



University of  
Zurich <sup>UZH</sup>

URPP Language and Space

# TeDDi Sample: Text Data Diversity Sample for Language Comparison and Multilingual NLP

Steven Moran<sup>1</sup>, Christian Bentz<sup>2</sup>, **Ximena Gutierrez-Vasques**<sup>3</sup>, Olga Sozinova<sup>3</sup>, Tanja Samardzic<sup>3</sup>

University of Neuchatel<sup>1</sup>

University of Tübingen<sup>2</sup>

URPP Language and Space, University of Zürich<sup>3</sup>

LREC 2022



FONDS NATIONAL SUISSE  
SCHWEIZERISCHER NATIONALFONDS  
FONDO NAZIONALE SVIZZERO  
SWISS NATIONAL SCIENCE FOUNDATION



# 1. What is TeDDi Sample?

Motivation

Goal

Data collection and curation

Data availability

# 2. Demo



## What is TeDDi Sample?

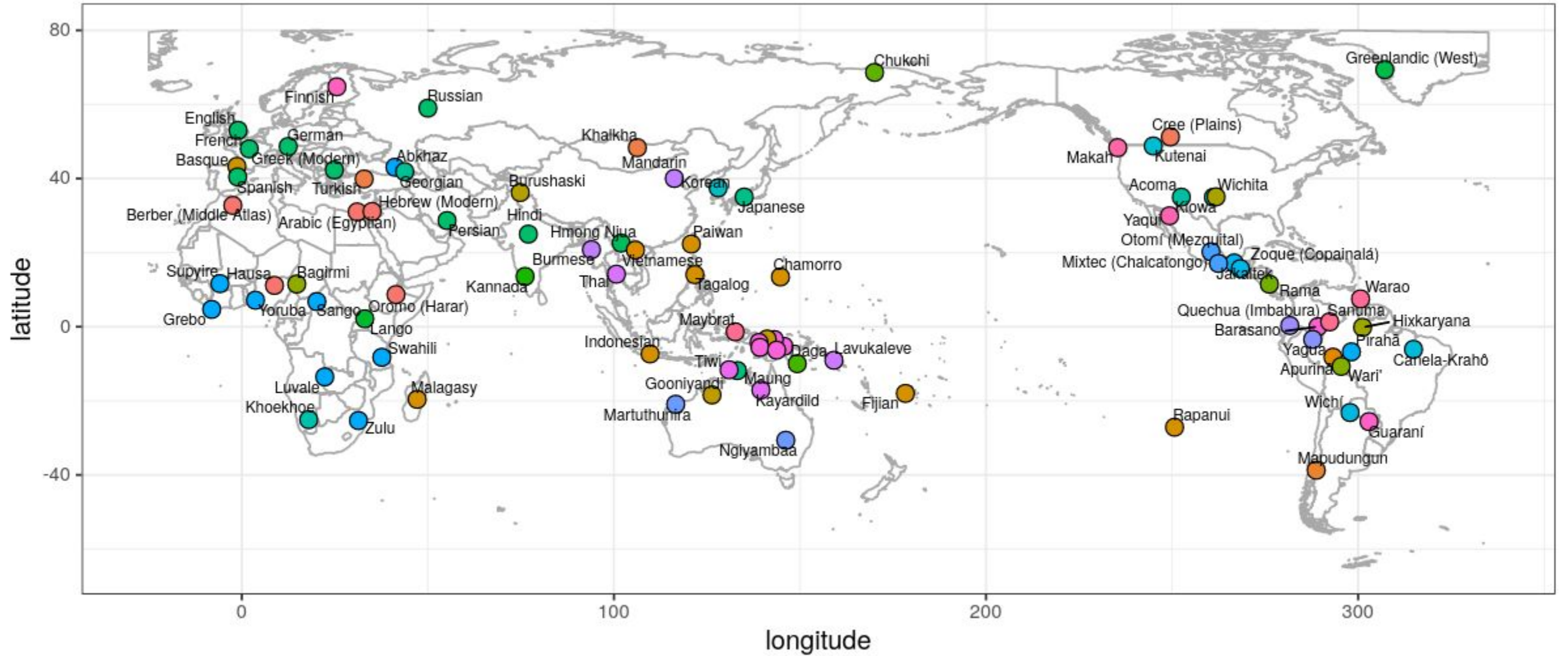
A diversity sample of **text data**. It features:



## What is TeDDi Sample?

A diversity sample of **text data**. It features:

- Text corpora for languages on the typological diversity sample in the [World Atlas of Language Structures \(WALS\)](#)





## What is TeDDi Sample?

A diversity sample of **text data**. It features:

- Text corpora for languages on the typological diversity sample in the [World Atlas of Language Structures \(WALS\)](#)
- Currently **89 languages**. More than **20k texts**



## What is TeDDi Sample?

A diversity sample of **text data**. It features:

- Text corpora for languages on the typological diversity sample in the [World Atlas of Language Structures \(WALS\)](#)
- Currently **89 languages**. More than **20k texts**
- **Open-source** corpus processing tools



## Motivation

- Access to data from minority and under-resourced languages leads to deeper and more complete understanding of language
- Multilingual datasets are becoming increasingly important in NLP





## Motivation

**How to select languages?**



## Motivation

### How to select languages?

#### Linguistics perspective

- More weight on representing a wide range of language families and areas, structural features

#### NLP Perspective

- Favors languages for which text data is readily available online.



