

ud2rrg: Transforming Universal Dependencies into RRG Trees

Kilian Evang Simon Petitjean

2020-02-26

TreeGraSP Meeting #4

Outline

Introduction

The basic algorithm

Clause linkage

RRGbank results

ud2rrg for RRGparbank

RRGparbank results

Outline

Introduction

The basic algorithm

Clause linkage

RRGbank results

ud2rrg for RRGparbank

RRGparbank results

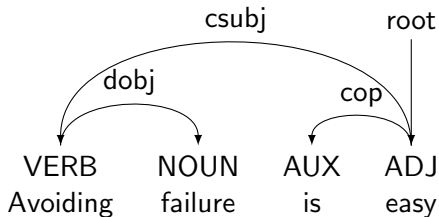
What are Universal Dependencies?

- framework for consistent annotation of parts of speech, syntactic dependencies across languages
- treebanks for 70+ languages

“Manning’s law”

1. UD needs to be satisfactory on linguistic analysis grounds for individual languages.
2. UD needs to be good for linguistic typology, i.e., providing a suitable basis for bringing out cross-linguistic parallelism across languages and language families.
3. UD must be suitable for rapid, consistent annotation by a human annotator.
4. UD must be suitable for computer parsing with high accuracy.
5. UD must be easily comprehended and used by a non-linguist, whether a language learner or an engineer with prosaic needs for language processing. We refer to this as seeking a habitable design, and it leads us to favor traditional grammar notions and terminology.
6. UD must support well downstream language understanding tasks (relation extraction, reading comprehension, machine translation, ...).

UD: example



Why convert UD to RRG?

- a single algorithm for a large number of languages, treebanks
- similar in philosophy to RRG
- minimal theoretical assumptions, good starting point

