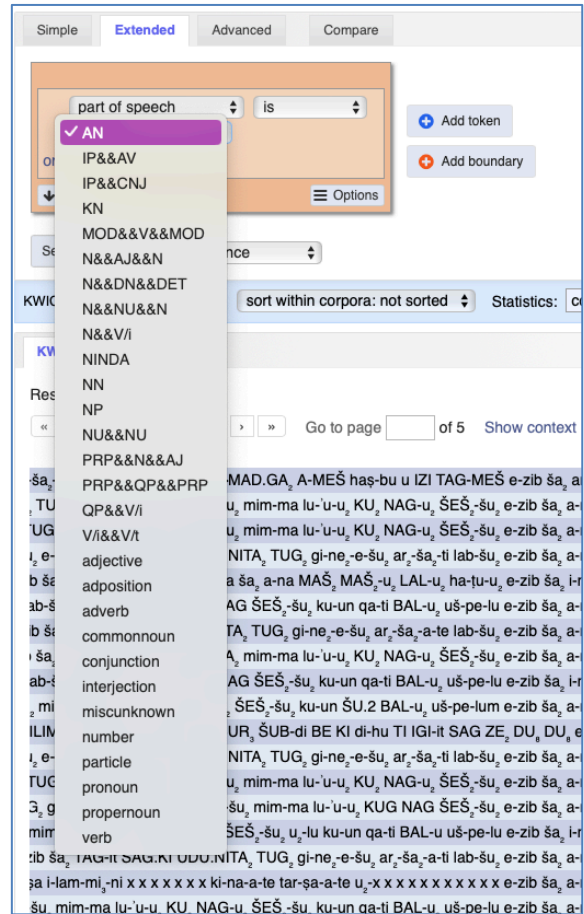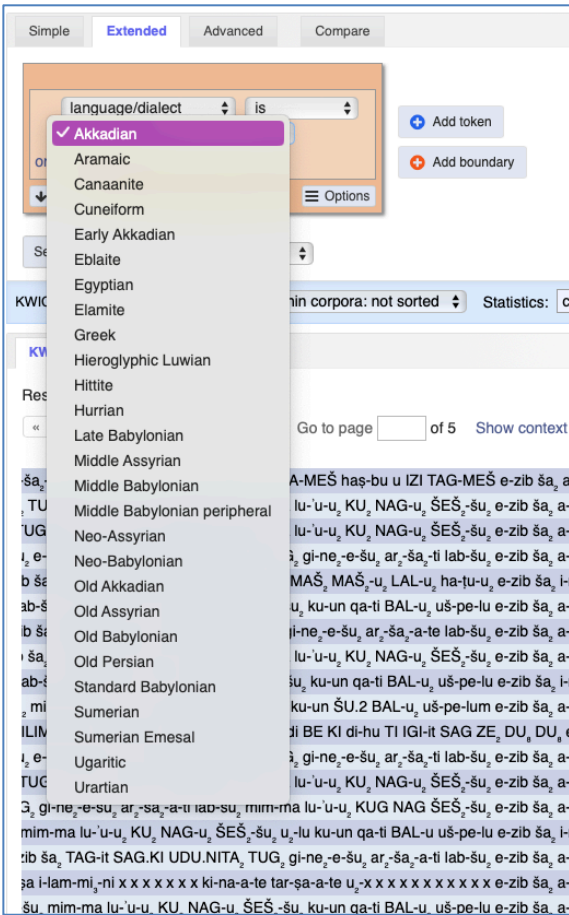- text languages, i.e. languages assigned to the whole text

---

[1] Note that the way compound words are written in different projects in Oracc varies. The joined words which in Oracc have separate translations, word classes, etc. have in Korp been joined with "&&". Sometimes compound words in Oracc have been defined as one word and the parts have been joined with "-", e.g. *EN-MU.MU*, baseform *bēl-zakār-šumi*. Sometimes the parts of a compound word are all defined separately.