## Goal

The aim of TeDDi is to facilitate the use of text-based quantitative methods for analyzing linguistic diversity in both [linguistic research](#) and [NLP](#):



## Goal

The aim of TeDDi is to facilitate the use of text-based quantitative methods for analyzing linguistic diversity in both [linguistic research](#) and [NLP](#):

- Approximate the diversity of languages across the world by means of text samples
- Complement the existing knowledge about the structure of languages, which mostly consists of high-level feature-value pairs stored in linguistic databases



## Data collection and curation



Based on [WALS 100-language sample](#)

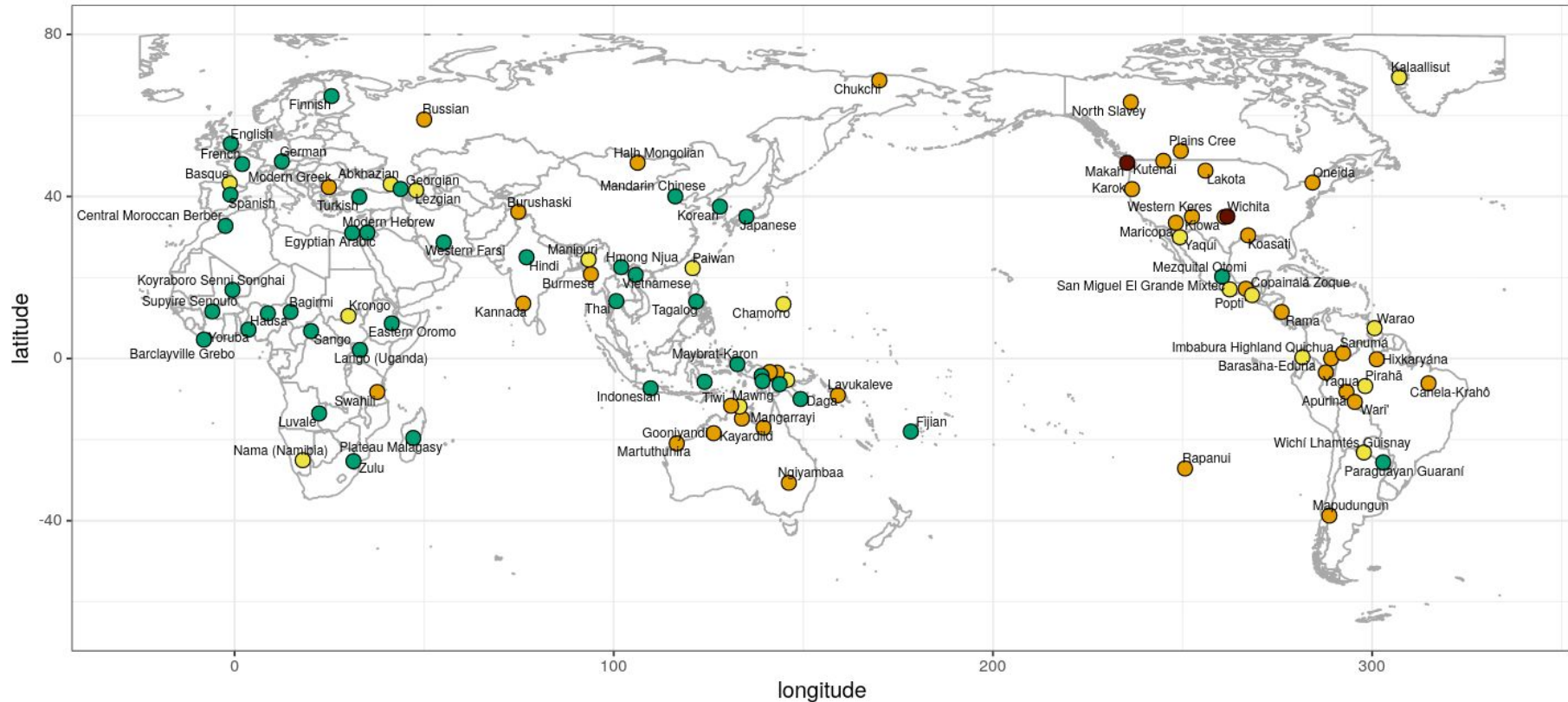
Maximizes [genealogical](#) (language family) and [areal](#) (geographic) diversity

Selection is [independent](#) of [text data availability](#)



# Data collection and curation

Endangerment Status ● endangered ● extinct ● safe ● vulnerable





## Data collection and curation

### Existing resources

- Project Gutenberg, Open Subtitles (Lison and Tiedemann, 2016), The Parallel Bible Corpus (Mayer and Cysouw, 2014), the Universal Declaration of Human Rights