Outline

Introduction

The basic algorithm

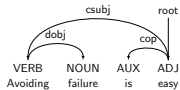
Clause linkage

RRGbank results

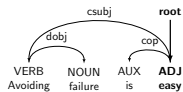
ud2rrg for RRGparbank

RRGparbank results

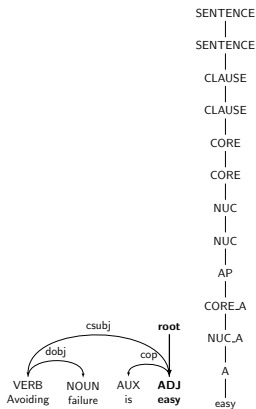
Example



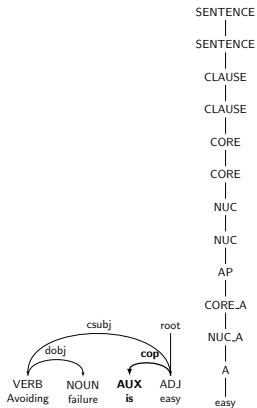
Example



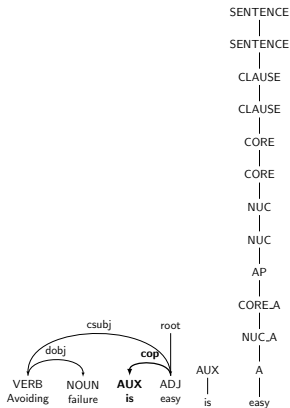
Example



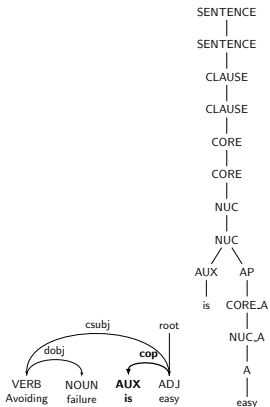
Example



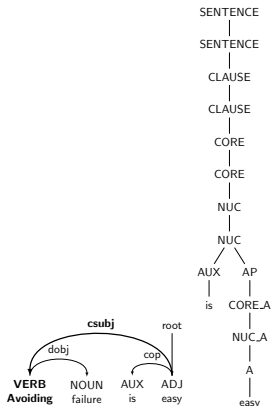
Example



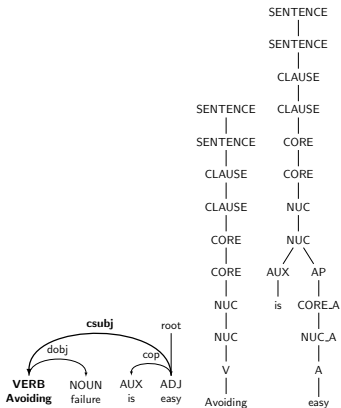
Example



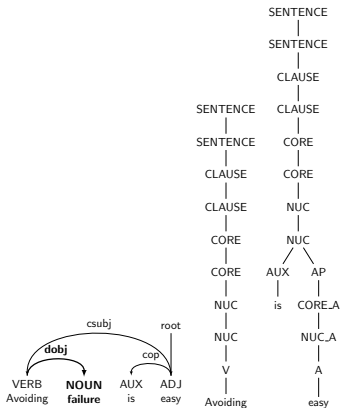
Example



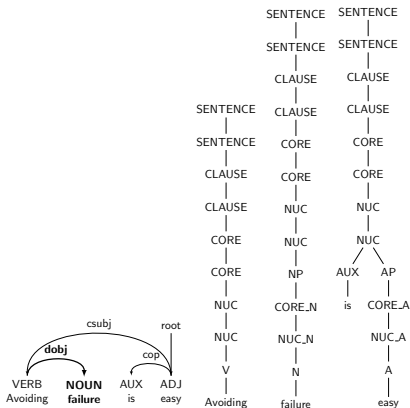
Example



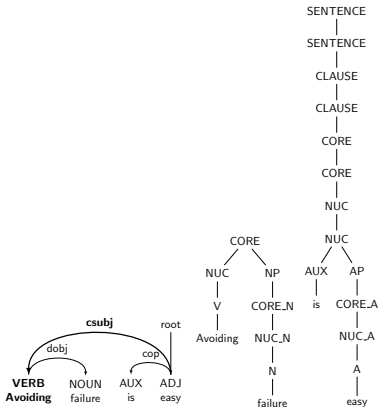
Example



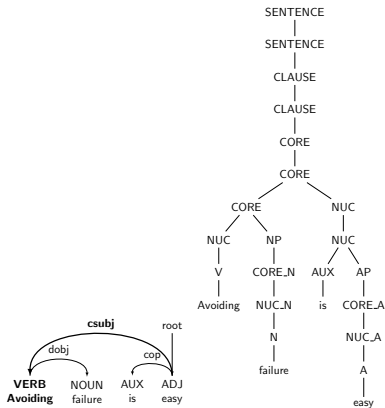
Example



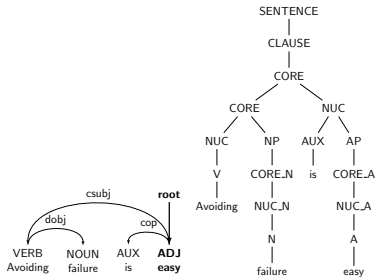
Example



Example



Example



The basic algorithm: overview

To convert a UD (sub)tree to RRG,

1. based on part of speech and lemma of root, choose an RRG *fragment*
2. for each child,
 - 2.1 convert it recursively,
 - 2.2 based on dependency relation, link the result into the current tree:
 - for csubj, nsubj, dobj, ...: simple attachment (trim to appropriate layer and attach at appropriate node),
 - for xcomp, ccomp, ...: clause linkage ({clause,core,nuclear} {cosubordination,subordination}),
3. conflate redundant nodes.

Outline

Introduction

The basic algorithm

Clause linkage

RRGbank results

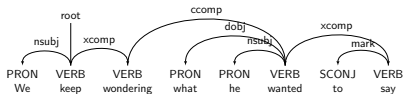
ud2rrg for RRGparbank

RRGparbank results

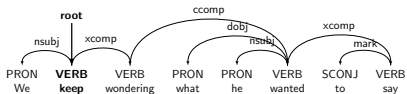
Clause linkage rules so far

- for xcomp dependencies:
 - phase verb (*begin, keep, stop...*) heading participle → nuclear cosubordination
 - otherwise → core cosubordination
- for ccomp dependencies:
 - discourse verb (*say, think, wonder...*) → clause subordination
 - otherwise → core subordination
- implemented for English by Naima Grebe using VerbNet

