

Document Overview

- This is a documentation overview of how to use Batchalign2 to do syntactic and morphological parsing with examples in Italian using the [VoxCommunis spoken Italian corpus](#)

Batchalign2 Overview

- Batchalign2 is a command-line pipeline developed under the TalkBank project for automatic morphological tagging and Universal Dependencies (UD) syntactic parsing across multiple languages, including Italian.
- It takes linguistic transcripts in the CHAT (.cha) format and produces two key annotation layers:
 - %mor: morphological analysis (lemmas, parts of speech, and inflectional features)
 - %gra: syntactic analysis following the Universal Dependencies (UD) standard.
- Batchalign2 integrates the **Stanza** NLP models and the **UD framework** into the TalkBank/CLAN ecosystem, enabling consistent morphosyntactic annotation of multilingual corpora.
- It automates tokenization, morphological tagging, lemmatization, and dependency parsing, storing all results back in CHAT format—making the output directly usable with CLAN tools for further linguistic or developmental analysis.
- Functions and Commands:
 - asr: ASR!
 - morphosyntax: PoS and dependency analysis
 - fa: Forced Alignment (requires utterance-level timings already)

Universal Dependencies Overview

- **Universal Dependencies (UD)** is an international framework for the consistent annotation of grammar—covering parts of speech, morphological features, and syntactic dependencies—across human languages.
- UD is an open, collaborative project involving over 600 contributors, maintaining more than 200 treebanks in 150+ languages.
- Its goal is to create a cross-linguistically consistent representation of grammatical relations, allowing comparable syntactic analysis across languages like Italian, Mandarin, English, and many others.

What we need to do in terms of using CLAN and Batchalign2 to do morphological and syntactic analysis and parsing:

1. **Download CLAN (MacBook or Windows, or for a more detailed overview of CLAN, please see → [here](#))**
<https://dali.talkbank.org/clan/>

2. Download Batchalign/Environment Setup (Open MacBook Terminal and paste the following)

```
# Install UV (a lightweight Python environment manager)
curl -LsSf https://astral.sh/uv/install.sh | sh

# Install Batchalign2 for Python 3.11
UV_PYTHON=3.11 uv tool install batchalign

# Ensure ~/.local/bin is in your PATH
echo 'export PATH="$HOME/.local/bin:$PATH"' >> ~/.zshrc
source ~/.zshrc

# Check installation
batchalign --help
```

3. Data Preprocessing (R Studio) → also see [txt-to-cha.R](#)

Aim data cleaning: Txt file → cha file

- A valid .cha file must follow CHAT conventions: every utterance on its own line and starting with a speaker code (e.g., *ADU:, *CHI:, *MOT:).

Tips:

- A Cha file can be automatically converted to a txt file and vice versa (just **rename** the a cha file to a txt file! **Example:** Sample.cha → Sample.txt)
- **@Languages:** must include the language code (e.g., ita, zho, eng, spa).
 - The Italian model will be loaded automatically from Stanza/UD based on **@Languages: ita**.
- Sentences must end with a space and punctuation (., !, or ?).
- Delete empty or malformed tokens like <>, which cause lexer errors.

Things we need to do for the VoxCommunis dataset:

1. Add speaker tier at the beginning of each sentence (e.g., *ADU:)
2. Remove all other symbols (e.g., @, *, % or <>:)
3. Ensure a space before punctuations (e.g., . ! or ?)
4. Add a final period if missing
5. Make sure to add the following at the beginning of the cha file
 - @UTF8
 - @Begin
 - @Languages: ita
 - @Participants: ADU Adult
 - @ID: ita|troncamento|ADU||||Adult||

6. Make sure to add the following at the end of the cha file
- @End

Successful Example:

```
@UTF8 → Must
@Begin → Must
@Languages: ita → Must
@Participants: ADU Adult → Must
@ID: ita|troncamento|ADU||||Adult|| → Must
@Media: 10192025, audio → optional
@Date: 19-OCT-2025 → optional
*ADU: Il libro ha suscitato molte polemiche a causa dei suoi contenuti . → Must
@End → Must
```

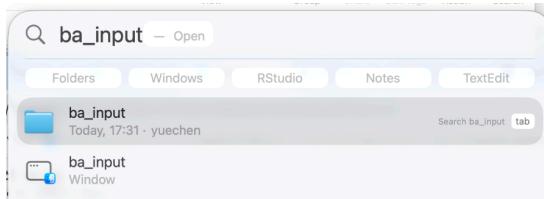


```

1 @UTF8
2 @Begin
3 @Languages: ita
4 @Participants: ADU Adult
5 @ID: ita|troncamento|ADU||||Adult||
6 @Media: 10192025, audio
7 @Date: 19-OCT-2025
8 *ADU: Il libro ha suscitato molte polemiche a causa dei suoi contenuti .
9 @End
10

```

Sample Input file (where you should put it? Search ba_input!):