The search field underneath will change into a dropdown list of possible values if you choose
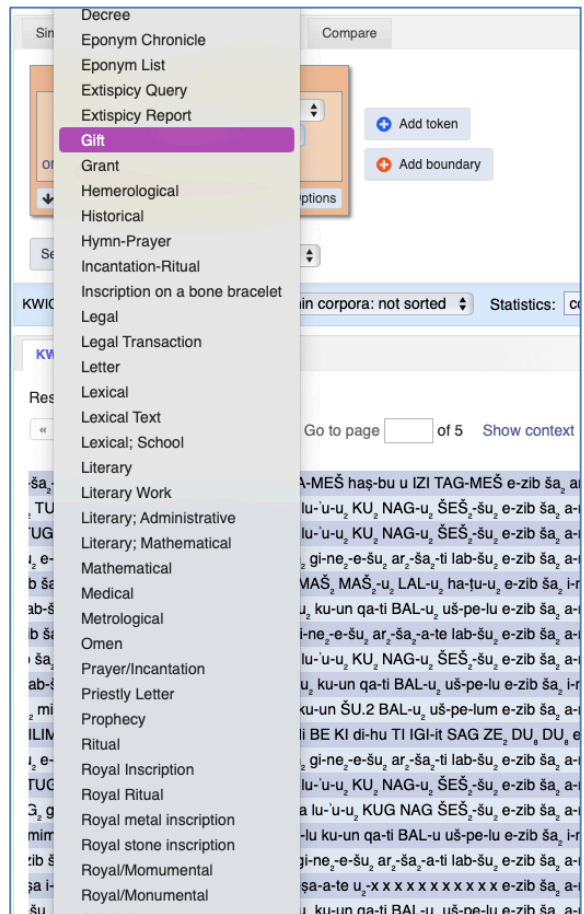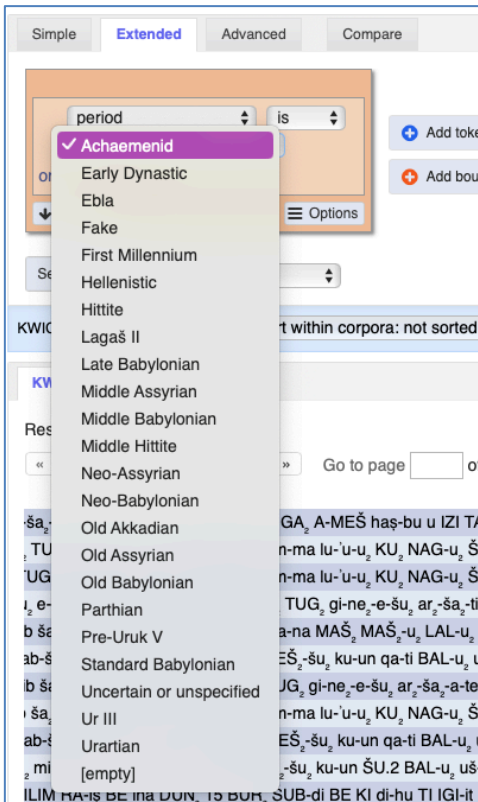
- part-of-speech

- language/dialect



- period



- genre

In some of the other cases, you get a list when you click the text box and start writing. In other cases, such as the CDLI-number, you have to write the value you are looking for in the field. In the sidebar, you can find hints on what the different attributes can contain.

**In all cases,** you can choose whether you want that the search term
- is
- is not

in the results.
When there is no dropdown list, you can also specify if you want that the searched word
- begins with
- contains
- ends with

the given value.
or that it
- is
- is not

the regular expression (regexp) given.

The regular expression is explained in the last section Advanced Search and is only of interest to those who want to do more complicated searches. It is possible to do complicated searches already using the Extended search even without the regular expressions (see next section).

**Note** that although Korp in Extended search gives you an option (right of the search button) of searching "within sentence", in practice the searches in Oracc corpora are all within the whole text.
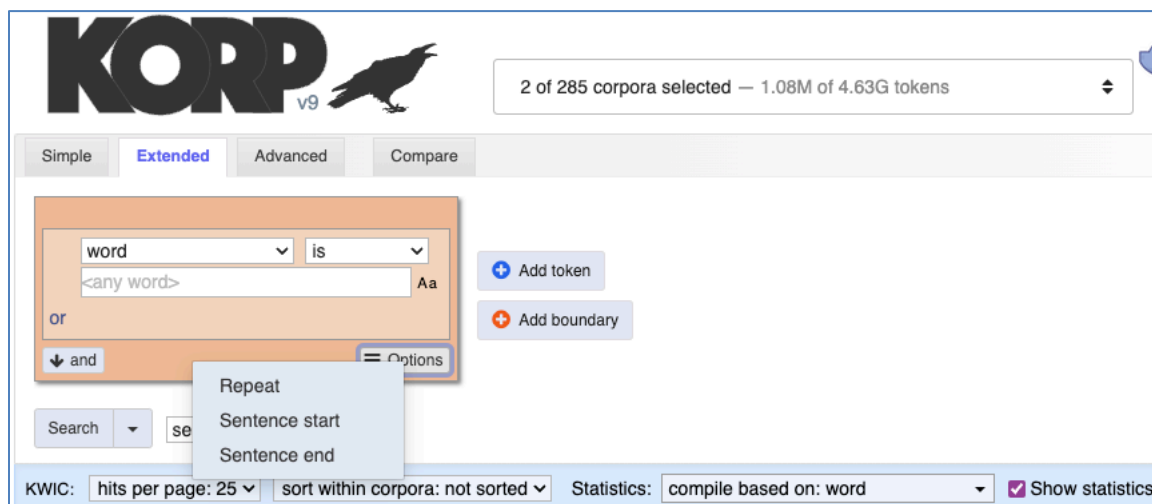
## 5. MAKING MORE COMPLICATED QUERIES IN THE EXTENDED SEARCH