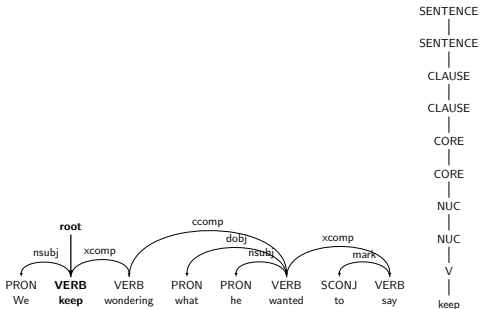
Example



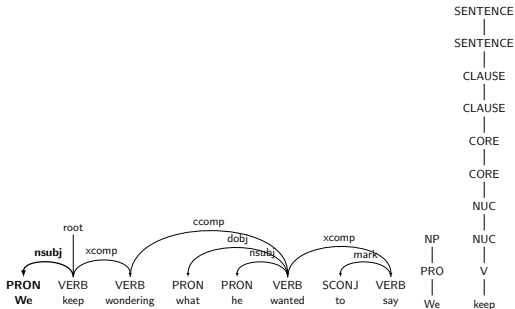
Example



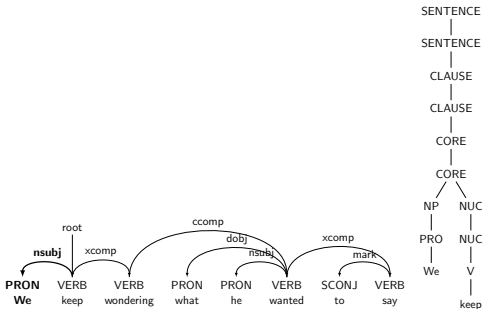
Example



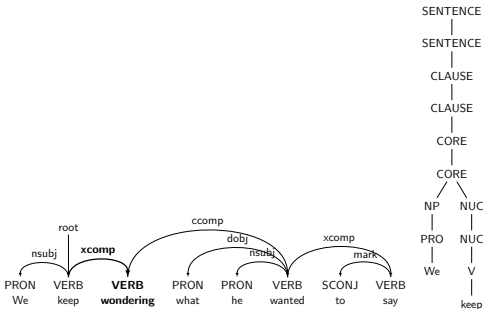
Example



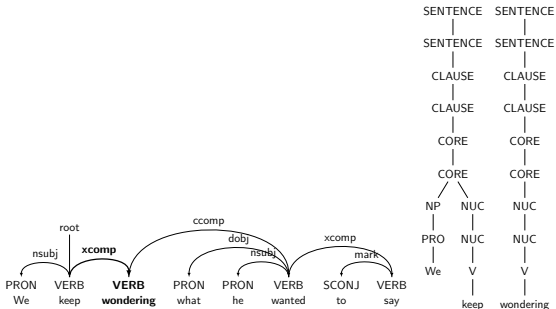
Example



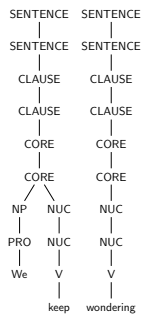
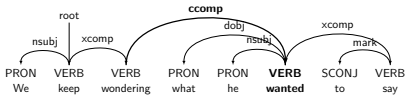
Example



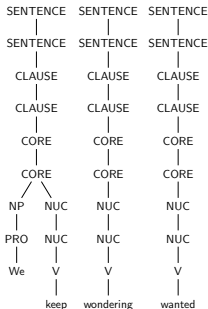
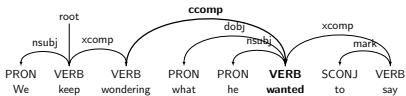
Example



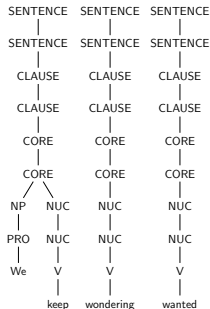
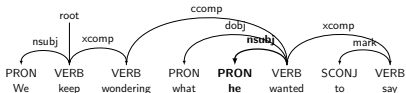
Example



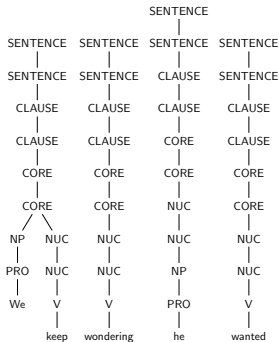
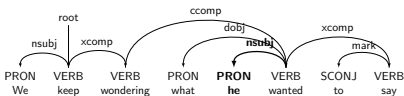
Example



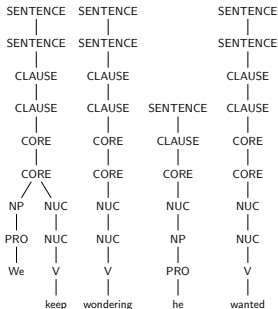
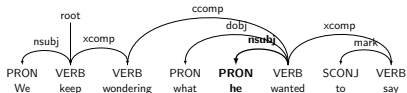
Example



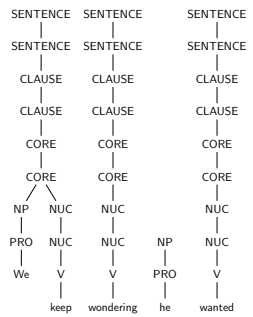
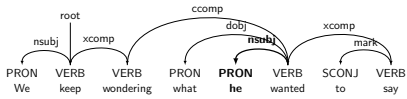
Example



