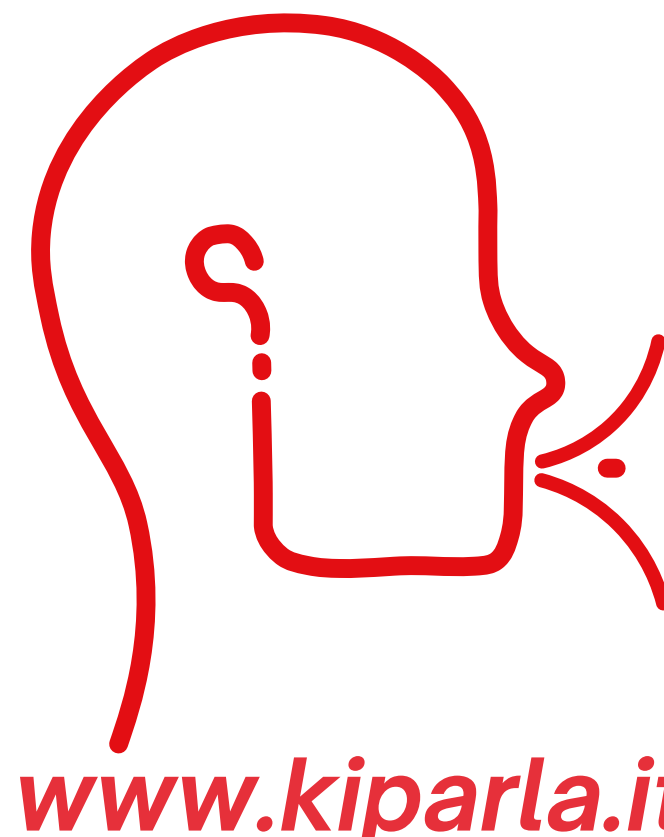# Introducing KIParla Forest:
## *seeds for a UD annotation of interactional syntax*

Ludovica Pannitto[1], Eleonora Zucchini[1], Silvia Ballarè[1], Cristina Bosco[2], Caterina Mauri[1], Manuela Sanguinetti[3]

[1]Alma Mater Studiorum - University of Bologna,[2]University of Turin,[3]University of Cagliari

SYNTAX FEST Ljubljana August 26-29 2025

EXPERIMENTAL LAB
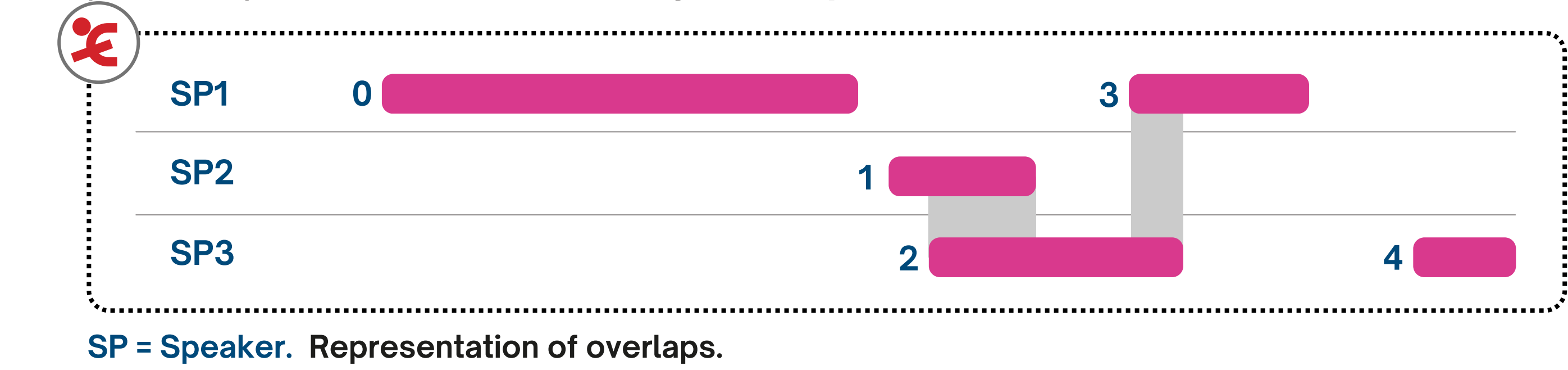lilec.lab@unibo.it

UniDive

www.kiparla.it

## From KIParla to *KIParla Forest*

The KIParla resource [1,2,3] is a corpus of spoken italian manually transcribed following Jeffersonian notation for Conversation Analysis (CA) [4]. All summed up, the KIParla counts ca. **228 hours** of recordings and approximately **2M transcribed tokens**.

```
0 bene (.) allora e::h (.) >non siete
tutti e due di< bologna quindi
1 [solo] io
2 [n~] (.) no s[olo lui]
3 [di bolo]gna,
4 napoli
```

```
0 well (.) then e::h (.) >you're not
both from< bologna then
1 [just] me
2 [n~] (.) no j[ust him]
3 [from bolo]gna,
4 napoli
```