## Data collection and curation

### Existing resources

- Project Gutenberg, Open Subtitles (Lison and Tiedemann, 2016), The Parallel Bible Corpus (Mayer and Cysouw, 2014), the Universal Declaration of Human Rights

### Manually collected translations, transcriptions, and grammatical annotations





## Data collection and curation

### Existing resources

- Project Gutenberg, Open Subtitles (Lison and Tiedemann, 2016), The Parallel Bible Corpus (Mayer and Cysouw, 2014), the Universal Declaration of Human Rights

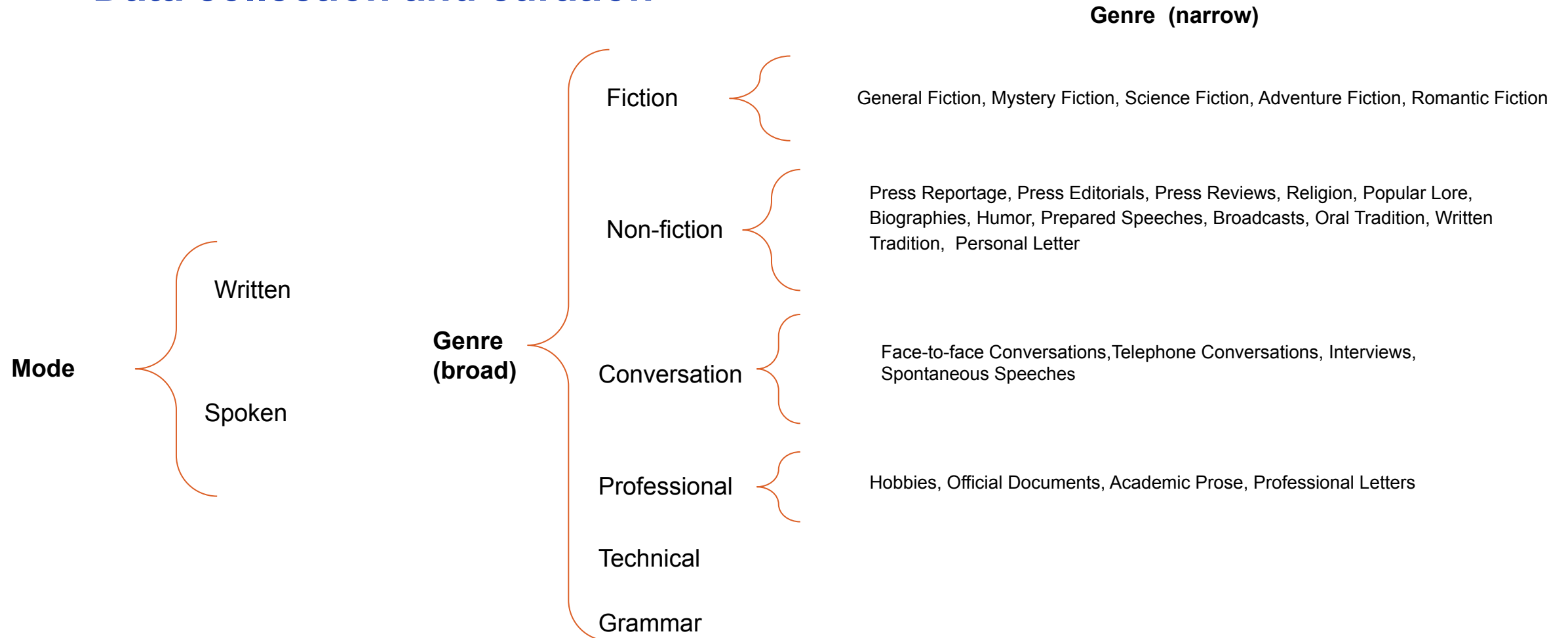
### Manually collected translations, transcriptions, and grammatical annotations

### Sampling (for rich-resource languages)

- We do not include all of the texts available. Instead, we create smaller samples limiting the maximal size of a text unit to 50k tokens of contiguous text