Click the ≡ Options button on the right lower corner of the orange search box to find more options
- **Repeat** – find words that occur several times right after each other (for how to use this, see Search example 2)
- **Sentence start** – find word instances that are only at the beginning of a text
- **Sentence end** – find word instances that are only at the end of a text

These sentence boundaries will open a Boundary unit box before or after the search box

These can also be added by clicking the ⊕ Add boundary button and choosing ←**first** or **last**→



The orange box is the heart of the Extended search and represents one word.

Click **or** in the left bottom corner of the search box and specify another word to get the attestations of both

- e.g. find all words where word (i.e. transliteration) is either *lugal* or *LUGAL*

Click the ⬇ **and** button in the lower part of the box and specify an attribute that the word must have to narrow down your search (note that searching for the same attribute in both options will not give any results)

- e.g. find words the normalization (i.e. transcription) of which is *kaspa* and the word (i.e. transliteration) is not KU₃.BABBAR

Click the ⊕ Add token sign to the right of the box and specify another word that has to follow the first one

- e.g. base form (i.e. dictionary form) is *šarru* + base form is *dannu*

*All these options can be combined as many times as you want!*

*You can remove additional words by clicking the x in the top left corner of an orange box. The word attributes added by clicking or-button can be removed by clicking the – on the left side of the box. These options are only available when words or attributes have been added.*

**Search example 1**

Words (i.e. transliterations) that begin with the letter "a" and where the part of speech is verb and the language/dialect is Akkadian.

**DO:**
  word       begins with     a
**and**
  part of speech     is  verb (from the dropdown)
**and**
  language/dialect  is  Akkadian (from the dropdown)
**Search**

**Search example 2**

Texts from the Neo-Assyrian period featuring base form (i.e. dictionary form) of the divine name "Aššur" and the translation "king" with no more than 8 words in between.

**DO:**

base form   is   Aššur

**and**

oracc_pos_subcategory     is   DN Divine Name

**and**

period        is   Neo-Assyrian

**+Add token**

word      is   <any word>

Options: Repeat

Repeat 0 to 8 times

**+Add token**

translation (base form) is   king

**Search**

## Search example 3

Word (i.e. transliteration) "DINGIR-MEŠ" **or** "DINGIR" followed by (+Add token) any adjective, i.e. part of speech is adjective (from the dropdown).



## Search example 4

Base form (i.e. dictionary form) "eqlu" in texts where provenance is Uruk (start writing Uruk), period Hellenistic (from the dropdown), and genre Legal (from the dropdown).

## 6. STATISTICS

The Statistics tab to the right of the table called KWIC in the search results gives the number of occurrences for each matched word both in all results (Total) and within individual corpus/dataset. For the statistics to be calculated the box before the words **Show statistics** in the blue bar (starting with the word KWIC) below the search box **must be checked** before performing the search.

- The number of occurrences is shown as relative frequencies per million tokens, a common measure in corpus linguistics.
- The numbers in parentheses are the absolute frequencies (i.e. the number of occurrences).
- The default view shows the statistics of the transliteration(s) of the word searched regardless of what attributes were searched for.
- You can change what attribute(s) are considered by selecting the attributes you want from "Statistics: (compile based on:)" in the blue bar and then clicking Search. You can even choose several at once (when you change the selection, hit Search again):
  - e.g. search for all occurrences of words the translation of which contains "love" and base the statistics, for example, on base form, part-of-speech, and language/dialect to see what combinations there are in the results.
- You can sort the statistics according to any column by clicking the heading of that column. The columns with words will be sorted alphabetically and the numbers numerically. Clicking a second time will reverse the sorting.