The corpus is currently segmented into **transcription units (TUs)**, based on pseudo-prosodic hints. TUs may **overlap.**



SP = Speaker. Representation of overlaps.

## Speech-specific metadata

```
# sent_id
# text
# jefferson_text
# audio_url
```

At sentence level, metadata include the **sentence id** (aligned to conversation), original Jefferson transcription and url to audio (available for research purposes). Further metadata is retained in a separate *json* file available in the repository.

Each token retains (in MISC):
- **speaker ID**
- **boundary-related** features,
- **language variation** features (foreign languages + dialects)
- and **pseudonymization** information

```
AlignBegin=xxx(ms) and AlignEnd=xxx(ms)
UnitBoundary=Yes
Lang=(NO_ISO_CODE|iso-code)
Anonymized=Yes
```

- Information derived from CA annotation:
  - **intonation** pattern
  - **prolonged sounds**
  - **volume** and **pace**
  - **short pauses** and **prosodic links**

```
Intonation=(Rising|WeaklyRising|Descending)
Prolongation=Yes
PauseAfter=Yes
Volume=(High|Low)
Pace=(Fast|Slow)
ProsodicLink=Yes
```

- information about **overlapping speech** (reference to overlapping tokens)

## Morphosyntactic information

**Tokenization, lemmatization**

Pauses are removed and transformed into a feature in MISC.

**Multi-word** tokens keep CA features.

Interrupted words are lemmatized as their complete version when context informative enough + specific feature

| text | form | upos | lemma | MISC |
|---|---|---|---|---|
| casa (.) | casa | NOUN | casa | PauseAfter= Yes |

**PoSTagging**

*Conservative approach:* tag main category of each word

```
basta lit. 3sg of bastare, to
suffice tagged as VERB - en. 'stop'

tipo lit. type tagged as NOUN -
en. 'for example', 'like'
```
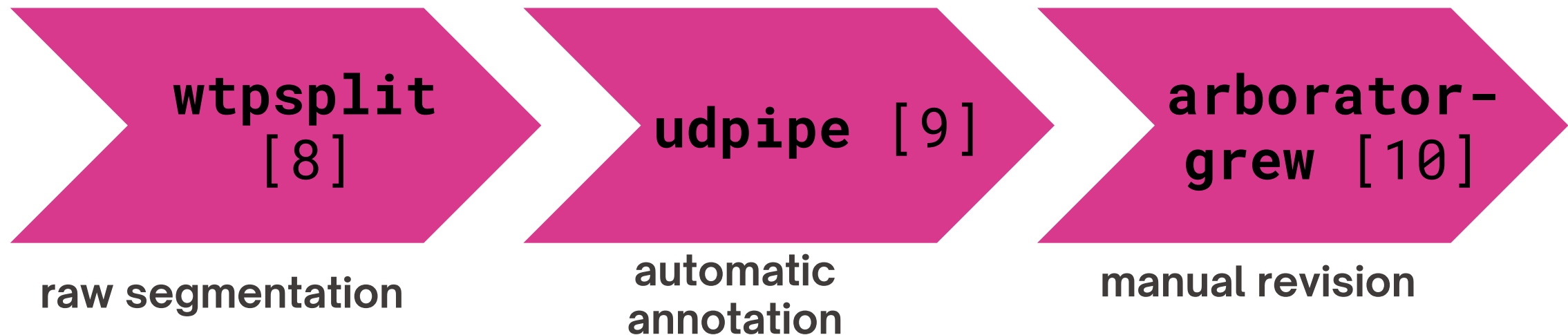
| form | upos | lemma | MISC | *gloss* |
|---|---|---|---|---|
| c' | PRON | ci | | there |
| era | VERB | essere | | was |
| so~ | ADV | solo | Interrupted=Yes | on~ |
| c' | PRON | ci | | there |
| era | VERB | essere | | was |
| solo | ADV | solo | | only |
| casa | NOUN | casa | | house |
| mia | ADJ | mio | | my |

Inter-annotator agreement
- Cohen's κ > 0.87
- most disagreement on CCONJ and ADV

**Pipeline**

wtpsplit [8] → udpipe [9] → arborator-grew [10]

raw segmentation → automatic annotation → manual revision

## SpLAn-UD

Increased attention to the syntactic annotation of spoken varieties within the Universal Dependencies framework is attested by the fact that the number of treebanks including or completely dedicated to spoken language is on the rise.

