# SwissSLi: the Multi-parallel Sign Language Corpus for Switzerland
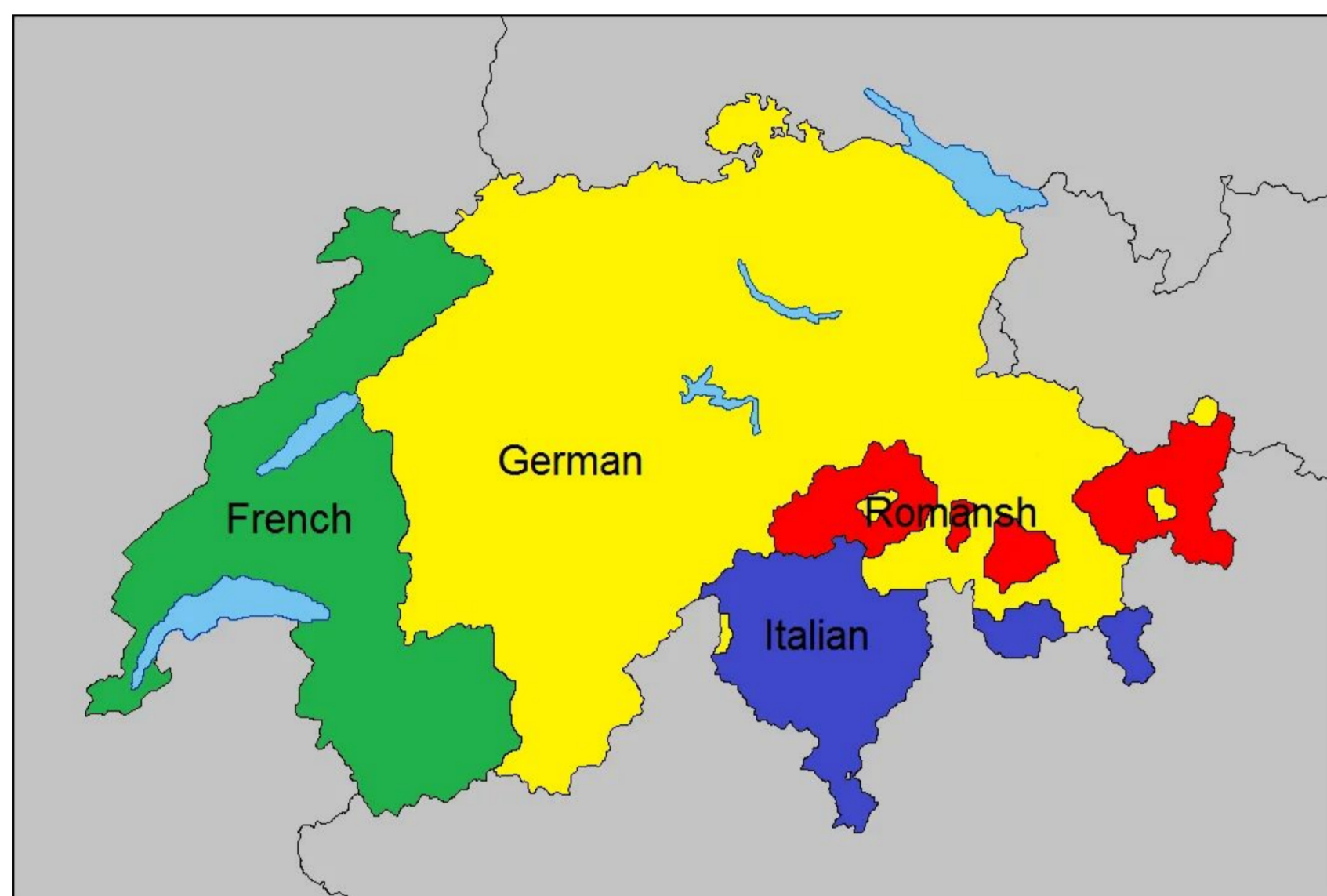
**LREC-COLING 2024**

