

# Towards automatically extracting morphosyntactical error patterns from L1-L2 parallel dependency treebanks

18th Workshop on Innovative Use of NLP  
for Building Educational Applications

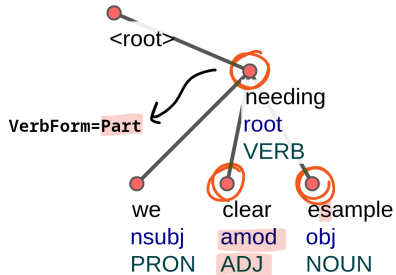
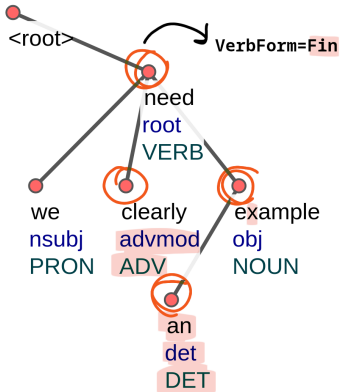
Arianna Masciolini, Elena Volodina and Dana Dannélls

Språkbanken Text, University of Gothenburg



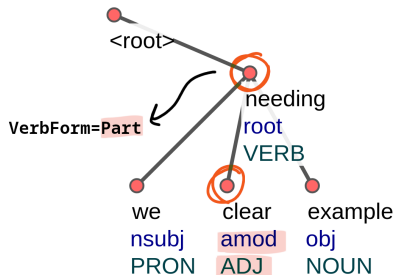
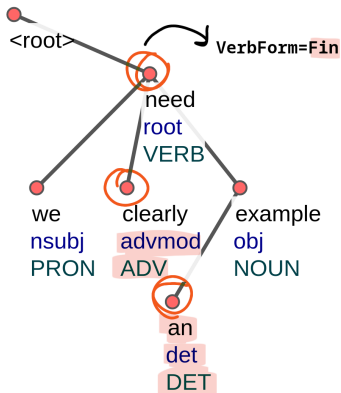
- ❖ learner sentences || correction hypotheses
- ❖ no explicit error labelling,  
just morphosyntactical annotation
- ❖ main design goal: interoperability → UD

# Example



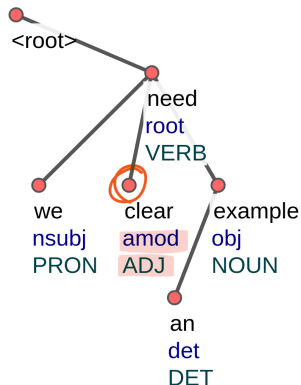
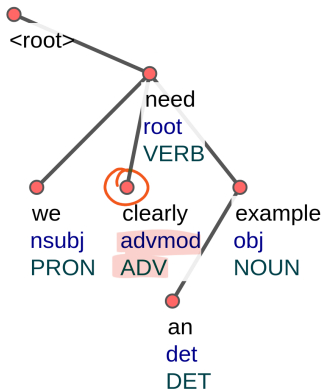
<we clearly needing an example, we clear needing \_ esample>

# Example



*<we clearly needing an example, we clear needing \_ example>*

# Example



{we clearly need an example, we clear need an example}



- ❖ find instances of specific error patterns → L2-UD query engine<sup>1</sup>
- ❖ automatically classify syntactical errors → SErCL<sup>2</sup>

---

<sup>1</sup> Masciolini, 2023

<sup>2</sup> Choshen et al., 2020



- ❖ find instances of specific error patterns → L2-UD query engine<sup>1</sup>
- ❖ automatically classify syntactical errors → SErCL<sup>2</sup>
- ❖ **extract machine-readable error patterns**

---

<sup>1</sup> Masciolini, 2023

<sup>2</sup> Choshen et al., 2020



1. **error detection**: align L1-L2 sentences and filter discrepant alignment
2. **pattern generation**: convert pairs of UD subtrees into machine-readable error patterns



# Step 1: error detection



## concept-alignment

Public

Syntax-based Concept Alignment for Machine Translation

☆ Star



machine-translation

universal-dependencies

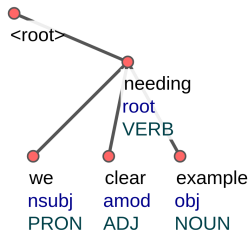
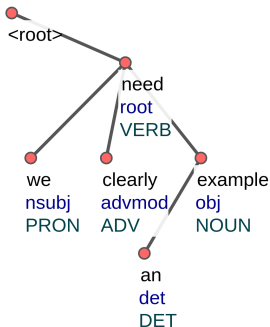
grammatical-framework

● TeX ☆ 1 🍷 1 📄 BSD 2-Clause "Simplified" License Updated on Mar 23

[github.com/harisont/concept-alignment](https://github.com/harisont/concept-alignment)