4. Parsing: Commands used for Batchalign2 (MacBook Terminal)

```
# check input file
sed -n '1,25p' ~/ba_input/troncamento.cha
```

```
# run batch aligner
batchalign morphotag ~/ba_input ~/ba_output

# check output file %mor / %gra
grep -n "^\%mor:\|^%gra:" ~/ba_output/troncamento.cha | head
```

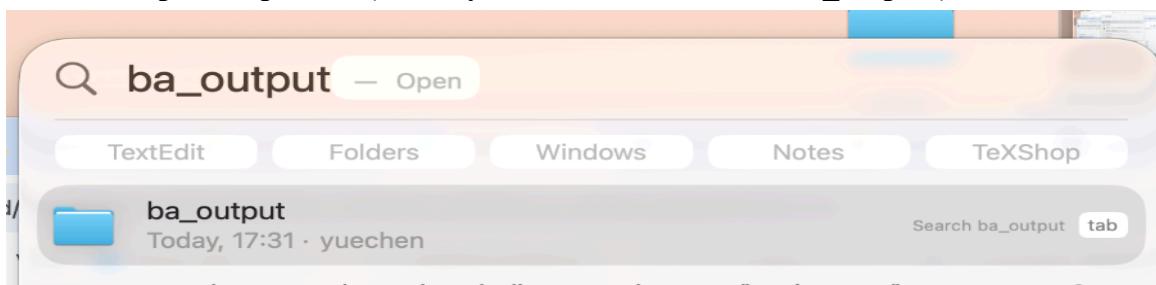
What you will see if it's successfully running:

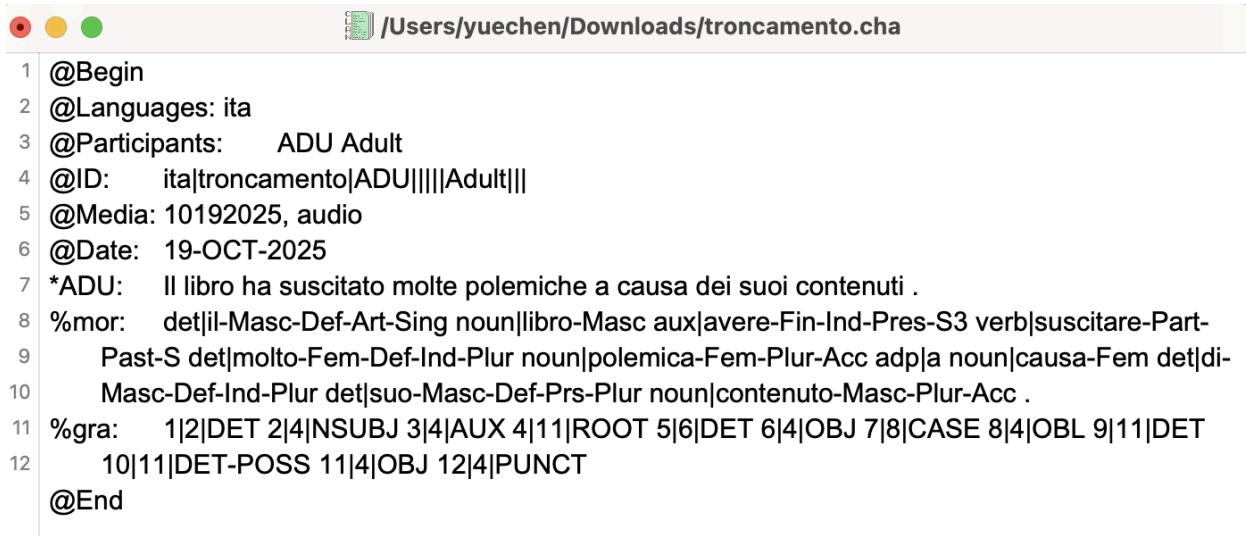
● ● ● ⚡ yuechen — python -m batchalign morphotag ~/ba_input ~/ba_output — 8...

```
/Users/yuechen/.local/share/uv/tools/batchalign/lib/python3.11/site-packages/pr
atio/utilities/utils.py:9: UserWarning: pkg_resources is deprecated as an API. S
ee https://setuptools.pypa.io/en/latest/pkg_resources.html. The pkg_resources pa
ckage is slated for removal as early as 2025-11-30. Refrain from using this pack
age or pin to SetupTools<81.
    from pkg_resources import resource_filename
/Users/yuechen/.local/share/uv/tools/batchalign/lib/python3.11/site-packages/pya
nnote/audio/core/io.py:212: UserWarning: torchaudio._backend.list_audio_backends
    has been deprecated. This deprecation is part of a large refactoring effort to
    transition TorchAudio into a maintenance phase. The decoding and encoding capabi
    lities of PyTorch for both audio and video are being consolidated into TorchCode
    c. Please see https://github.com/pytorch/audio/issues/3902 for more information.
    It will be removed from the 2.9 release.
        torchaudio.list_audio_backends()

Mode: morphotag; got 1 transcript to process from /Users/yuechen/ba_input:
Downloading https://raw.githubusercontent.com/stanfordnlp/stanza-resources/main/
Downloading
https://raw.githubusercontent.com/stanfordnlp/stanza-resources/main/resources_1.
11.0.json: 436kB [00:00, 242MB/s]
: troncamento.cha ----- 3% 0:32:56 Running: Morpho-Syntax
```