You can see the texts of an individual line in the statistics by clicking a word on that line. You will get a new tab with the KWIC view of those results within your original search.


### Statistics example

Search for the base form *dannu* and base the statistics on word (= transliteration), translation (of base form), and part of speech. Click Search. Go to Statistics tab and order the data according to translation by clicking the heading *translation*.

| word | translation (base form) | part of speech | Total | ADsD: Astr... | AEMW: Ak... | Akkadian ... | ARIo: Ach... | ATAE: Arc... | blms: Bilin... | BT |
|---|---|---|---|---|---|---|---|---|---|---|
| dan-nu | (large) vat | commonnoun | 0.9 (6) | 0 (0) | 0 (0) | 0 (0) | 0 (0) | 0 (0) | 0 (0) | 0 (0 |
| dan-nu-tu | (large) vat | commonnoun | 0.1 (1) | 0 (0) | 0 (0) | 0 (0) | 0 (0) | 0 (0) | 0 (0) | 0 (0 |
| {DUG}dan-nu | (large)vat | commonnoun | 0.1 (1) | 0 (0) | 0 (0) | 0 (0) | 0 (0) | 0 (0) | 0 (0) | 0 (0 |
| dan-nu | strong | adjective | 130.9 (903) | 52.6 (12) | 11.1 (2) | 0 (0) | 170.6 (2) | 201.3 (11) | 84.5 (5) | 0 (0 |
| dan-nu-ti | strong | adjective | 40.1 (277) | 0 (0) | 0 (0) | 0 (0) | 0 (0) | 36.6 (2) | 16.9 (1) | 0 (0 |
| dan-ni | strong | adjective | 24.2 (167) | 0 (0) | 0 (0) | 0 (0) | 0 (0) | 54.9 (3) | 0 (0) | 0 (0 |
| KALAG | strong | adjective | 16.2 (112) | 437.9 (100) | 0 (0) | 0 (0) | 0 (0) | 36.6 (2) | 0 (0) | 0 (0 |
| dan-nu-te | strong | adjective | 8.1 (56) | 0 (0) | 0 (0) | 0 (0) | 0 (0) | 18.3 (1) | 0 (0) | 0 (0 |
| dan-nu-ti-šu₂ | strong | adjective | 7.2 (50) | 0 (0) | 0 (0) | 0 (0) | 0 (0) | 0 (0) | 0 (0) | 0 (0 |
| KAL.MEŠ | strong | adjective | 6.7 (46) | 0 (0) | 0 (0) | 0 (0) | 0 (0) | 0 (0) | 0 (0) | 0 (0 |
| da-nu₄-tim | strong | adjective | 5.7 (39) | 0 (0) | 0 (0) | 0 (0) | 0 (0) | 0 (0) | 0 (0) | 0 (0 |
| da-num | strong | adjective | 5.5 (38) | 0 (0) | 0 (0) | 0 (0) | 0 (0) | 0 (0) | 0 (0) | 0 (0 |
| dan-na-at | strong | adjective | 4.9 (34) | 0 (0) | 0 (0) | 0 (0) | 0 (0) | 0 (0) | 0 (0) | 0 (0 |
| KALAG-MEŠ | strong | adjective | 4.8 (33) | 0 (0) | 0 (0) | 0 (0) | 0 (0) | 0 (0) | 0 (0) | 0 (0 |
| · | strong | adjective | 4.5 (31) | 0 (0) | 0 (0) | 0 (0) | 0 (0) | 0 (0) | 0 (0) | 0 (0 |
| dan-na | strong | adjective | 3.9 (27) | 0 (0) | 5.6 (1) | 0 (0) | 0 (0) | 0 (0) | 16.9 (1) | 0 (0 |

# 7. MAP

The Statistics tab also allows viewing the provenance of the documents with the searched words in a map. You can use it for all kinds of searches, but it is best suited for comparing words or word forms and their use in different parts of the ancient Near East (but remember that a text was not necessarily written in the place it was found in).

**To use the map, make a search with the Show statistics selected**. It is best, but not necessary, to compile based on just one attribute such as base form. Then **in the Statistics tab select the words you want to see in the map**. You can easily select all rows by selecting the checkbox in the heading row. If you have searched for a transliteration or based your search on several attributes, you might want to choose only some of the rows to see on the map. **Click the Show map and then select Show map**.