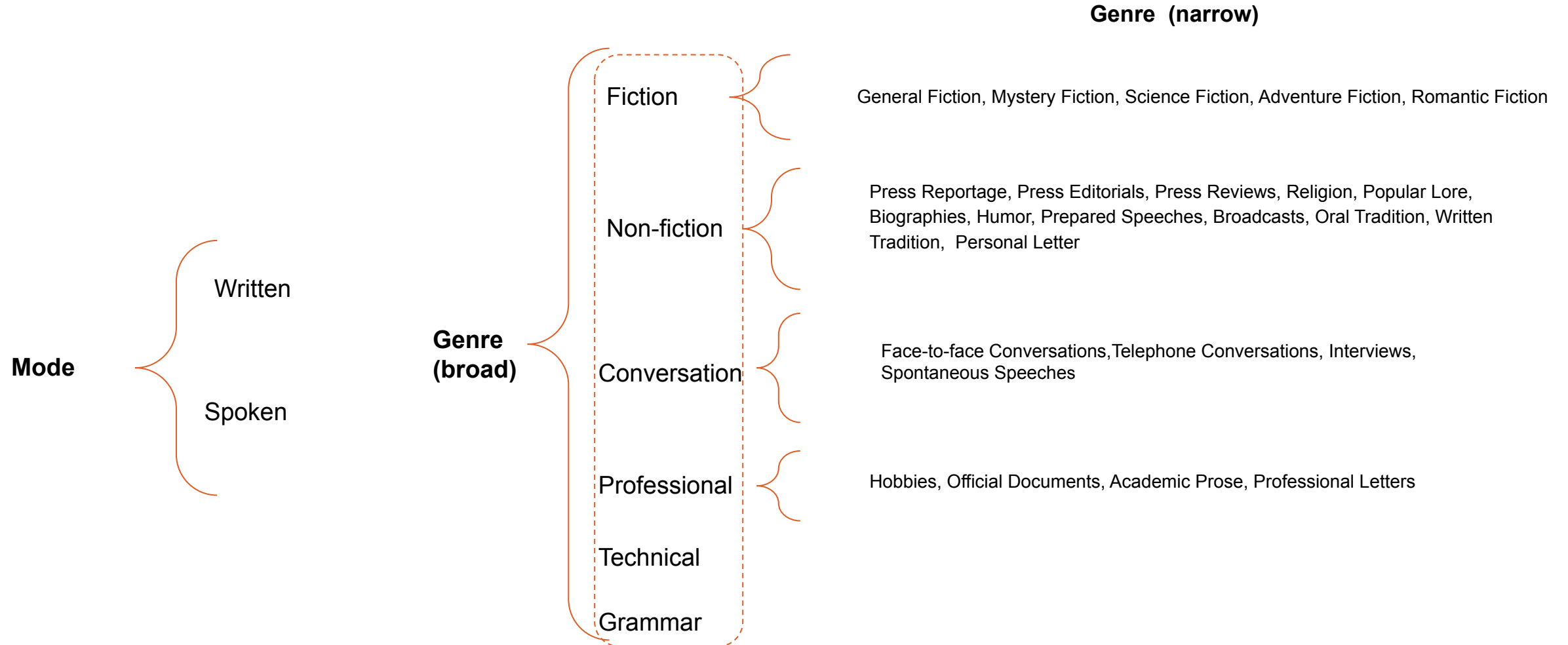


# Data collection and curation





# Data collection and curation





## Data collection and curation

Genre	Langs*	Tokens	Scriptst†
conversation	10	15,835	1
fiction	12	36,811,339	7
grammar	5	1271	1
nonfiction	73	101,588,748	13
professional	40	80,092	15
<b>Total</b>	<b>89</b>	<b>ca. 138 million</b>	<b>16</b>

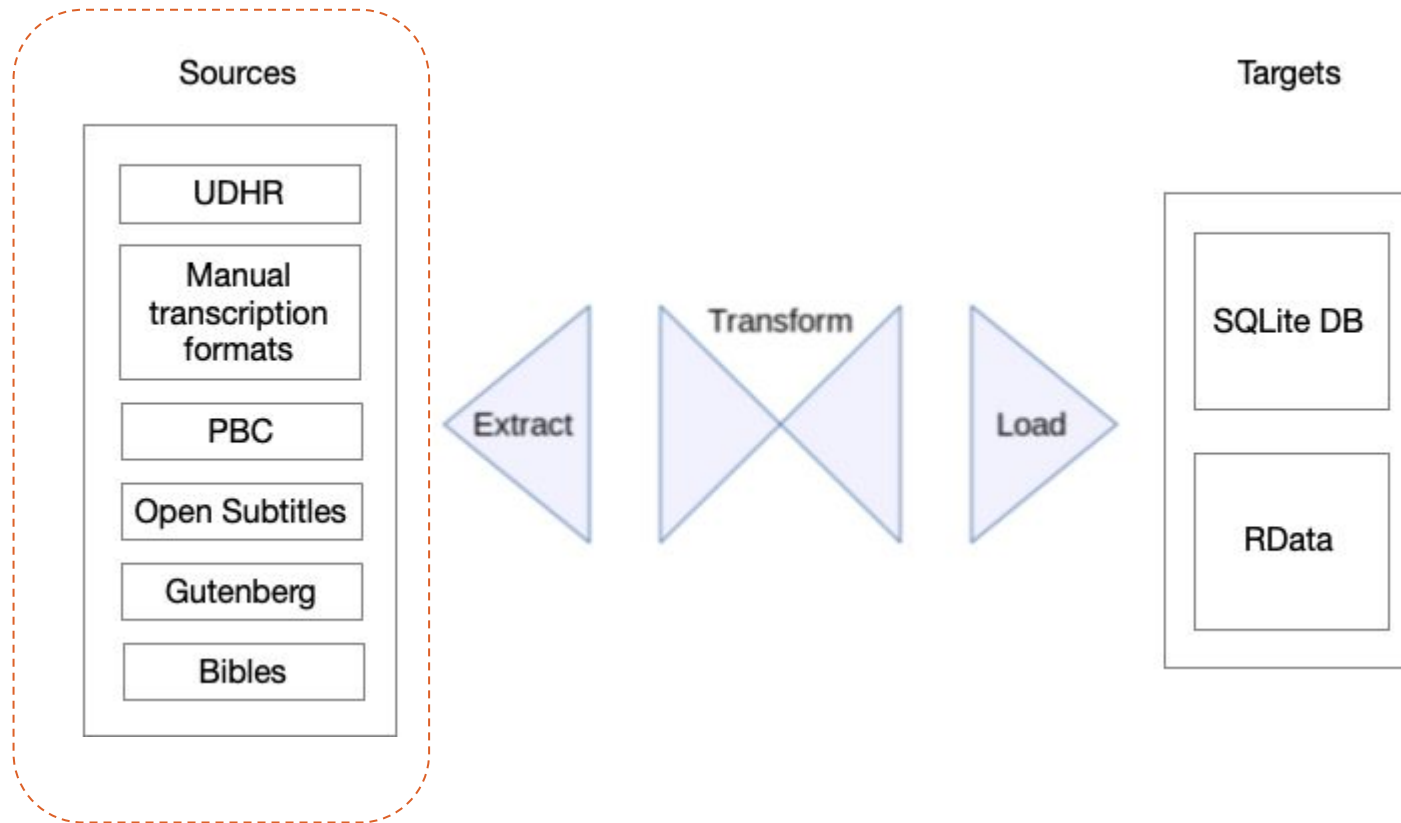
- 89 different languages, stemming from 58 language families (according to WALS)