Sample Output file (where you can find it? Search ba_output!):





```

1 @Begin
2 @Languages: ita
3 @Participants: ADU Adult
4 @ID: ita|troncamento|ADU|||||Adult|||
5 @Media: 10192025, audio
6 @Date: 19-OCT-2025
7 *ADU: Il libro ha suscitato molte polemiche a causa dei suoi contenuti .
8 %mor: det|il-Masc-Def-Art-Sing noun|libro-Masc aux|avere-Fin-Ind-Pres-S3 verb|suscitare-Part-
9 Past-S det|molto-Fem-Def-Ind-Plur noun|polemica-Fem-Plur-Acc adp|a noun|causa-Fem det|di-
10 Masc-Def-Ind-Plur det|suo-Masc-Def-Prs-Plur noun|contenuto-Masc-Plur-Acc .
11 %gra: 1|2|DET 2|4|NSUBJ 3|4|AUX 4|11|ROOT 5|6|DET 6|4|OBJ 7|8|CASE 8|4|OBL 9|11|DET
12 10|11|DET-POSS 11|4|OBJ 12|4|PUNCT
@End

```

5. Understanding the Output (sentence level) → check [it_vxc_spkr17_with_parsing.tsv](#):

Sample sentence: Il libro ha suscitato molte polemiche a causa dei suoi contenuti.

English Translation: The book has sparked much controversy because of its contents.

Morphological layer (%mor) (also see the column named [morphonological_analysis](#) in the [it_vxc_spkr17_with_parsing.tsv](#) file):

```
%mor: det|il-Masc-Def-Art-Sing noun|libro-Masc aux|avere-Fin-Ind-Pres-S3
verb|suscitare-Part-Past-S det|molto-Fem-Def-Ind-Plur noun|polemica-Fem-Plur-Acc
adp|a noun|causa-Fem det|di-Masc-Def-Ind-Plur det|suo-Masc-Def-Prs-Plur
noun|contenuto-Masc-Plur-Acc .
```

What does the Morphological layer (%mor) tire do?

- Breaks the sentence down word-by-word and encodes part of speech + grammatical features (gender, number, tense, etc.).
 - Lemma (dictionary form): e.g. libro, suscitare
 - Part of speech: noun, verb, det, adp, etc.
 - Morphosyntactic features:
 - Gender: Masc/Fem
 - Number: Sing/Plur
 - Tense: Pres, Past, Fut
 - Mood: Ind (indicative), Subj (subjunctive)
 - Person: S1, S2, S3 (subject 1st, 2nd, 3rd person)
 - Case: Nom (nominative), Acc (accusative), etc.

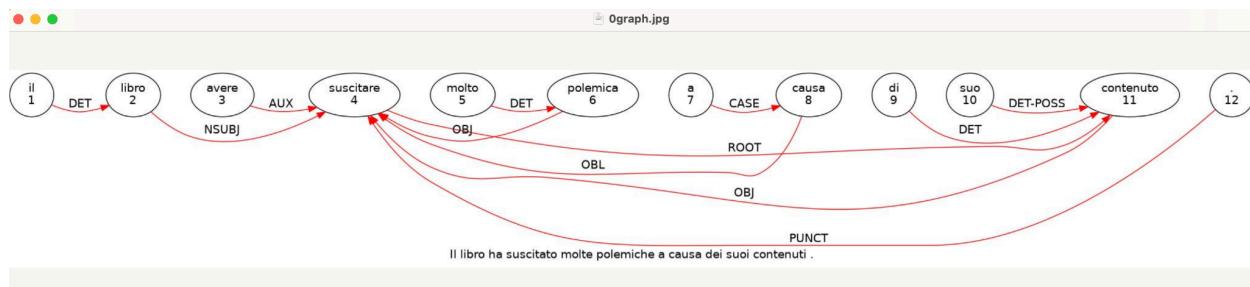
Syntactic layer (%gra) (also see the column named **syntactic_parsing** in the `it_vxc_spkr17_with_parsing.tsv` file):

```
%gra: 1|2|DET 2|4|NSUBJ 3|4|AUX 4|11|ROOT 5|6|DET 6|4|OBJ 7|8|CASE 8|4|OBL  
9|11|DET 10|11|DET-POSS 11|4|OBJ 12|4|PUNCT
```

What does the Syntactic layer (%gra) tire do?

- These lines show token IDs (word order, such as word 1, word 2, word 3), dependency heads, and UD relations (e.g., NSUBJ = subject, OBL = oblique).