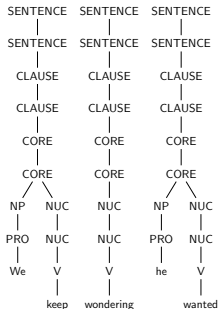
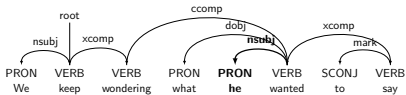
Example



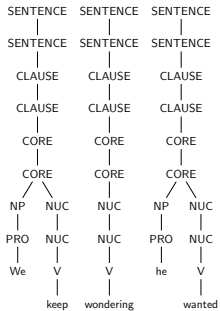
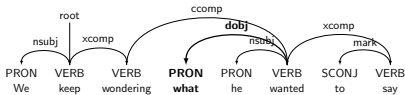
Example



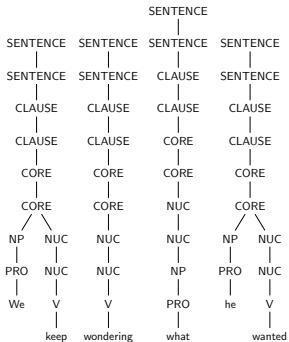
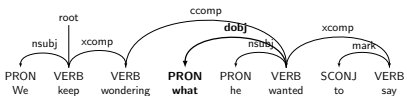
Example



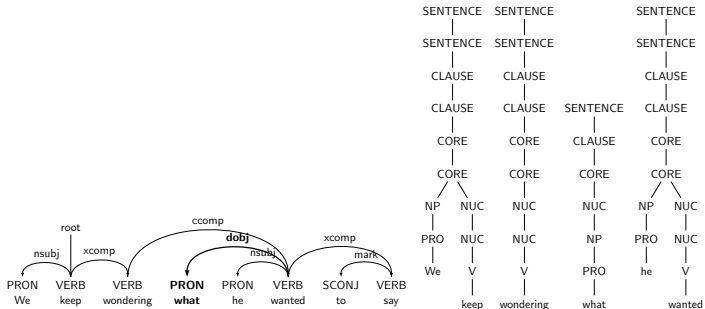
Example



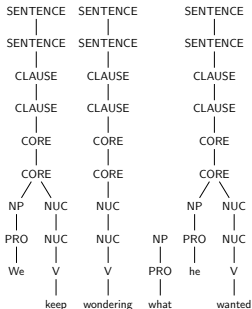
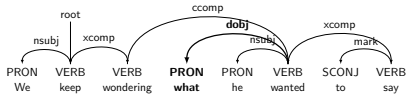
Example



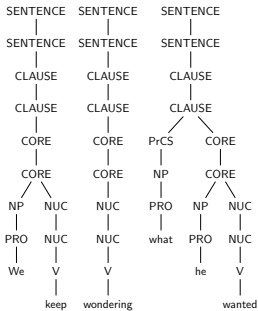
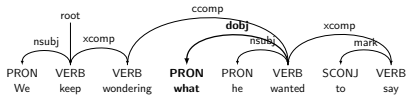
Example



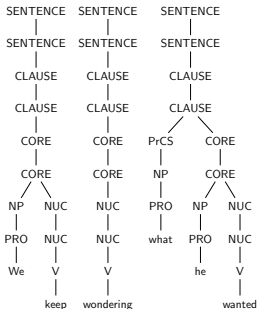
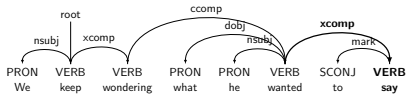
Example



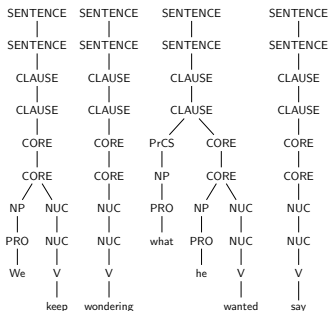
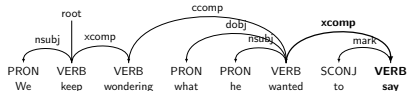
Example



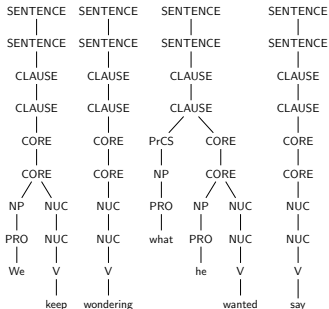
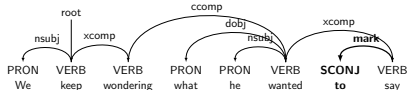
Example