Treebank curators took different directions in the creation of their resources, which could impact on derived measures or performance on downstream tasks [5,6]
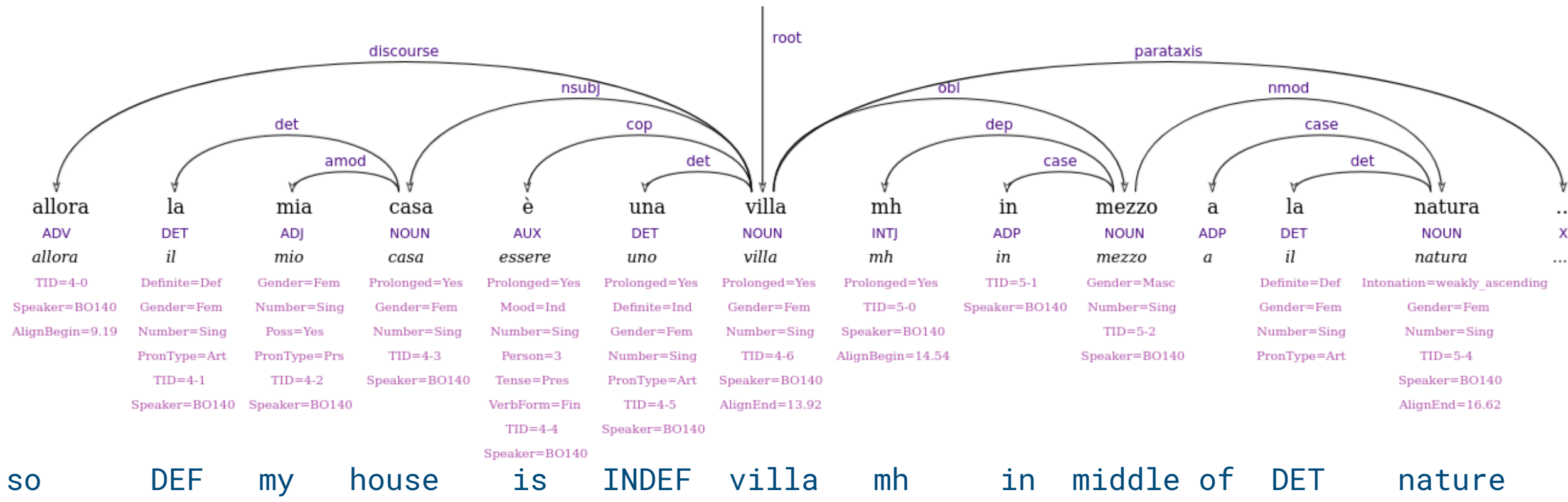
*Join!*

Task 1.5

UniDive

## Data sample

First release (Nov 2025) aims at ~22K tokens selected based on **type of interaction** (free turn-taking, partially free turn-taking, rigid turn-taking, close to no interaction).

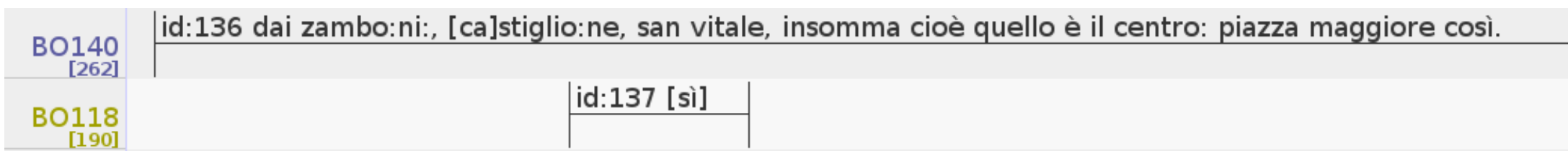| | | conversation id | info |
|---|---|---|---|
| free | table conversations + free interactions | BOA3017 | 4551 tokens, 4 participants |
| partial | semi-structured interviews | PBB004 | 5898 tokens, 3 participants |
| | | BOD2018 | 4634 tokens, 2 participants |
| close to none | lessons | TOD1005bis | 6788 tokens, 1 participant |