Example of the Universal Dependencies (UD) dependency tree:



Overall, what this shows:

- Each oval = a token (word), labeled with
 - its lemma (dictionary form), and
 - its index (1–12 here).
- Each red arrow = a dependency relation between two words:
 - The arrow points from the head (governing word) to its dependent (modified word).
 - The label on the arrow (e.g., NSUBJ, AUX, OBJ, OBL, etc.) indicates the syntactic relation type.

Final sample [TSV file output](#):

it_vxc_spkr17_with_parsing

	morphological_analysis	syntactic_parsing
1	det-Masc-Def-Art-Sing noun[libro-Masc aux[javere-Fin-Ind-Pres-3 verb]suscitare-Part-Past-S det]Masc-Def-Ind-Pur noun[polimica-Fem-Plur Acc ad]pa noun[causa-Fem det]di-Masc-Def-Ind-Pur det[suo-Masc-Def-Ind-Pur noun[ad]dilirio-Def-Jnf-Fem-Def-Art-Sing noun[sede-Fem ad]episcopale-S1 auxjessere-Fin-Ind-Pres-33 auxjessere-Part-Past-S adv]immediatamente adj[sgetto]-S1 det]Masc-Def-Ind-Pur noun[chiostro-Masc-Plur cmjcm noun[ospedale-Masc-Plur ccnje noun[chiesa-Fem-Plur .	1 DET 2 DET 2 NSUBJ 3 4 AUX 4 1 ROOT 5 6 DET 6 4 OBJ 7 8 CASE 8 4 OBJ 9 1 DET 1 2 CASE 2 3 ADVMOD 3 4 DET 4 9 NSUBJ 5 4 AMOD 6 9 AUX 7 8 COP 8 9 ADVMOD 9 auxjessere-Fin-Ind-Past-33 det]Masc-Def-Art-Sing noun[fondatore-Masc ad]pi det]Masc-Def-Ind-Pur noun[chiostro-Masc-Plur cmjcm noun[ospedale-Masc-Plur ccnje noun[chiesa-Fem-Plur .
2	ad]pi det]Masc-Def-Ind-Pur noun[voto-Masc ad]sesso-S1 ?	1 DET 2 3 ROOT 3 2 AMOD 4 2 PUNCT
3	ad]pi det]Masc-Def-Ind-Pur noun[janno-Masc-Plur cmjcm pron[e]gli-Prs-33 adjpi verb]tornare-Int-S adjpi prop]n[India ad]pi verb[accogliere-Inf-S adj]abro-P1 noun[insegnamento	1 DET 2 1 FLAT-NAME 3 1 PUNCT
4	prop]n[Salvation prop]n[ue	1 DET 2 2 FLAT-NAME 3 1 PUNCT
5	ad]pi det]questo-Masc-Def-Sing noun[modo-Masc cmjcm prop]Decio verb]tottenham-Fin-Ind-Past-33 det]Masc-Def-Art-Sing noun[potere-Masc-Acc adj]imperiale-S1 .	1 DET 2 3 CASE 2 3 DET 3 6 OBL 4 3 PUNCT 5 6 NSUBJ 6 9 ROOT 7 8 MARK 8 6 XCOMP 9
6	prop]n[Sparta prop]Novara verb]acquisire-Fin-Ind-Pres-33 det]Masc-Def-Art-Sing noun[modo-Masc-Acc ad]pi verb]giocare-Inf-3 adjpi ad]rimo-S1 noun[categoria-Fem .	1 NSUBJ 2 1 FLAT-NAME 3 1 ROOT 4 5 DET 5 9 NSUBJ 6 5 AMOD 7 8 MARK 8 5 ADV 9 advjessere-Fin-Ind-Past-33 det]Masc-Def-Art-Sing noun[modo-Masc-Acc ad]pi verb]continuire-Inf-3 adjpi verb]vivere-Inf-5 adjpi descritta-Fem-Def-Art-Sing noun[canzone-Fem-1 .
7	1 DET 2 2 FLAT-NAME 3 1 PUNCT	
8	1 DET 2 3 CASE 2 3 DET 3 9 OBL 4 3 PUNCT 5 6 NSUBJ 6 9 ROOT 7 8 DET 8 6 OBJ 9 8 AMOI 10 DET 2 4 NSUBJ 3 2 PUNCT 4 8 NSUBJ 5 6 CONJ 7 8 AUX 8 1 ROOT 9 10 M 11 DET 2 5 NSUBJ 3 6 PUNCT 4 9 NSUBJ 5 7 CONJ 6 8 AUX 7 9 MARK 8 10 ADV 9 advjessere-Masc cmjcm prop]Kygo conje prop]Shear aux[javere-Fin-Ind-Pres-3 verb]prop]pome-Part-Past-S adjpi verb]continuire-Inf-3 adjpi verb]vivere-Inf-5 adjpi descritta-Fem-Def-Art-Sing noun[canzone-Fem-1 .	1 DET 2 2 FLAT-NAME 3 1 PUNCT
12	1 DET 2 3 CASE 2 3 DET 3 9 OBL 4 3 PUNCT 5 6 NSUBJ 6 4 CONJ 7 8 AUX 8 1 ROOT 9 10 M 11 DET 2 5 NSUBJ 3 6 PUNCT 4 9 NSUBJ 5 7 CONJ 6 8 AUX 7 9 MARK 8 10 ADV 9 advjessere-Masc cmjcm prop]Stephen prop]brun .	1 DET 2 1 FLAT-NAME 3 1 PUNCT