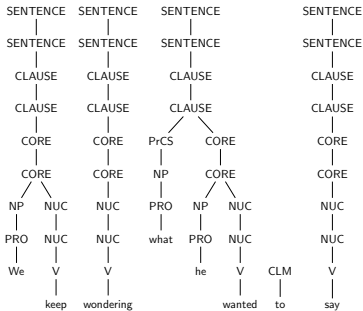
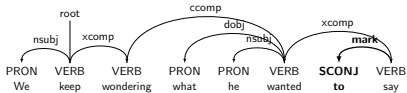
Example



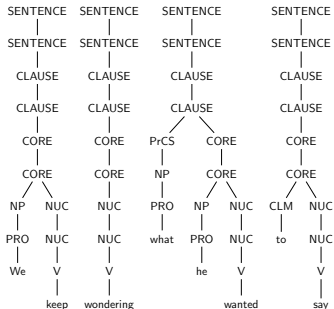
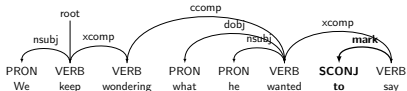
Example



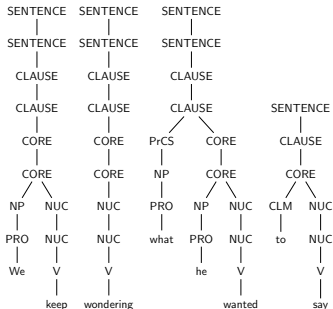
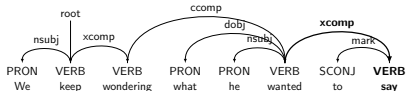
Example



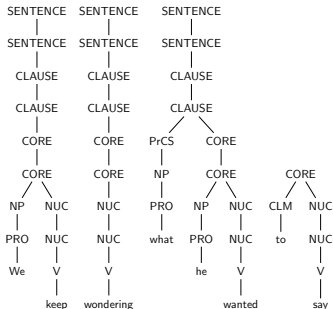
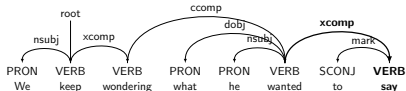
Example



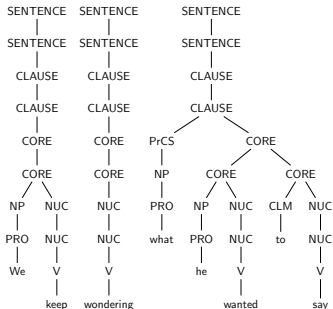
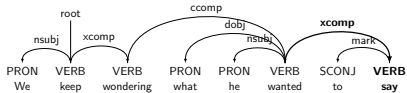
Example



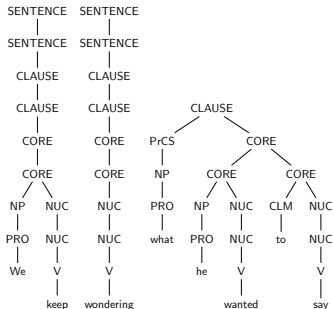
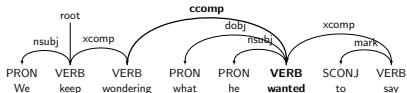
Example



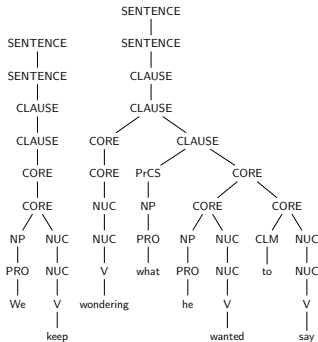
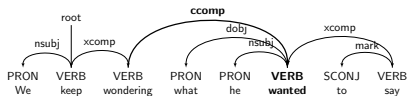
Example



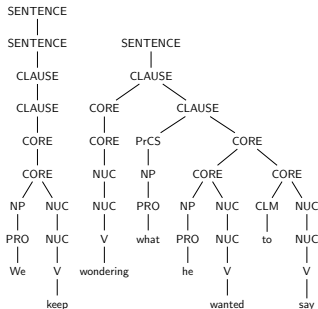
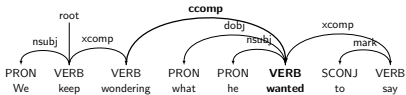
Example