**Map example**

Search for the base forms of *uqnû (i.e. lapis lazuli)*, *elmēšu* (amber), *annaku* (tin), or *siparru* (bronze)
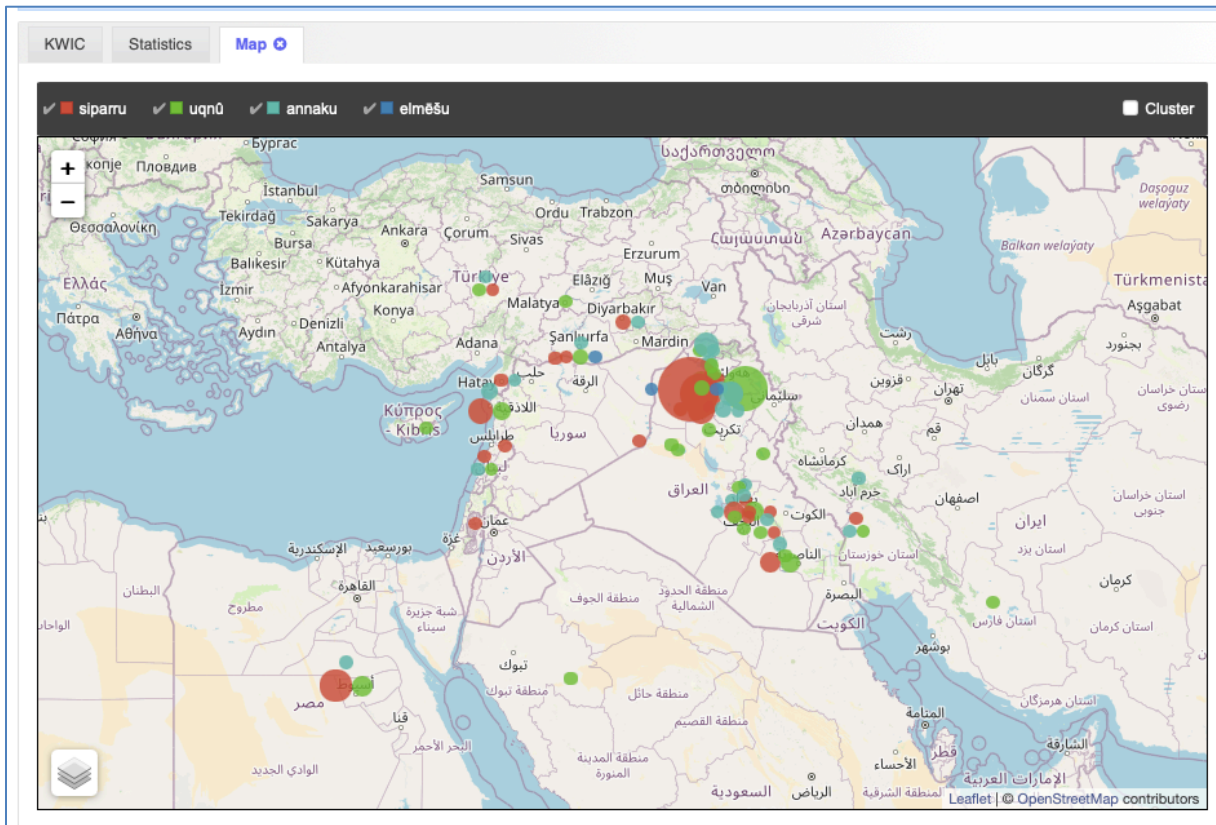


Example by Aleksi Sahala.

In the map view, you can select and unselect the words shown in the black bar above the map. By checking the box next to the word Cluster in the right upper corner, you will see provenances for each site as a bar diagram. By hovering with the mouse over a place/circles/bars you will see statistics for the words in that location.



Map where the box next to Cluster is ticked and only 3 of the words is selected (in the black bar). The mouse is hovering over the place Uruk revealing the statistics for the selected words in that location.

## 8. COMPARING RESULTS

You can save searches for comparison with each other.
When you have typed in your search, click the arrow in the search button and give your search a name. When you have at least two saved queries, go to the Compare tab (above the search box, to the right of the simple and extended search tabs) and choose from the dropdown lists:
- queries to compare
- what attribute the comparison is based on

Click Compare.
The results will be shown in a new tab, thus your latest search is still active.