13	1 DET 2 3 CASE 2 3 DET 3 9 OBL 4 3 PUNCT 5 6 NSUBJ 6 4 CONJ 7 8 AUX 8 1 ROOT 9 10 M 11 DET 2 5 NSUBJ 3 6 PUNCT 4 9 NSUBJ 5 7 CONJ 6 8 AUX 7 9 MARK 8 10 ADV 9 advjessere-Masc cmjcm prop]Ursino .	1 DET 2 1 FLAT-NAME 3 1 PUNCT
14	1 DET 2 3 CASE 2 3 DET 3 9 OBL 4 3 PUNCT 5 6 NSUBJ 6 4 CONJ 7 8 AUX 8 1 ROOT 9 10 M 15 DET 2 4 NSUBJ 3 2 PUNCT 4 8 NSUBJ 5 6 CONJ 6 7 AUX 7 8 MARK 8 9 ADV 16 DET 2 5 NSUBJ 3 6 PUNCT 4 9 NSUBJ 5 7 CONJ 6 8 AUX 7 9 MARK 8 10 ADV 17 DET 2 6 NSUBJ 3 7 PUNCT 4 10 NSUBJ 5 8 CONJ 6 9 AUX 7 10 MARK 8 11 ADV 18 DET 2 7 NSUBJ 3 8 PUNCT 4 11 NSUBJ 5 9 CONJ 6 10 AUX 7 11 MARK 8 12 ADV 19 DET 2 8 NSUBJ 3 9 PUNCT 4 12 NSUBJ 5 10 CONJ 6 11 AUX 7 12 MARK 8 13 ADV 20 DET 2 9 NSUBJ 3 10 PUNCT 4 13 NSUBJ 5 11 CONJ 6 12 AUX 7 13 MARK 8 14 ADV 21 DET 2 10 NSUBJ 3 11 PUNCT 4 14 NSUBJ 5 12 CONJ 6 13 AUX 7 14 MARK 8 15 ADV 22 DET 2 11 NSUBJ 3 12 PUNCT 4 15 NSUBJ 5 13 CONJ 6 14 AUX 7 15 MARK 8 16 ADV 23 DET 2 12 NSUBJ 3 13 PUNCT 4 16 NSUBJ 5 14 CONJ 6 15 AUX 7 16 MARK 8 17 ADV 24 DET 2 13 NSUBJ 3 14 PUNCT 4 17 NSUBJ 5 15 CONJ 6 16 AUX 7 17 MARK 8 18 ADV 25 DET 2 14 NSUBJ 3 15 PUNCT 4 18 NSUBJ 5 16 CONJ 6 17 AUX 7 18 MARK 8 19 ADV 26 DET 2 15 NSUBJ 3 16 PUNCT 4 19 NSUBJ 5 17 CONJ 6 18 AUX 7 19 MARK 8 20 ADV 27 DET 2 16 NSUBJ 3 17 PUNCT 4 20 NSUBJ 5 18 CONJ 6 19 AUX 7 20 MARK 8 21 ADV 28 DET 2 17 NSUBJ 3 18 PUNCT 4 21 NSUBJ 5 19 CONJ 6 20 AUX 7 21 MARK 8 22 ADV 29 DET 2 18 NSUBJ 3 19 PUNCT 4 22 NSUBJ 5 20 CONJ 6 21 AUX 7 22 MARK 8 23 ADV 30 DET 2 19 NSUBJ 3 20 PUNCT 4 23 NSUBJ 5 21 CONJ 6 22 AUX 7 23 MARK 8 24 ADV 31 DET 2 20 NSUBJ 3 21 PUNCT 4 24 NSUBJ 5 22 CONJ 6 23 AUX 7 24 MARK 8 25 ADV 32 DET 2 21 NSUBJ 3 22 PUNCT 4 25 NSUBJ 5 23 CONJ 6 24 AUX 7 25 MARK 8 26 ADV 33 DET 2 22 NSUBJ 3 23 PUNCT 4 26 NSUBJ 5 24 CONJ 6 25 AUX 7 26 MARK 8 27 ADV 34 DET 2 23 NSUBJ 3 24 PUNCT 4 27 NSUBJ 5 25 CONJ 6 26 AUX 7 27 MARK 8 28 ADV 35 DET 2 24 NSUBJ 3 25 PUNCT 4 28 NSUBJ 5 26 CONJ 6 27 AUX 7 28 MARK 8 29 ADV 36 DET 2 25 NSUBJ 3 26 PUNCT 4 29 NSUBJ 5 27 CONJ 6 28 AUX 7 29 MARK 8 30 ADV 37 DET 2 26 NSUBJ 3 27 PUNCT 4 30 NSUBJ 5 28 CONJ 6 29 AUX 7 30 MARK 8 31 ADV 38 DET 2 27 NSUBJ 3 28 PUNCT 4 31 NSUBJ 5 29 CONJ 6 30 AUX 7 31 MARK 8 32 ADV 39 DET 2 28 NSUBJ 3 29 PUNCT 4 32 NSUBJ 5 30 CONJ 6 31 AUX 7 32 MARK 8 33 ADV 40 DET 2 29 NSUBJ 3 30 PUNCT 4 33 NSUBJ 5 31 CONJ 6 32 AUX 7 33 MARK 8 34 ADV 41 DET 2 30 NSUBJ 3 31 PUNCT 4 34 NSUBJ 5 32 CONJ 6 33 AUX 7 34 MARK 8 35 ADV 42 DET 2 31 NSUBJ 3 32 PUNCT 4 35 NSUBJ 5 33 CONJ 6 34 AUX 7 35 MARK 8 36 ADV 43 DET 2 32 NSUBJ 3 33 PUNCT 4 36 NSUBJ 5 34 CONJ 6 35 AUX 7 36 MARK 8 37 ADV 44 DET 2 33 NSUBJ 3 34 PUNCT 4 37 NSUBJ 5 35 CONJ 6 36 AUX 7 37 MARK 8 38 ADV 45 DET 2 34 NSUBJ 3 35 PUNCT 4 38 NSUBJ 5 36 CONJ 6 37 AUX 7 38 MARK 8 39 ADV 46 DET 2 35 NSUBJ 3 36 PUNCT 4 39 NSUBJ 5 37 CONJ 6 38 AUX 7 39 MARK 8 40 ADV 47 DET 2 36 NSUBJ 3 37 PUNCT 4 40 NSUBJ 5 38 CONJ 6 39 AUX 7 40 MARK 8 41 ADV 48 DET 2 37 NSUBJ 3 38 PUNCT 4 41 NSUBJ 5 39 CONJ 6 40 AUX 7 41 MARK 8 42 ADV 49 DET 2 38 NSUBJ 3 39 PUNCT 4 42 NSUBJ 5 40 CONJ 6 41 AUX 7 42 MARK 8 43 ADV	1 DET 2 1 FLAT-NAME 3 1 PUNCT