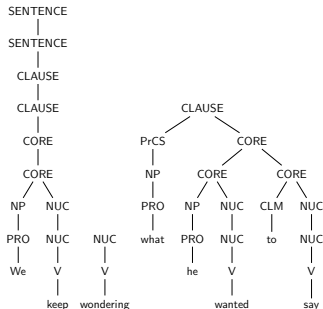
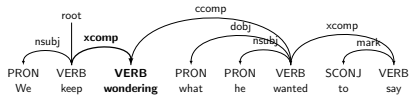
Example



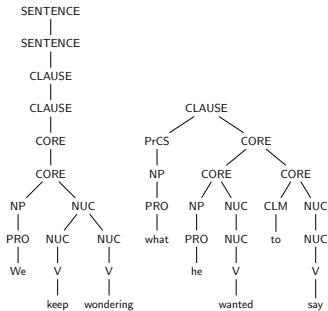
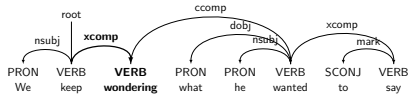
Example



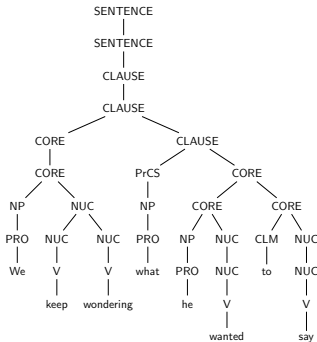
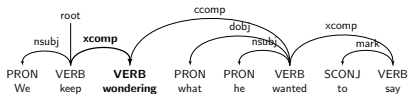
Example



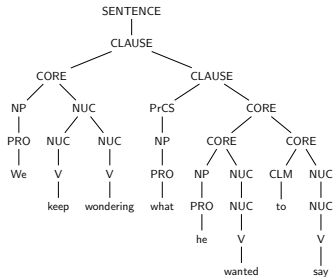
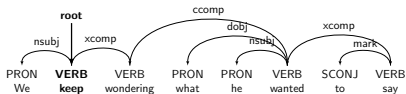
Example



Example



Example



Outline

Introduction

The basic algorithm

Clause linkage

RRGbank results

ud2rrg for RRGparbank

RRGparbank results

Current results

RRGbank development data

language	English
source	Penn Treebank WSJ
sentences	2097
avg. sentence length	8.98

ptb2rrg vs. ud2rrg

	ptb2rrg.py	ud2rrg.py
coverage	100%	96.28%+
EVAlB f-score	92.74	85.99%

Outline

Introduction

The basic algorithm

Clause linkage

RRGbank results

ud2rrg for RRGparbank

RRGparbank results

ud2rrg for RRGparbank

- RRGparbank: annotate multilingual data
- Multext-East “1984” Corpus: English + translations into languages of eastern Europe
- Gold tokens + POS tags + features + alignments
- No gold constituency trees + language specific transformation rules → ptb2rrg hard to adapt
- Adapt ud2rrg for different kinds of UD inputs

New challenges

- Targeted languages: English, German, Hungarian, Russian, Persian
- Input: dependencies generated with udpipes
- Sometimes only raw text (Russian)
- Output of udpipes: specific POS tags, more/less precise dependency labels
- Unequal quality of parses

English data

- Similar to the input of ud2rrg for rrgbank
- No gold constituency parse → less information (marked as argument, PTB POS tags, etc)

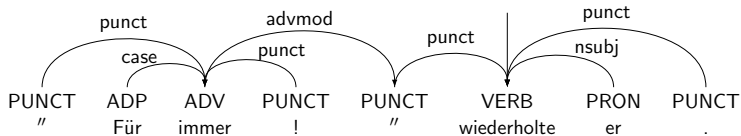
Russian data

- no gold tokens / POS tags: use raw data

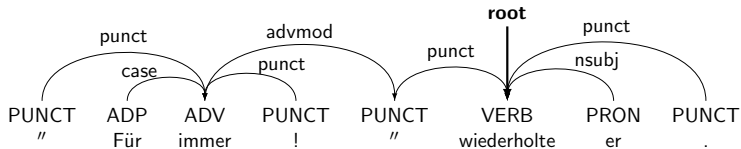
German data

- Problems: more resource related than language related
- Low quality of dependency parsing
- Improve udpipe parsing: use subset of features
- Improve ud2rrg robustness: reassign dependents / update root