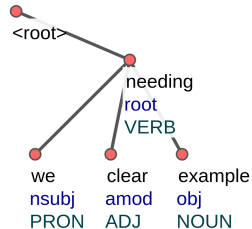
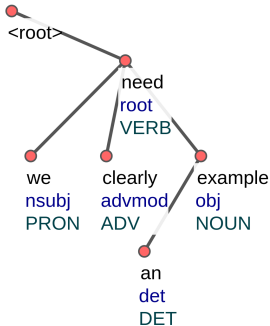
- ❏ extracts **subtree alignments** from parallel UD treebanks
- ❏ syntax-based but language-agnostic
- ❏ designed to generate translation lexica, but easy to adapt to the L1-L2 case

# Alignments



- ❑ ⟨we clearly need an example, we clear needing example⟩, ⟨need, needing⟩
- ❑ ⟨we, we⟩
- ❑ ⟨clearly, clear⟩
- ❑ ⟨an example, example⟩, ⟨example, example⟩

# Errors



- ❑  $\langle \text{we clearly need an example, we clear needing example} \rangle^*$ ,  
 $\langle \text{need, needing} \rangle^*$
- ❑  $\langle \text{we, we} \rangle$
- ❑  $\langle \text{clearly, clear} \rangle^*$
- ❑  $\langle \text{an example, example} \rangle^*$ ,  $\langle \text{example, example} \rangle$

## Step 2: pattern generation

# Query languages for UD trees



- ❖ several options to choose from
  - ❖ PML-TQ, Grew-match, UDAPI...
- ❖ decided on gf-ud's embedded query language
  - ❖ sufficiently expressive and user-friendly
  - ❖ easy to use as a library



---

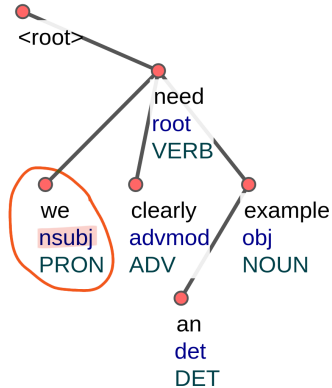
| <b>pattern type</b>            | <b>example</b>                         |
|--------------------------------|--|
| single-token patterns          | DEPREL "nsubj"                         |
| tree patterns                  | TREE (POS "NOUN") [DEPREL "det"]       |
| sequence patterns              | SEQUENCE [DEPREL "advmod", POS "VERB"] |
| logical operators <sup>3</sup> | OR [POS "NOUN", POS "PRON"]            |

---

---

<sup>3</sup> AND, OR, NOT

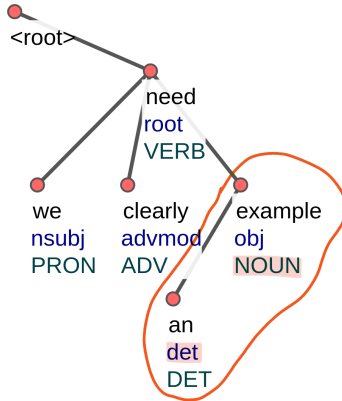
# Single-token patterns



DEPREL "nsubj"

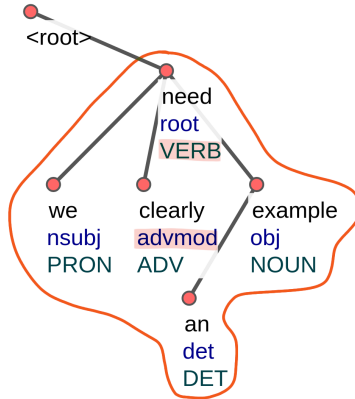


# Tree patterns



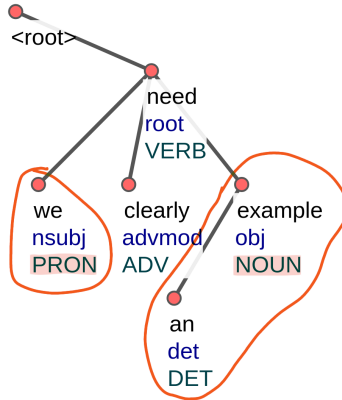
TREE (POS "NOUN") [DEPREL "det"]