R code of how to merge the CHA file to the TSV file:

- See [Cha-to-TSV.R](#)

6. Understanding the Output (token level) → [check it_vxc_spkr17_tokens_level.tsv](#)

In addition to the sentence-level TSV file (**it_vxc_spkr17_with_parsing.tsv**), it is often useful to create a token-level long table, where each row corresponds to a single token (word) rather than a whole sentence. This format is more convenient for downstream quantitative analyses (e.g., POS and lemma frequency, dependency patterns, argument structure, and syntactic complexity measures).

Overview workflow:

- **Input: it_vxc_spkr17_with_parsing.tsv**
 - Contains one row per sentence, including:

- sentence (orthographic form)
- morphonological_analysis (the %mor: tier)
- syntactic_parsing (the %gra: tier)
- **Output: it_vxc_spkr17_tokens_level.tsv**
 - Contains one row per token, with the following core columns:
 - sent_id: sentence index (links back to the sentence-level TSV)
 - token_id: token index within the sentence
 - mor: raw %mor token (e.g., noun|libro-Masc)
 - gra: raw %gra token (e.g., 2|4|NSUBJ)
 - upos: part of speech (e.g., noun, verb, det)
 - lemma: dictionary form (e.g., libro, suscitare)
 - feats: all remaining morphological features (e.g., Masc-Plur-Acc)
 - id_gra: token index from the %gra tier
 - head: dependency head index
 - deprel: dependency relation (e.g., NSUBJ, OBJ, OBL)
 - Any additional sentence-level metadata (e.g., speaker ID, condition, task type) can be joined back by sent_id.

How is the token-level file created (conceptual steps)?

- **Assign a sentence ID:**
 - Add a column sent_id to the sentence-level TSV (e.g., 1, 2, 3, ...).
- **Split the %mor and %gra strings into token vectors:**
 - For each sentence, split morphonological_analysis and syntactic_parsing on spaces, yielding:
 - %mor: det|il-Masc-Def-Art-Sing, noun|libro-Masc, aux|avere-Fin-Ind-Pres-S3, ...
 - %gra: 1|2|DET, 2|4|NSUBJ, 3|4|AUX, ...
 - Tokens are aligned by position (1st %mor token with 1st %gra token, etc.).
- **Expand to one row per token:**
 - For each sentence sent_id, create rows:
 - token_id = 1, 2, 3, ...
 - mor = the corresponding %mor token
 - gra = the corresponding %gra token
- **Parse %mor into POS, lemma, and features:**
 - Split each mor token at the first |:
 - upos = part of speech (e.g., noun, verb, det, adp)
 - stem_feats = remaining string (e.g., libro-Masc, suscitare-Part-Past-S)
 - Then split stem_feats at the first -:
 - lemma = dictionary form (e.g., libro, suscitare)