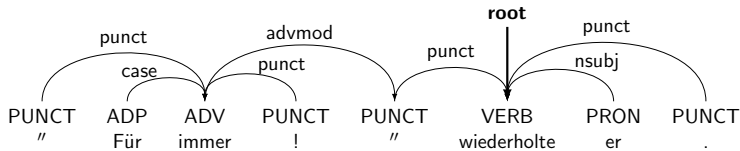
Reassigning dependents: example



Reassigning dependents: example

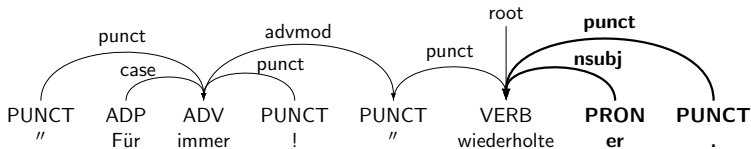


Reassigning dependents: example



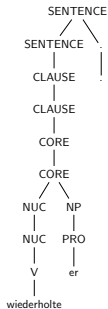
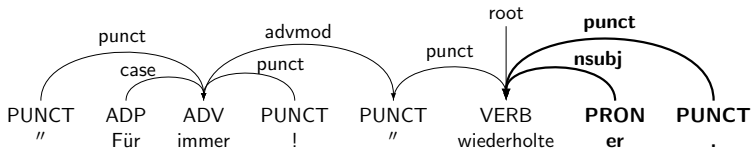
SENTENCE
|
SENTENCE
|
CLAUSE
|
CLAUSE
|
CORE
|
CORE
|
NUC
|
NUC
|
V
|
wiederholte

Reassigning dependents: example

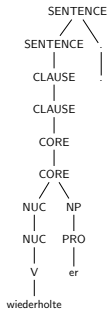
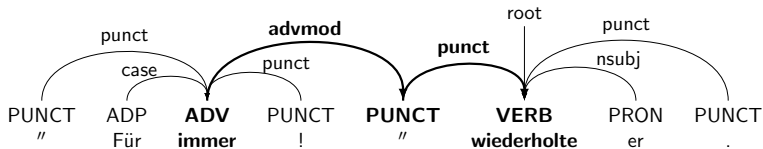


SENTENCE
|
SENTENCE
|
CLAUSE
|
CLAUSE
|
CORE
|
CORE
|
NUC
|
NUC
|
V
|
wiederholte

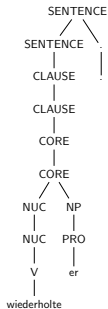
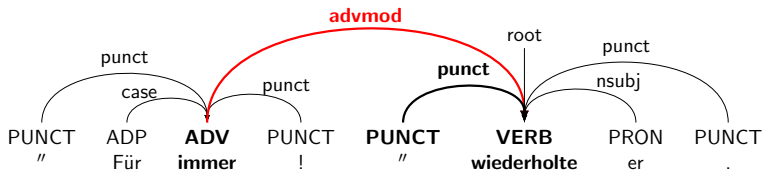
Reassigning dependents: example



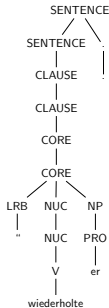
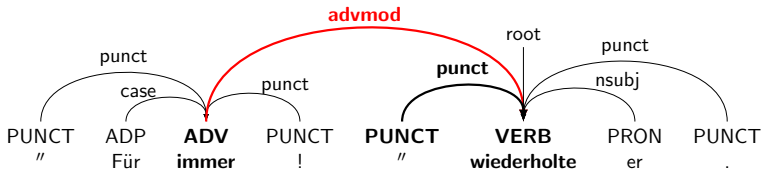
Reassigning dependents: example



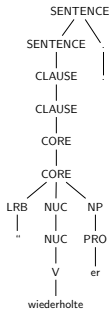
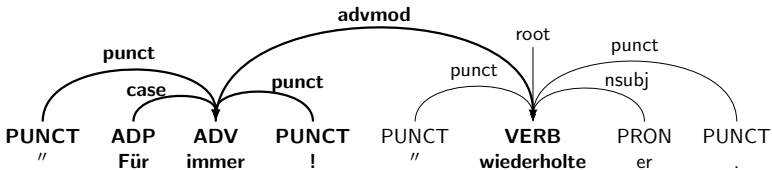
Reassigning dependents: example



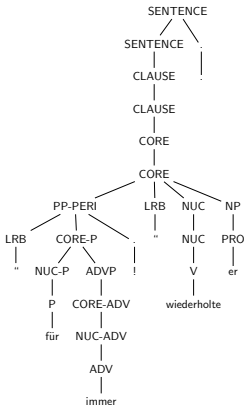
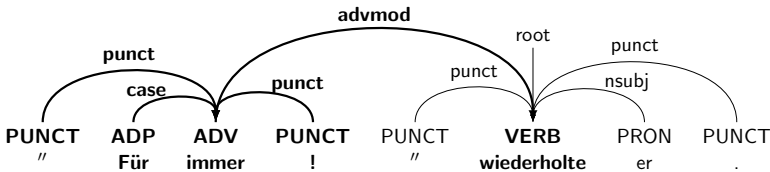
Reassigning dependents: example