## Segmentation and syntax

TU boundaries are not reliable as sentence boundaries: we identify maximal-units **based on syntactic dependencies**, in order to focus on **interactional syntax** and **cross-speaker syntactic affordances**
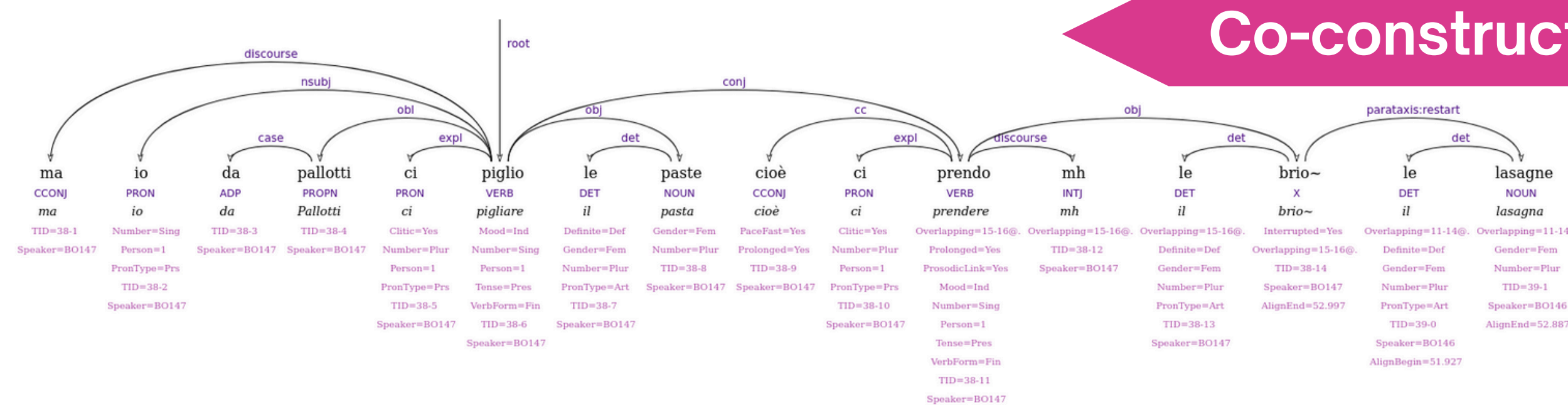


**Special cases:**
- **connectives with discourse functions:**
  - unit boundary is postulated if no relation with previous portion exists
- **feedback phenomena:**
  - no unit boundary is postulated if speech flow is uninterrupted



Co-construction

## REFERENCES

[1] https://kiparla.it
[2] C. Mauri, et al., Kiparla corpus: a new resource for spoken italian, in *Proceedings of the 6th Italian Conference on Computational Linguistics CLiC-it* (2019).
[3] S. Ballarè, C. Mauri, et al.. La creazione del corpus KIParla: criteri metodologici e prospettive future. *Rivista italiana di dialettologia* (2020)
[4] G. Jefferson, et al., Glossary of transcript symbols with an introduction, *Conversation analysis* (2004)
[5] K. Dobrovoljc. Spoken language treebanks in universal dependencies: An overview. In *Proceedings of the Thirteenth Language Resources and Evaluation Conference* (2022)
[6] S. Kahane, et al. Annotation guidelines of UD and SUD treebanks for spoken corpora: A proposal. In *Proceedings of the 20th International Workshop on Treebanks and Linguistic Theories* (2021)
[7] C. Bosco et al. The Evalita 2014 Dependency Parsing task. in *Proceedings of the First Italian Conference on Computational Linguistics CLiC-it* (2014)
[8] B. Minixhofer et al. Where's the point? self-supervised multilingual punctuation-agnostic sentence segmentation. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics* (2023)
[9] M. Straka. 2024. Universal dependencies 2.15 models for UDPipe 2 (2024-11-21). *LINDAT/CLARIAH-CZ digital library at the Institute of Formal and Applied Linguistics* (2024).
[10] G. Guibon et al., When collaborative treebank curation meets graph grammars. In *Proceedings of the Twelfth Language Resources and Evaluation Conference* (2020)