\*According to ISO-639-3 codes.

†According to ISO-15924 codes.

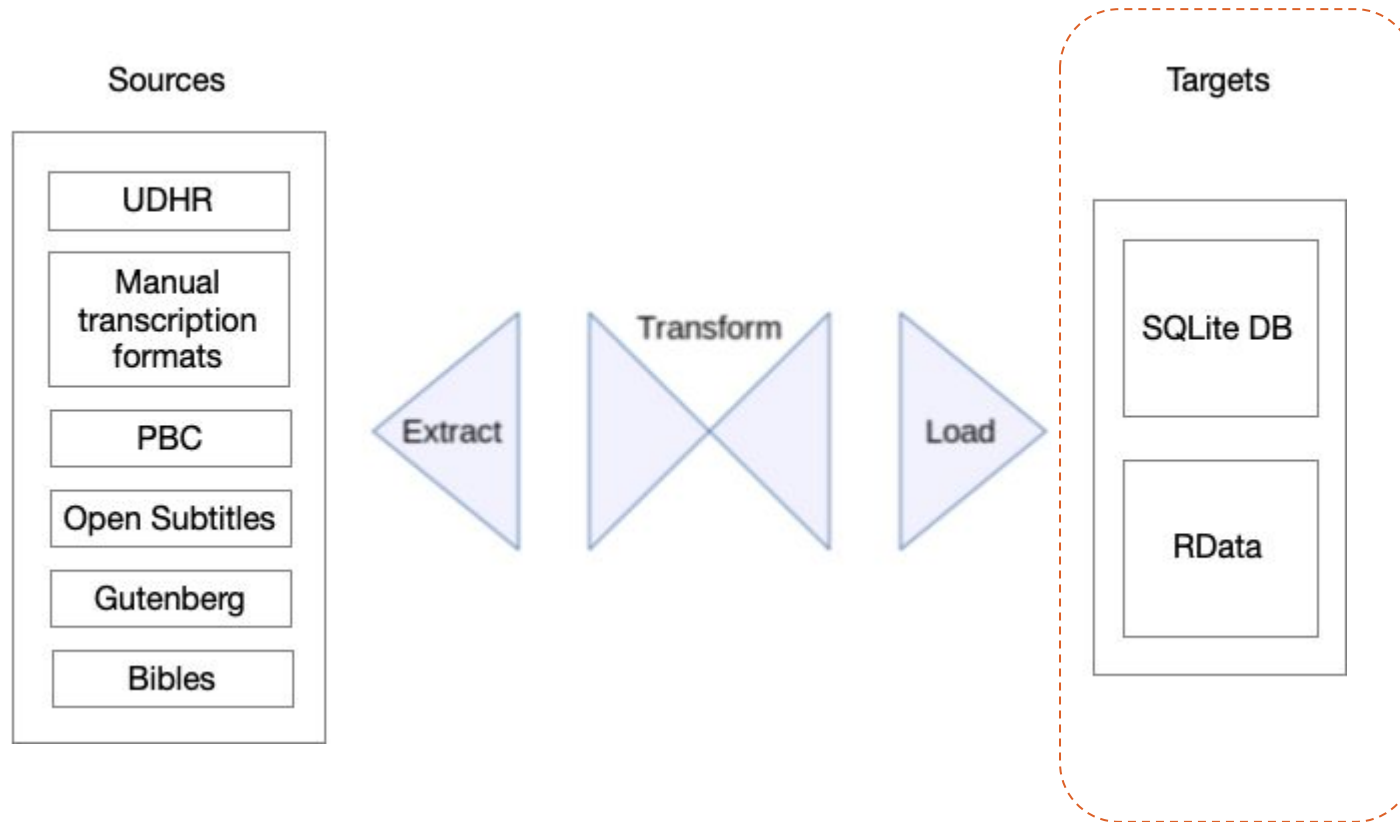


## Data availability



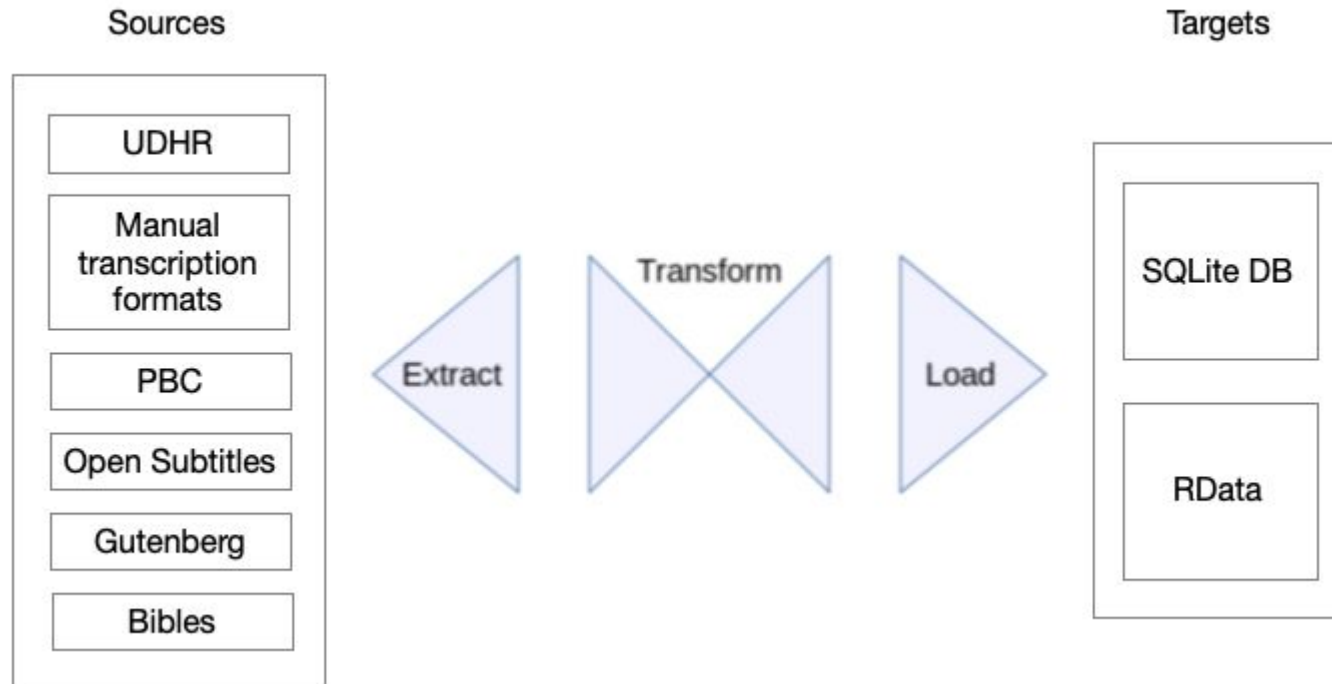