**Zifan Jiang** (jiang@cl.uzh.ch), Anne Göhring, Amit Moryossef, Rico Sennrich, Sarah Ebling
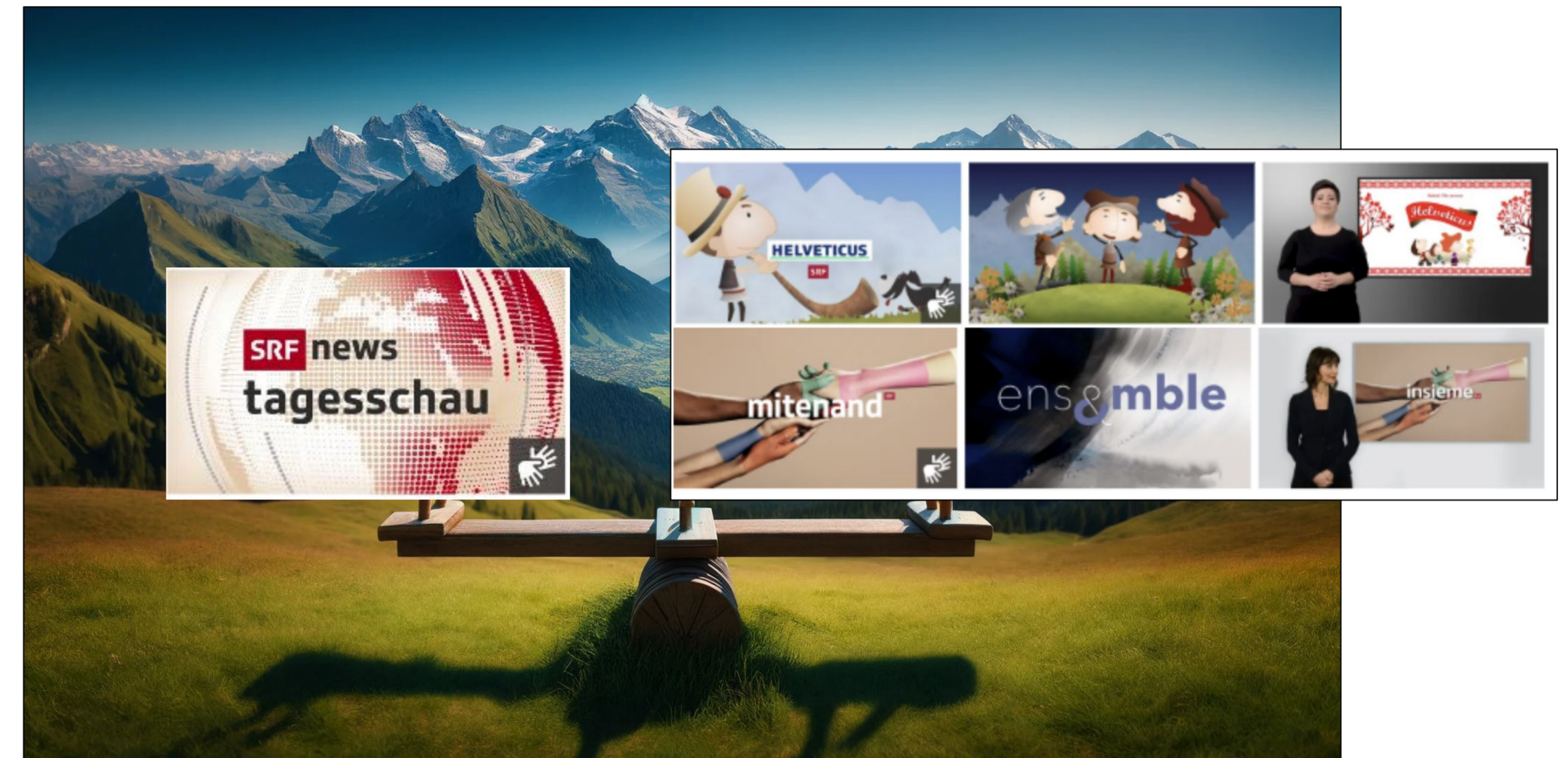Department of Computational Linguistics, University of Zurich

## Motivation

> Include sign languages in natural language processing (Yin et al., 2021), as they **are** fully natural languages.
> Limited data is a main challenge, as shown by *WMT-SLT* (Müller et al., 2022, 2023).
> The multilingual landscape in Switzerland offers the possibility for a **multi-parallel** sign language corpus.
> We focus on TV programs **translated by deaf signers offline**, in total ~30 hours in 3 sign and 3 spoken languages.

