

2.1 Sub-corpora

As explained in the [section](#) about the creation of the sub-corpora, you can work with either the full corpus WUS or you can select different sub-corpora that mainly depend on the main language within the chat. You find the list of sub-corpora in the bottom left in ANNIS.

The list of sub-corpora is also a good starting point to get information about available fields for your query, to get examples and statistics.

Tokens and messages per sub-corpus

Next to the name of each sub-corpus, you see the number of messages (marked as "Texts") and tokens. You can use these figures for statistics.

Please note: If you work with corpora where not all participants gave their [permission](#) to use their texts, the figure for tokens is off because messages without permission were replaced by texts like *redactedQ12tokens55characters* . These texts count as tokens, too. If you need statistics that depend on the number of tokens in a (sub-)corpus, you are advised to work with corpora with the extension [_DEMOG](#).

Information about the corpus

When you press on the small <i> for information to the right of each (sub-)corpus name, you receive more information about the corpus. More specifically:

- Some statistic information about the corpus including a link at the bottom for reference
- Information about the version to be quoted in publications.
- If you need specific information about an individual chat, you can select the chat instead of the corpus in the top left to get information such as number of messages, number of speakers, etc. This is also an easy way to see which chats are integrated in this sub-corpus.

Corpus information for WUS_FRA_DEMOG (ID: 51867)

Metadata		Available annotations		
Select corpus/document: WUS_FRA_DEMOG		Node Annotations		
		Edge Annotations		
Name	Value	Edge Types		
editors	Anne Göhring, Beni Ruel, Simone Ueberwasser	Meta Annotations		
longName	whatsUpSwitzerland FRA demographics subcorpus	Name	Example (click to use query)	URL
project	whatsUpSwitzerland	total_msg	node & metac:total_msg="14789"	⚡
releaseDate	20190815	consent_speakers	node & metac:consent_speakers="2"	⚡
shortName	WUS_FRA_DEMOG	contains_eng	node & metac:contains_eng="true"	⚡
version	7.0	demographics	node & metac:demographics="2"	⚡
		doc	node & metac:doc="chat195"	⚡
		content_msg	node & metac:content_msg="6838"	⚡
		empty_msg	node & metac:empty_msg="0"	⚡
		media_msg	node & metac:media_msg="69"	⚡
		lang_less_than_100	node & metac:lang_less_than_100="eng"	⚡
		encrypted_msg	node & metac:encrypted_msg="0"	⚡
		lang_100_and_more	node & metac:lang_100_and_more="fra"	⚡
		speakers	node & metac:speakers="2"	⚡
		no_consent_msg	node & metac:no_consent_msg="0"	⚡

Link to corpus: http://linguistik.sms.uzh.ch/8080/annis-gui/#_c=V1VTX02SQV9ERU1PRw

Figure 1: annotations for a (sub-)corpus

On the right side of the information window, you see which annotations are available to be queried for the selected sub-corpus.

- You have two categories of information: Node Annotations are attributes on either token and message level, that we considered to be basic units. Meta Annotations contain information at the chat level; most of the meta annotations indicate sizes (e.g. total number of messages in a given chat) and were automatically computed.
- To the right of the name of the annotation, there is an example query for that specific annotation. If you click on that text, a sample query is entered into the query field in the main screen. This is the easiest way to generate queries, since you can always modify it in the query field. Example: if you click on `Lang_100_and_more`, an example like `node & meta::lang_100_and_more="fra, eng"` is entered into the query field. This query would search for messages in chats with more than 100 messages in French and in English. More precisely: "node" will fetch also all tokens that are in such chats; if you want to distinguish between messages and tokens, you should explicitly query for one or the other: `tok & ...` ; or `msg & ...`

List of chats in the sub-corpus

By clicking on the little piece of paper next to the information `<i>` in the list of sub-corpora, you get a list of all chats in this specific sub-corpus.

From here, you can click on *complete chat view* to view this whole chat (without any annotations). once in this list of messages, you can always click on an individual message ID to see that message with annotations.

If you click on the little `<i>` at the very right of the list of chats, you see all the meta information about this chat.

From:

<https://whatsup.linguistik.uzh.ch/> -

Permanent link:

https://whatsup.linguistik.uzh.ch/02_browsing/01_sub_corpora?rev=1587118274

Last update: **2022/06/27 09:21**