## Data availability





## Data availability

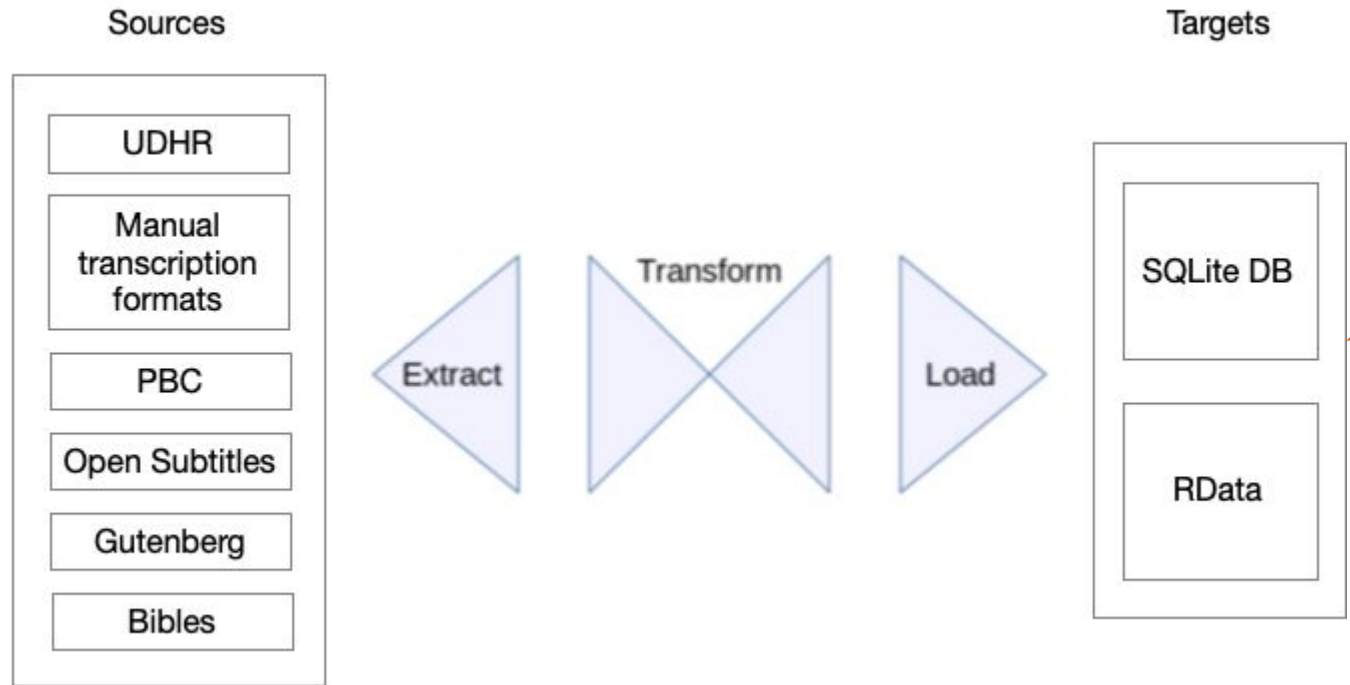


Source text files available  
directly in the repo

Generate these locally using  
the scripts in the repo