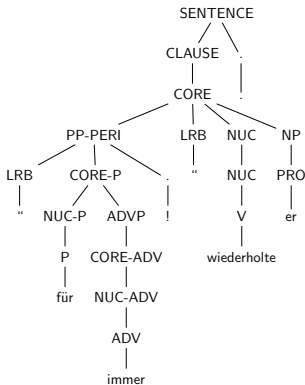
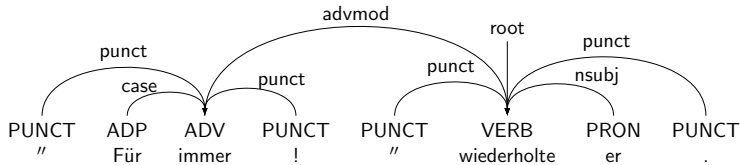
Reassigning dependents: example



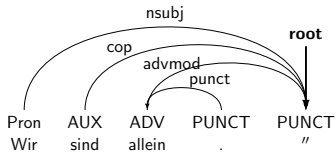
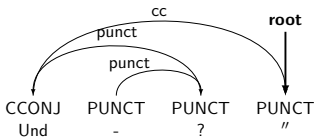
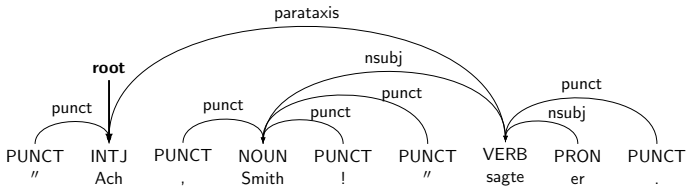
Reassigning dependents: example



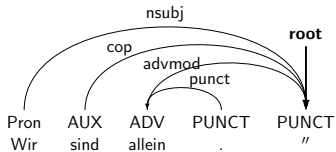
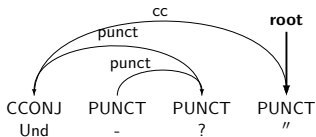
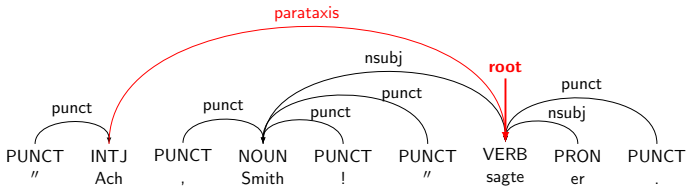
Reassigning dependents: example



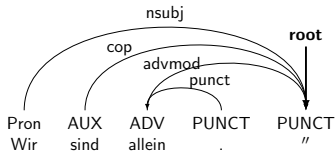
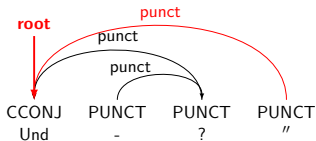
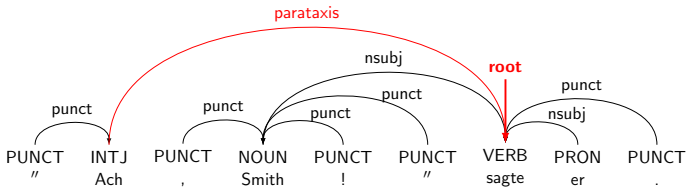
Reassigning roots: examples



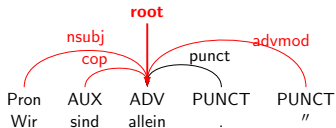
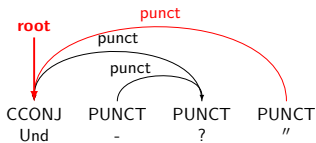
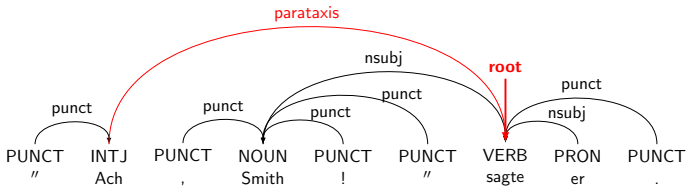
Reassigning roots: examples



Reassigning roots: examples



Reassigning roots: examples



Outline

Introduction

The basic algorithm

Clause linkage

RRGbank results

ud2rrg for RRGparbank

RRGparbank results

Current results

ud2rrg for 3 languages

	German (gold+silver)	English (silver)	Russian (gold+silver)
EVALB f-score	76.49%	70.62%	84.52%

Bibliography I