# Sequence patterns



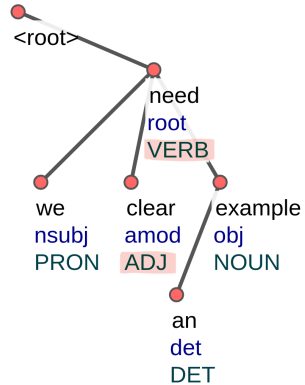
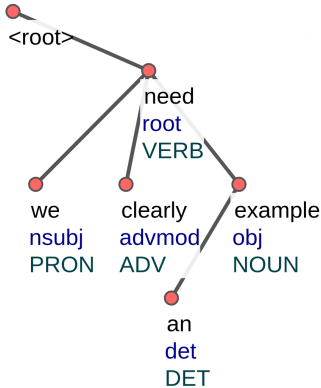
SEQUENCE [DEPREL "advmod", POS "VERB"]

# Logical operators



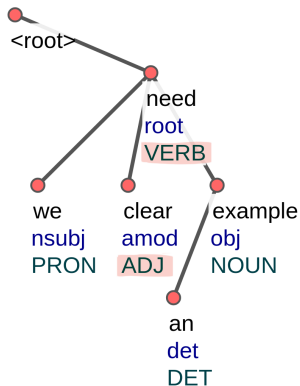
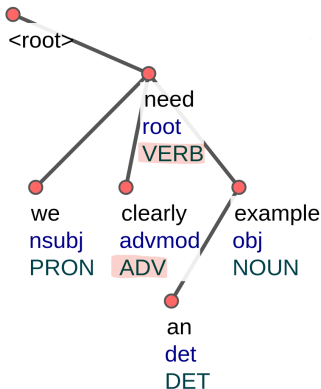
OR [POS "NOUN", POS "PRON"]

# L1-L2 UD patterns



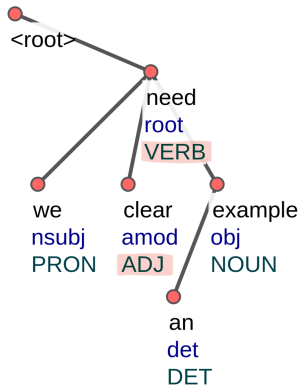
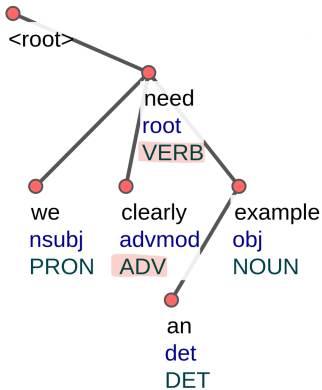
TREE\_ (POS "VERB") [POS "ADJ"]

# L1-L2 UD patterns



`<TREE_ (POS "VERB") [POS "ADV"], TREE_ (POS "VERB") [POS "ADJ"]>`

# L1-L2 UD patterns



TREE\_ (POS "VERB") [POS "{ADV -> ADJ}"]

## 0. Automatically generated L1-L2 pattern

TREE

```
(AND [FORM "need", LEMMA "need", POS "VERB", XPOS "VBP", DEPREL "root",
      FEATS "Mood=Ind|Number=Plur|Person=1|Tense=Pres|VerbForm=Fin"])
[AND [FORM "we", LEMMA "we", POS "PRON", XPOS "PRP", DEPREL "nsubj"
      FEATS "Case=Nom|Number=Plur|Person=1|PronType=Prs"],
     AND [FORM {"clearly" -> "clear"}, LEMMA {"clearly" -> "clear"},
          POS {"ADV" -> "ADJ"}, XPOS {"RB" -> "JJ"},
          FEATS {"_" -> "Degree=Pos"}, DEPREL {"advmod" -> "amod"}],
```

TREE

```
(AND [FORM "example", LEMMA "example", POS "NOUN",
      XPOS "NN", FEATS "Number=Sing", DEPREL "obj"])
[AND [FORM "an", LEMMA "a", POS "DET", XPOS "DT",
      FEATS "Definite=Ind|PronType=Art", DEPREL "det"]
```

]

]

## 1. Filtering by CoNLL-U field

Ignoring FORM, LEMMA, XPOS and DEPREL:

TREE

```
(AND [POS "VERB",  
      FEATS "Mood=Ind|Number=Plur|Person=1|Tense=Pres|VerbForm=Fin"])  
[AND [POS "PRON", FEATS "Case=Nom|Number=Plur|Person=1|PronType=Prs"],  
     AND [POS {"ADV" -> "ADJ"}, FEATS {"_" -> "Degree=Pos"}]],  
TREE  
  (AND [POS "NOUN", FEATS "Number=Sing"])  
  [AND [POS "DET", FEATS "Definite=Ind|PronType=Art"]]  
]
```



## 2. Removal of never-discrepant fields

In all alignments, FEATS is either identical both in the L1 and in the L2 or absent in one of the components:

TREE

```
(AND [POS "VERB"])  
[AND [POS "PRON"],  
  AND [POS {"ADV" -> "ADJ"}],  
  TREE (AND [POS "NOUN"]) [AND [POS "DET"]]  
]
```