### Comparison example

Search for two consecutive words with the translations (sense) "king" and (+Add token) "strong" (bear in mind that adjectives follow nouns in the Akkadian word order). After the search results have loaded, click the arrow in the search button, name the search Strong_king (in the box that opens) and click Save.

Search for two consecutive words with the translations (sense) "king" and "great". Name it Great_king and Save.

The Compare tab will now have the number 2 on its title (the number is higher if you have saved more searches). Go to the Compare tab and choose to Compare the saved query Strong_king with Great_king and compile based on transcription. Click Compare.



Comparison of the number of instances of the different transcriptions of the phrases 'strong king' and 'great king' in Oracc in korp 2021 (excluding DCCLT, EPSD2, and ETCSRI).

There is a button to the right of the comparison options to Delete saved searches. You do not have to do this, since you can name different searches with different names and use them also later by choosing them from the dropdown once you are in the Compare tab. If you close the window (or even the browser), the cookies will remember your saved searches when you come back to the *Other languages* section of Korp (i.e. not Finnish, Swedish, or Parallel).

## 9. EXPORTING THE RESULTS

You can export the results of your search in many different formats (the file formats supported in most cases are csv (colon-separated values), tsv (tab-separated values), and xls (excel)):

ANNOT_ = text as a table, one token (with annotations) per row
REF_ = a list with the transliteration of the text and the text attributes
SENTENCE_ = one text per row with transliteration, base forms, and text attributes + more
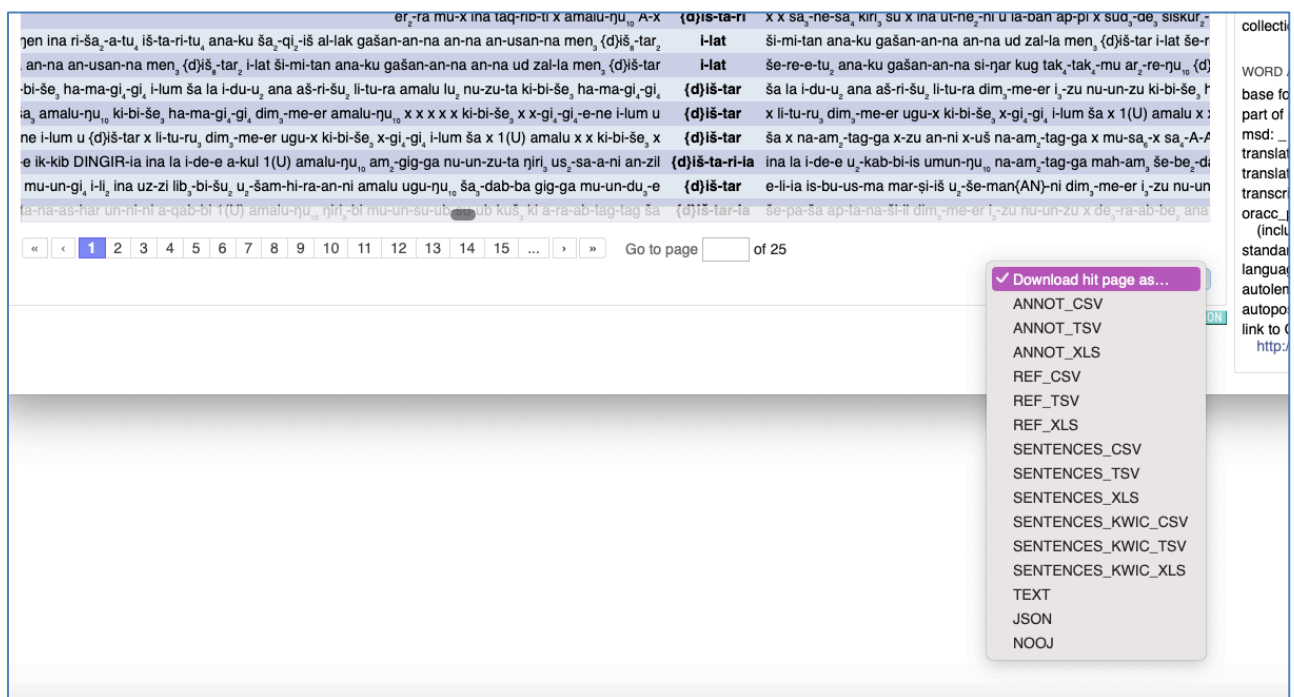SENTENCE_KWIC_ = one text per row with transliteration, base forms, and text attributes +
TEXT = transliteration, one text per row
JSON = Structured format with text attributes and annotated words in attribute_name:attribute format
NOOJ = xml format for use in the NooJ-tool

Click the *Download hit page as…* in the lower right corner of the results area and from the dropdown select the desired format and file format by clicking a row. You do not have to do anything else! It may take a while before you see the prompt for where you want to save the file (in Windows) or the file starts downloading to your Downloads folder (in Mac).