- feats = all remaining morphological features (merged into one string for now).
- **Parse %gra into dependency indices and labels**
 - Split each gra token on |:
 - id_gra = token index (1, 2, 3, ...)
 - head = index of the head token in the sentence
 - deprel = UD dependency relation (e.g., NSUBJ, OBJ, OBL, AUX, ROOT).
- **Optionally join back sentence-level metadata**
 - Using sent_id, merge any additional columns from it_vxc_spkr17_with_parsing.tsv (e.g., speaker, condition, register) into the token-level table.

A minimal R implementation of this pipeline is provided in the accompanying script (see [Cha-to-TSV.R](#) and the token-level extension). The token-level long table is the recommended format for:

- Computing POS, lemma, and feature frequencies
- Examining subject/object/oblique distributions across verbs
- Deriving syntactic complexity measures (e.g., number of clauses, dependency length)
- Conducting mixed-effects modeling and other quantitative analyses at the token or sentence level.

Final sample [TSV file output](#) (token level):

1	1	det-Masc-Def-Art-Sing	det	il	Masc-Def-Art-Sing	1	2	DET	common_voice_i_23000167.m3p	017292020540700d118d7f395a37e9fd703b1768170884409
1	2	noun(lbno)-Masc	noun	lbno	Masc	2	4	NSUBJ	common_voice_i_23000167.m3p	017292020540700d118d7f395a37e9fd703b1768170884409
1	3	aux(ave)-Fin-Ind-Pres-S3	aux	ave	Fin-Ind-Pres-S3	3	4	AUX	common_voice_i_23000167.m3p	017292020540700d118d7f395a37e9fd703b1768170884409
1	4	verb(suscitate)-Part-Part-S	verb	suscitate	Part-Part-S	4	11	ROOT	common_voice_i_23000167.m3p	017292020540700d118d7f395a37e9fd703b1768170884409
1	5	det(molto)-Fem-Def-Ind-Pur	det	molto	Masc-Def-Ind-Pur	5	6	DET	common_voice_i_23000167.m3p	017292020540700d118d7f395a37e9fd703b1768170884409
1	6	noun(polemico)-Fem-Plur-Acc	noun	polemico	Masc-Plur-Acc	6	4	OBL	common_voice_i_23000167.m3p	017292020540700d118d7f395a37e9fd703b1768170884409
1	7	adp(j-	adp	a	NA	7	8	CASE	common_voice_i_23000167.m3p	017292020540700d118d7f395a37e9fd703b1768170884409
1	8	noun(causa)-Fem	noun	causa	Fem	8	4	OBL	common_voice_i_23000167.m3p	017292020540700d118d7f395a37e9fd703b1768170884409
1	9	det(Masc-Def-Ind-Pur)	det	di	Masc-Def-Ind-Pur	9	11	DET	common_voice_i_23000167.m3p	017292020540700d118d7f395a37e9fd703b1768170884409
1	10	det(ju)-Masc-Def-Plur	det	suo	Masc-Def-Plur	10	11	DET-POSS	common_voice_i_23000167.m3p	017292020540700d118d7f395a37e9fd703b1768170884409
1	11	noun(contenuto)-Masc-Plur-Acc	noun	contenuto	Masc-Plur-Acc	11	4	OBJ	common_voice_i_23000167.m3p	017292020540700d118d7f395a37e9fd703b1768170884409
1	12	.	NA	NA	NA	12	4	PUNCT	common_voice_i_23000167.m3p	017292020540700d118d7f395a37e9fd703b1768170884409
2	1	adv(jino)	adv	fino	NA	1	2	CASE	common_voice_i_202404040.m3p	c03ab6962c0920540700d118d7f395a37e9fd7009247ec47ba3d331
2	2	adv(jallinizio)	adv	dalinizio	NA	2	9	ADV/MOD	common_voice_i_202404040.m3p	c03ab6962c0920540700d118d7f395a37e9fd7009247ec47ba3d331
2	3	det(ji)-Fem-Def-Art-Sing	det	il	Masc-Def-Art-Sing	3	4	DET	common_voice_i_202404040.m3p	c03ab6962c0920540700d118d7f395a37e9fd7009247ec47ba3d331
2	4	noun(seco)-Fem	noun	seco	Fem	4	9	NSUBJ	common_voice_i_202404040.m3p	c03ab6962c0920540700d118d7f395a37e9fd7009247ec47ba3d331
2	5	adj(episopcale)-S1	adj	episopcale	S1	5	4	AMOD	common_voice_i_202404040.m3p	c03ab6962c0920540700d118d7f395a37e9fd7009247ec47ba3d331