## 3. Elimination of identical subpatterns

Removing identical subtrees<sup>4</sup>:

```
TREE_  
  (AND [POS "VERB"])  
  [AND [POS {"ADV" -> "ADJ"}]]
```

---

<sup>4</sup> optionally, identical roots can be replaced with the wildcard pattern TRUE too



## 4. Monolingual single-pattern simplifications

AND [p] is equivalent to p:

```
TREE_ (POS "VERB") [POS {"ADV" -> "ADJ"}]
```

# Preliminary evaluation

- ❖ evaluation through a **similar example retrieval task**:
  1. **extract** L1-L2 patterns from an error-correction input pair
  2. **query** an L1-L2 treebank with the extracted patterns
- ❖ interactive version available as prototype CALL application

## 2 datasets for linguistic acceptability judgments:

- ❑ only one error per sentence
- ❑ filtered: only morphosyntactical errors
- ❑ automatically parsed with UDPipe 2

| <b>name</b> | <b>language</b> | <b>size<sup>5</sup></b> | <b>description</b>                 |
|-------------|-----------------|-------------------------|------------------------------------|
| BLiMP       | English         | 14 996                  | artificially generated sentences   |
| DaLAJ       | Swedish         | 1 198                   | postprocessed L2 learner sentences |

---

<sup>5</sup> post-filtering

|         | <b>BLiMP</b> | <b>DaLAJ</b> |
|---------|--------------|--------------|
| $R^6$   | 82%          | 69%          |
| $R_+^7$ | 82%          | 63%          |

---

<sup>6</sup> retrieval rate

<sup>7</sup> successful retrieval rate



## L2-UD Public

Tools for working with UD treebanks of learner texts.

☆ Star

● Haskell  MIT License Updated on Apr 30

`github.com/harisont/l2-ud`

- ❏ novel approach to error pattern extraction
- ❏ preliminary, bilingual evaluation on LA datasets giving promising results
- ❏ interactive similar example retrieval pipeline available as prototype CALL application



- ❖ extraction method:
  - ❖ handle nonexistent word forms
  - ❖ deal with real-world L2 data:
    - non-morphosyntactical errors (spelling, lexical...)
    - multiple overlapping errors
- ❖ example retrieval application:
  - ❖ implement pattern selection/ranking
  - ❖ build UI
- ❖ use L1-L2 patterns from feedback comment generation
- ❖ improve automatic annotation of L2 sentences

# Thank you!

asynchronous questions/comments: [arianna.masciolini@gu.se](mailto:arianna.masciolini@gu.se)

- ❖ Leshem Choshen, Dmitry Nikolaev, Yevgeni Berzak, and Omri Abend. *Classifying syntactic errors in learner language*. In Proceedings of the 24th Conference on Computational Natural Language Learning, pages 97–107, Online, 2020. Association for Computational Linguistics
- ❖ Prasanth Kolachina and Aarnte Ranta. *From abstract syntax to Universal Dependencies*. Linguistic Issues in Language Technology, 13, 2016
- ❖ John Lee, Keying Li, and Herman Leung. *L1-L2 parallel dependency treebank as learner corpus*. In Proceedings of the 15th International Conference on Parsing Technologies, pages 44–49, Pisa, Italy, 2017. Association for Computational Linguistics

- ❖ Arianna Masciolini. *A query engine for L1-L2 parallel dependency treebanks*. In Proceedings of the 24th Nordic Conference on Computational Linguistics (NoDaLiDa), pages 574–587, Tórshavn, Faroe Islands, 2023. University of Tartu Library
- ❖ Arianna Masciolini and Aarne Ranta. *Grammar-based concept alignment for domain-specific Machine Translation*. In Proceedings of the Seventh International Workshop on Controlled Natural Language (CNL 2020/21), Amsterdam, Netherlands, 2021. Special Interest Group on Controlled Natural Language
- ❖ Aarne Ranta and Prasanth Kolachina. *From Universal Dependencies to abstract syntax*. In Proceedings of the NoDaLiDa 2017 Workshop on Universal Dependencies (UDW 2017), pages 107–116, Gothenburg, Sweden, 2017. Association for Computational Linguistics



- ❖ Elena Volodina, Yousuf Ali Mohammed, and Julia Klezl. *DaLAJ-a dataset for linguistic acceptability judgments for Swedish: Format, baseline, sharing*. arXiv preprint arXiv:2105.06681, 2021
- ❖ Alex Warstadt, Alicia Parrish, Haokun Liu, Anhad Mohananey, Wei Peng, Sheng-Fu Wang, and Samuel R. Bowman. *BLiMP: The benchmark of linguistic minimal pairs for English*. Transactions of the Association for Computational Linguistics, 8:377–392, 2020