The turquoise JSON button underneath the dropdown will open the JSON file in a new tab with all the annotated words of a text on one line.

You can also export the **Statistics** data. At the bottom of the Statistics tab select whether you want to export relative or absolute frequencies and whether you want to have the data in the csv or tsv format and hit **Generate export** and then **Export**. The JSON button here will open the statistics data in a new tab.

**The Extended search is sufficient for most searches**! The next section introduces the Advanced search function for even more complicated searches, and you might want to take a quick look at it especially if you are familiar with any programming languages. Or you can come back to it later if you feel like it. Anyway, **you should now be ready to use Oracc in Korp**. The best way to learn is by trying things out.
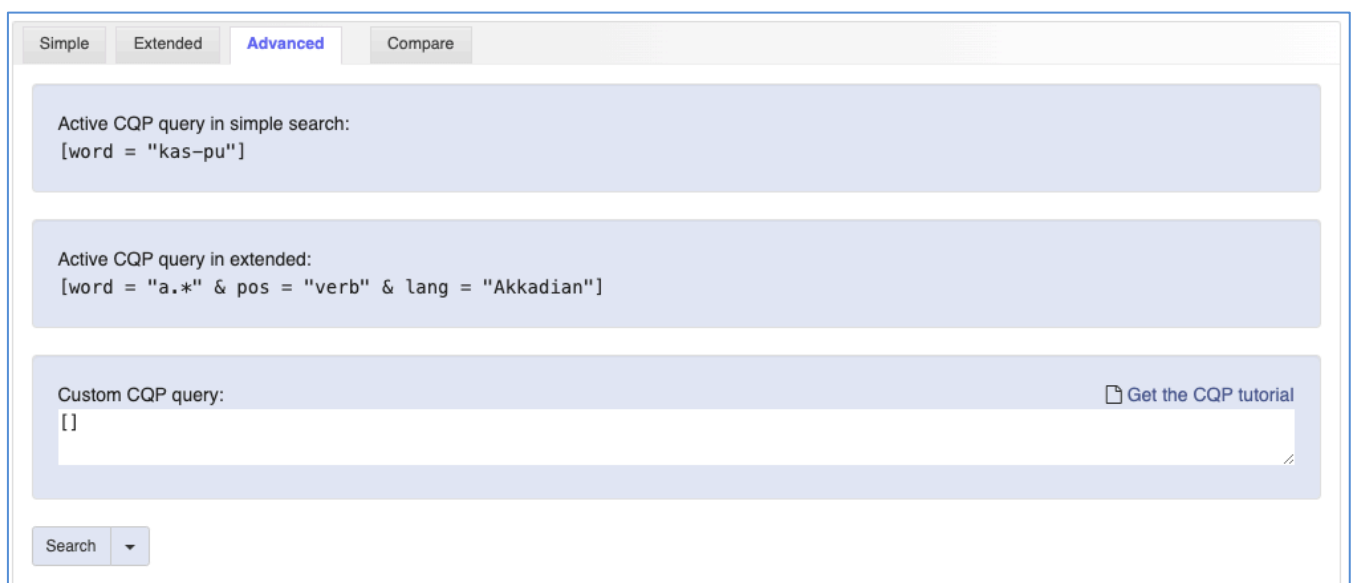Have fun!


## 10.    ADVANCED SEARCH

For a general guide for the advanced search in Korp see
https://www.kielipankki.fi/support/korp-advanced/

What you cannot do with Extended search is combine two searches. For example, in Search example 2 above, we performed a search in texts from the Neo-Assyrian period containing the base form of the divine name "Aššur" and the translation "king" with no more than 8 words in between. With the Advanced search, you can add to the same search the cases where the translation "king" is *before* the divine name "Aššur".

The Advanced search is performed by writing the query in the so-called CQP (Corpus Query Protocol) query language. If you want to use the Advanced search, it may take some learning.