## Background & Discussion

### 1. Languages in Switzerland



### 2. Taxonomy of Swiss TV Programs



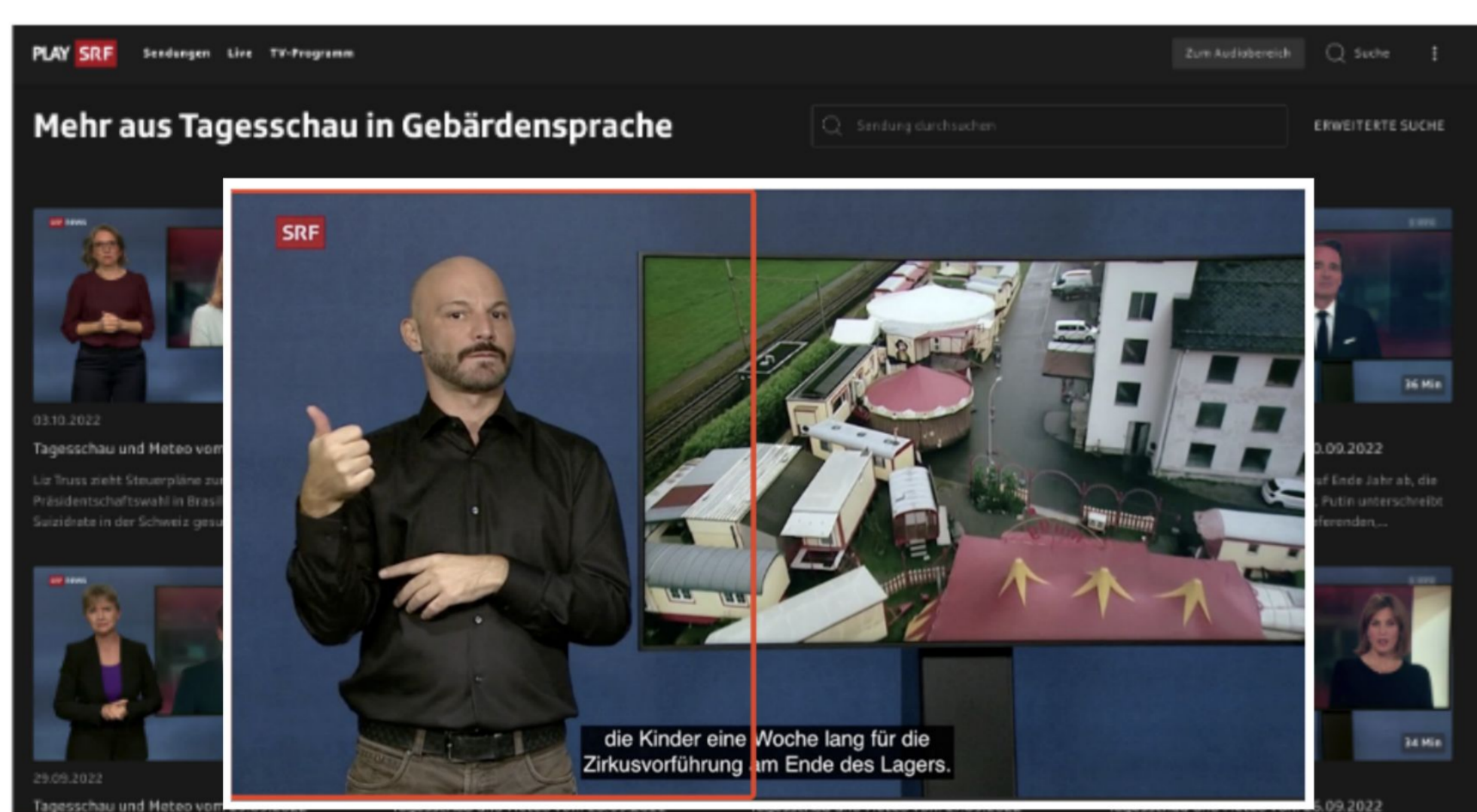### 3. Signing Mode and Data Quality

Sign language as source     TV interpretation     Offline translation



### 4. Data Quality, Quantity, and Licensing

Deaf community     Policy maker     Computer scientists



## Data & Processing

### 1. Raw Data from Program Websites



### 2. Video/Pose Processing

- Pose estimation
- Video segmentation

Figure taken from https://www.wmt-slt.com/data

Figure taken from Moryossef et al. (2023)



### 3. Subtitle Processing



Figure taken from the WMT-SLT findings paper

## Conclusion & Outlook

> The corpus is publicly available on the **SWISSUbase** data platform for research purposes under *CC BY-NC-SA 4.0*.
> Shared with the **informed consent** of the signers involved and a data-sharing agreement with the *Swiss Broadcasting Corporation*.
> A valuable asset for researchers in computer vision, natural language processing, and sign language linguistics.

**Flagship supported by**

Schweizerische Eidgenossenschaft
Confédération suisse
Confederazione Svizzera
Confederaziun svizra
Swiss Confederation

**Innosuisse – Swiss Innovation Agency**

**University of Zurich** UZH