2	6	aux(ju)-Fin-Ind-Pres-S3	aux	essere	Fin-Ind-Pres-S3	6	9	AUX	common_voice_i_202404040.m3p	c03ab6962c0920540700d118d7f395a37e9fd7009247ec47ba3d331
2	7	aux(ju)-Part-Part-S	aux	essere	Part-Part-S	7	9	COP	common_voice_i_202404040.m3p	c03ab6962c0920540700d118d7f395a37e9fd7009247ec47ba3d331
2	8	adv(immediatamente)	adv	immediatamente	NA	8	9	ADV/MOD	common_voice_i_202404040.m3p	c03ab6962c0920540700d118d7f395a37e9fd7009247ec47ba3d331
2	9	adj(loggetto)-S1	adj	soggetto	S1	9	12	ROOT	common_voice_i_202404040.m3p	c03ab6962c0920540700d118d7f395a37e9fd7009247ec47ba3d331
2	10	det(ju)-Fem-Def-Art-Sing	det	il	Masc-Def-Art-Sing	10	11	CASE	common_voice_i_202404040.m3p	c03ab6962c0920540700d118d7f395a37e9fd7009247ec47ba3d331
2	11	prop(jSanta)	propn	Santa	NA	11	9	OBL	common_voice_i_202404040.m3p	c03ab6962c0920540700d118d7f395a37e9fd7009247ec47ba3d331
2	12	prop(jSe)	propn	Sede	NA	12	11	FLAT-NNAME	common_voice_i_202404040.m3p	c03ab6962c0920540700d118d7f395a37e9fd7009247ec47ba3d331
2	13	.	NA	NA	NA	13	9	PUNCT	common_voice_i_202404040.m3p	c03ab6962c0920540700d118d7f395a37e9fd7009247ec47ba3d331
3	1	aux(ju)-Fin-Ind-Past-S3	aux	essere	Fin-Ind-Past-S3	1	3	COP	common_voice_i_24970935.m3p	057d3d822c12032d763ee7e754668e37053be4600fe625ec
3	2	det(ju)-Def-Art-Sing	det	il	Masc-Def-Art-Sing	2	3	DET	common_voice_i_24970935.m3p	057d3d822c12032d763ee7e754668e37053be4600fe625ec
3	3	noun(fondatore)-Masc	noun	fondatore	Masc	3	10	ROOT	common_voice_i_24970935.m3p	057d3d822c12032d763ee7e754668e37053be4600fe625ec
3	4	adp(ji)	adp	di	NA	4	6	CASE	common_voice_i_24970935.m3p	057d3d822c12032d763ee7e754668e37053be4600fe625ec
3	5	det(molto)-Masc-Def-Ind-Pur	det	molto	Masc-Def-Ind-Pur	5	6	DET	common_voice_i_24970935.m3p	057d3d822c12032d763ee7e754668e37053be4600fe625ec
3	6	noun(chiostro)-Masc-Plur	noun	chiostro	Masc-Plur	6	3	NMOD	common_voice_i_24970935.m3p	057d3d822c12032d763ee7e754668e37053be4600fe625ec
3	7	cm(jm)	cm	cm	NA	7	8	PUNCT	common_voice_i_24970935.m3p	057d3d822c12032d763ee7e754668e37053be4600fe625ec
3	8	noun(ospedale)-Masc-Plur	noun	ospedale	Masc-Plur	8	6	CONJ	common_voice_i_24970935.m3p	057d3d822c12032d763ee7e754668e37053be4600fe625ec
3	9	conj(je)	conj	e	NA	9	10	CC	common_voice_i_24970935.m3p	057d3d822c12032d763ee7e754668e37053be4600fe625ec
3	10	noun(chiesa)-Fem-Plur	noun	chiesa	Fem-Plur	10	6	CONJ	common_voice_i_24970935.m3p	057d3d822c12032d763ee7e754668e37053be4600fe625ec
3	11	.	NA	NA	NA	11	3	PUNCT	common_voice_i_24970935.m3p	057d3d822c12032d763ee7e754668e37053be4600fe625ec

7. Some CLAN Post-processing (optional, if you want to look at a different analysis)

- Once %mor and %gra tiers exist, you can analyze the Italian corpus with CLAN tools:
 - **kwal +t%mor ~/ba_output/troncamento.cha** # view morphology
 - **Command example:**
 - kwal +t%mor /Users/yuechen/ba_output/troncamento.cha
 - **kwal +t%gra +s"OBL" ~/ba_output/troncamento.cha** # search dependencies
 - **freq +t%mor -t* ~/ba_output/troncamento.cha** # POS frequency

8. Citation and Documentation

- Batchalign2: <https://talkbank.github.io/batchalign2/>
- [CLAN Manual](#), section on Batchalign and role conversion
- **Reference article:** [Liu, H. & MacWhinney, B. \(2024\). Morphosyntactic Analysis for CHILDES. Language Development Research, 4\(1\), 233–258.](#)