The easiest way to start learning CQP is to perform a Simple or Extended search and then check in the Advanced tab what it looks like in CQP. In the Advanced tab, you will see the Active CQP query of the latest Simple and Extended searches.



Advanced search tab after performing a Simple search "word (i.e. transliteration) is kas-pu" and the Extended search from Search example 1 above (Words (i.e. transliterations) that begin with the letter "a" and where the part of speech is verb and the language/dialect is Akkadian).

Brackets [] always represent an orange box in the search field, empty brackets stand for any word.

The attributes have slightly different names in CQP than in the Extended search:

|  | *CQP* | *Extended search* |
|---|---|---|
| **Word attributes** | word | word |
|  | lemma | base form |
|  | translation | translation (base form) |
|  | transcription | transcription |
|  | sense | translation (sense) |
|  | pos | part-of-speech |
|  | oraccpos | oracc_pos_subcategory |
|  | standard | standardized |
|  | lang | language/dialect |
| **Text attributes** | cdlinumber | CDLI number |
|  | genre | genre |
|  | provenance | provenance |
|  | period | _.text_period |
|  | subgenre | subgenre |
|  | language | text languages |

For example, in Search example 2 above, we performed a search in texts from the Neo-Assyrian period containing the base form of the divine name "Aššur" and the translation "king" with no more than 8 words in between. In the CQP language, the search looks like this:

[lemma = "Aššur" & oraccpos = "DN Divine Name" & _.text_period = "Neo-Assyrian"] []{0,8} [translation = "king"]

The CQP uses mathematical, boolean, and regular expressions used in most programming languages. For example:

[lemma = "Aššur"]         words where base form **is** *Aššur*
[lemma **!=** "Aššur"]        words where base form **is not** *Aššur*
[lemma = "**.***Aššur**.***"] words where base form **contains** *Aššur* ('**.***' stands for a letter zero or more times)
[lemma = "Aššur**.***"]        words where base form **starts with** *Aššur*
[lemma = "**.***Aššur"]        words where base form **ends with** *Aššur*
[lemma = "Aššur" **&** oraccpos = "DN Divine Name"]    words where dictionary form is *Aššur* **and** oracc subcategory of the part-of-speech is *DN Divine Name* (must be both to match)
[word = "lugal" | word = "LUGAL"]   words where transliteration is *lugal* **or** transliteration is *LUGAL* (both are a match)
[]{1,3}             one, two, or three words without specifying what word
[word = "LUGAL"]{2,2}     words where transliteration is *LUGAL* right after words where transliteration is also *LUGAL* (e.g. LUGAL LUGAL in the text)

A guide to regular expression used in Korp (from https://www.kielipankki.fi/support/korp-advanced/)

| . | any single symbol | |
|---|---|---|
| […] | a set or range of symbols: any of the symbols inside the brackets | [aeiouyäö] matches a single Finnish vowel symbol and [a-hw-z] all the letters from a to h and w to z. |
| [^…] | the complement of a set or range of symbols, none of the symbols inside the brackets | [^abcw-z] matches any symbol except the letters a, b, c, w, x, y, and z. |
| RS | concatenation: the substring matched by the expression R if followed by a substring matched by the expression S | [a-z][0-9] matches a lowercase letter followed by a digit |
| (…) | grouping | |
| R* | repeat zero or more times; R can be a single character, a set of characters, or parentheses containing a regular expression | a.* matches all strings starting with an a, while a(bc)*matches the strings a, abc, abcbc, abcbcbc etc. |
| R+ | repeat once or more | goo+d matches the strings good, goood, gooood etc. |
| R{n} | repeat exactly n times | |
| R{m,n} | repeat m to n times | |
| R? | optionality (repeat zero or one time) | favou?rite matches favorite and favourite |
| R\|S | alternatives; match R or S | apple\|orange matches the strings apple and orange; (read\|writ\|watch)ing matches reading, writing, watching |
| \c | escaping; escapes a special character | \. matches a literal full stop |

## Advanced search example

Texts from the Neo-Assyrian period featuring base form of the divine name "Aššur" and the translation "king" **or the other way round** (using the | ) with no more than 8 words (any words) in between.
[lemma = "Aššur" & oraccpos = "DN Divine Name"  & _.text_period = "Neo-Assyrian"]
[]{0,8}
[translation = "king"]
|
[translation = "king"]
[]{0,8}
[lemma = "Aššur" & oraccpos = "DN Divine Name" & _.text_period = "Neo-Assyrian"]