# Data availability



Source text files available directly in the repo

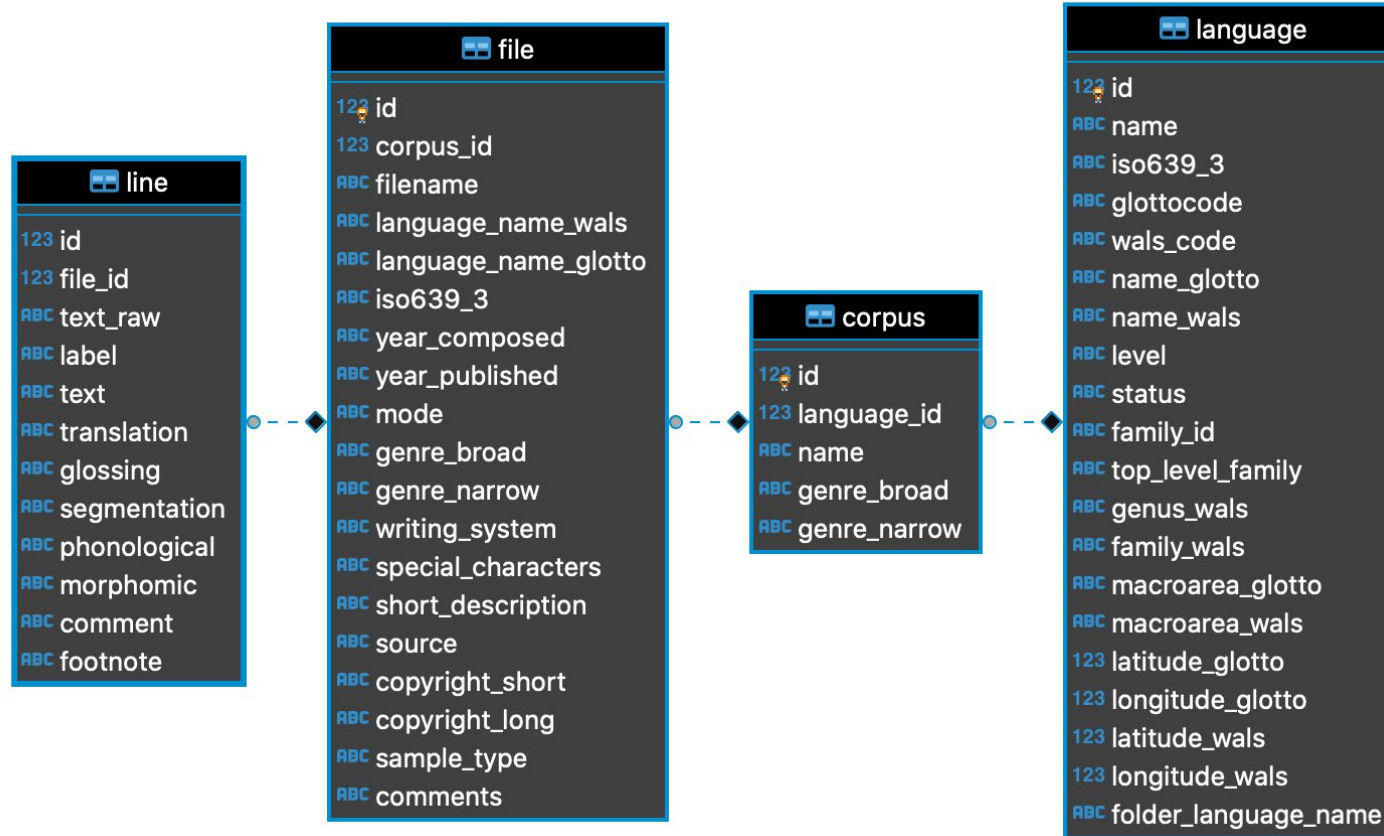
Generate these locally using the scripts in the repo

\*Database dumps are also available online





# Data availability



TeDDi's relational database schema



## Data availability

**Input corpora and the source code for processing them:**

[https://github.com/MorphDiv/TeDDi\\_sample](https://github.com/MorphDiv/TeDDi_sample)

Data exported into CLDF (Cross-Linguistic Data Formats)

[https://github.com/cldf-datasets/TeDDi\\_sample](https://github.com/cldf-datasets/TeDDi_sample)

CC BY-NC-SA 4.0 license\*



\*The particular copyright is given in the metadata header for each text

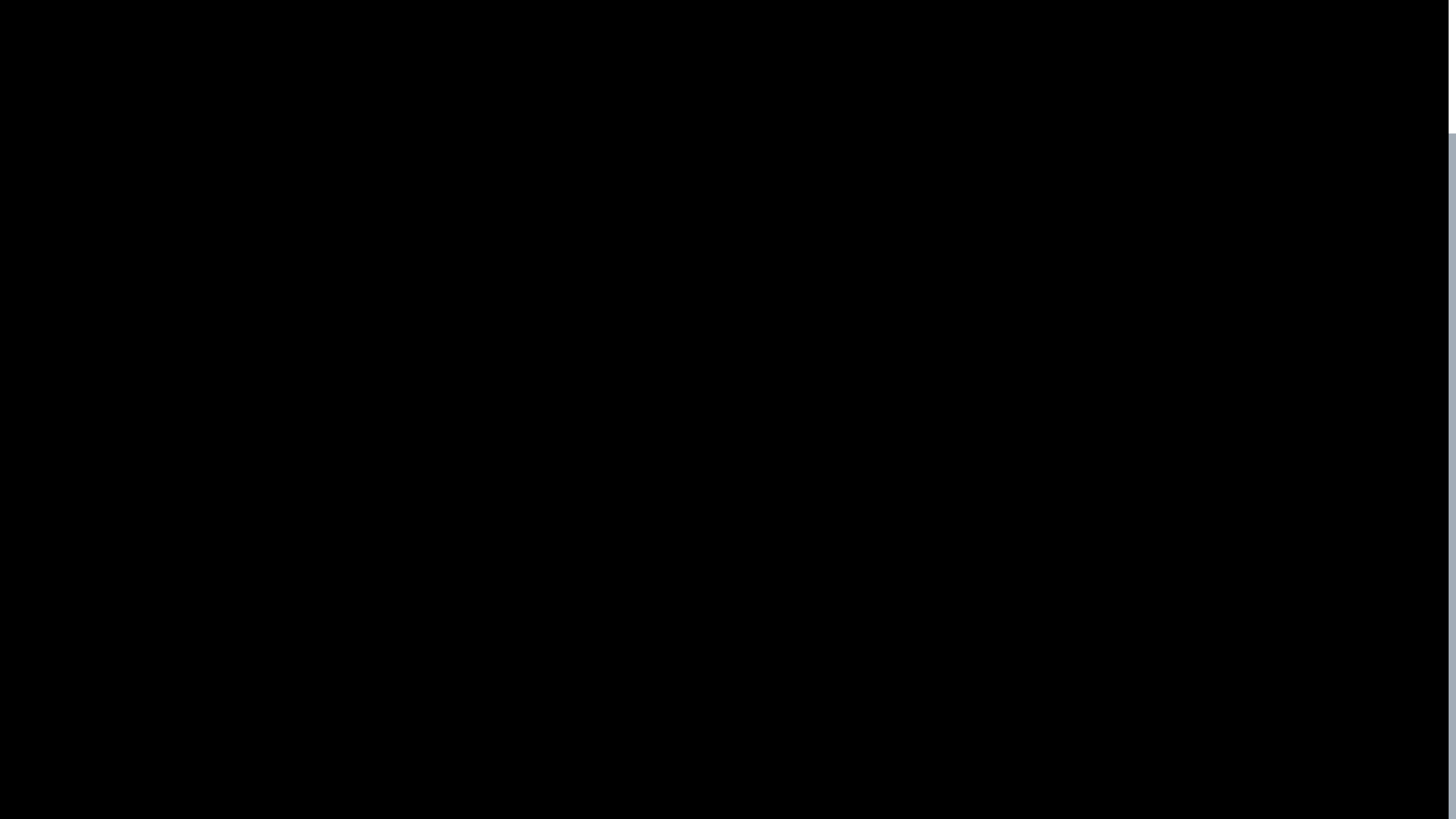
\* “takedown principle”, i.e., we can remove such material if contacted by people aggrieved by it (Seyfeddinipur et al., 2019, p. 554)



**University of  
Zurich** <sup>UZH</sup>

**URPP Language and Space**

**DEMO**



**Thank you!**

