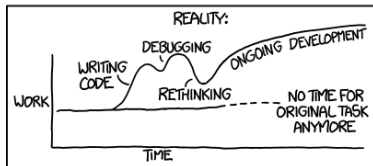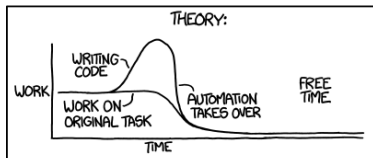# Udapi: Universal API for Universal Dependencies

**Martin Popel**, Zdeněk Žabokrtský, Martin Vojtek

Charles University, Faculty of Mathematics and Physics, Prague, Czechia
popel@ufal.mff.cuni.cz
Workshop on Universal Dependencies (UDW 2017), Gothenburg, 2017-05-22

https://xkcd.com/1319/

# Udapi – http://udapi.github.io/

- API and open-source framework for processing UD
- Python, Perl, Java
- Allows both fast prototyping and full applications
- Both command-line tool (udapy) and library
- Modularity, reusability, cooperation
- Based on 20-year experience with dep. treebanking, TrEd (tree editor), Treex/TectoMT (predecessor of Udapi)

## Parsing using UDPipe

```
echo "John loves Mary." | udapy \
    read.Sentences \
    tokenize.Simple \
    udpipe.Base model_alias=en tokenize=0 \
    write.Conllu
```

Output:

```
# sent_id = 1
# text = John loves Mary.
1 John   John PROPN NNP Number=Sing  2 _ _ _
2 loves  love VERB  VBZ Mood=Ind|... 0 _ _ _
3 Mary   Mary PROPN NNP Number=Sing  2 _ _ SpaceAfter=No
4 .      .    PUNCT .   _             2 _ _ _
```

# Parsing using UDPipe

```
echo "John loves Mary." | udapy \
    read.Sentences \
    tokenize.Simple \
    udpipe.Base model_alias=en tokenize=0 \
    write.Conllu
```

- Python command-line interface (called udapy)
- 4 processing units (called *blocks*)
- blocks may have parameters

# Parsing using UDPipe

```
echo "John loves Mary." | udapy \
    read.Sentences \
    tokenize.Simple \
    udpipe.Base model_alias=en tokenize=0 \
    write.Conllu
```

- Python command-line interface (called udapy)
- 4 processing units (called *blocks*)
- blocks may have parameters

# Parsing using UDPipe

```
echo "John loves Mary." | udapy \
    read.Sentences \
    tokenize.Simple \
    udpipe.Base model_alias=en tokenize=0 \
    write.Conllu
```

- Python command-line interface (called udapy)
- 4 processing units (called *blocks*)
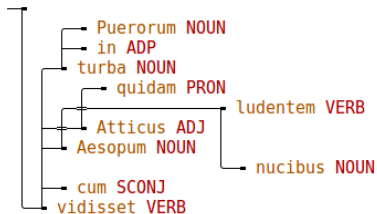- blocks may have parameters

# Parsing using UDPipe

```
echo "John loves Mary." | udapy \
    read.Sentences \
    tokenize.Simple \
    udpipe.Base model_alias=en tokenize=0 \
    write.Conllu
```

Shortcut:

```
echo "John loves Mary." | udapy -s \
    read.Sentences \
    udpipe.En
```

## Visualization – text-mode trees

```
cat latin-sample.conllu | udapy \
 write.TextModeTrees attributes=form,upos
```



```
udapy -T < latin-sample.conllu | less -R
```

# Visualization – text-mode trees in HTML

```
udapy -H < en-bugs.conllu > en-train-bugs.html
```



```
← → C ⌂ | ⓘ file:///en-train-bugs.html

bugs = ud.MarkBugs Error Overview:
        appos-chain        1
          det-upos         1
        punct-alpha        1
        punct-child        2
        punct-upos         2
         multi-obj        28
            TOTAL          35

docname = weblog-juancole.com_juancole_20051126063000_ENG_20051126_063000
# sent_id = weblog-typepad.com_ripples_20040407125600_ENG_20040407_125600-0062
# text = Take care, my friend, Linda

┌─ Take VERB root Bug=multi-obj
│ ┌─ care NOUN obj SpaceAfter=No
│ ├─ , PUNCT punct _
│ │ ┌─ my PRON nmod:poss _
│ ├─ friend NOUN obj SpaceAfter=No
│ ├─ , PUNCT punct _
│ └─ Linda PROPN vocative _

# sent_id = weblog-juancole.com_juancole_20030911085700_ENG_20030911_085700-0014
# text = On the one hand, it should pressure Musharraf to take off his uniform and
                                    ┌─ On ADP case _
                                    ├─ the DET det _
                                    ├─ one NUM nummod _
                                    ├─ hand NOUN obl SpaceAfter=No
        ┌─ , PUNCT punct _
        ├─ it PRON obj _
        ├─ should AUX aux _
        ├─ pressure VERB root Bug=multi-obj
        ├─ Musharraf PROPN obj _
                                ┌─ to PART mark _
                                ├─ take VERB xcomp _
                                ├─ off ADP compound:prt _
```

# Visualization – traditional-style trees in HTML

```
udapy write.Html < czeng.conllu > czeng.html
```

| Previous | 1 | 2 | 3 | 4 | 5 | 6 | 7 |
|          | 8 | 9 | 10 | 11 | 12 | ... |

Next    Save as SVG

[cs] Kde jsem.
[en] Where I am.

zone=cs
id=f000001-s1/cs

jsem
VERB
root

Kde          .
ADV       PUNCT
advmod     punct

zone=en
id=f000001-s1/en

am
VERB
root

Where    I     .
ADV    PRON  PUNCT
dep    nsubj  punct

lemma=be
Tense=Pres
VerbForm=Fin
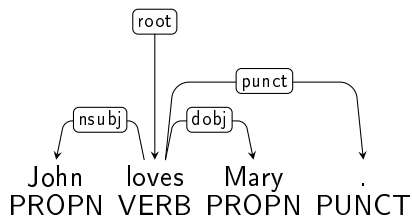
# Visualization – TikZ & LATEX

```
udapy write.Tikz < john.conllu > john.tex
```

```
\begin{dependency}
\begin{deptext}
% sent_id = 1
% text = John loves Mary.
John  \& loves \& Mary \& .    \\
PROPN \& VERB  \& PROPN \& PUNCT \\
\end{deptext}
\depedge{2}{1}{nsubj}
\deproot{2}{root}
\depedge{2}{3}{dobj}
\depedge{2}{4}{punct}
\end{dependency}
```

# Format conversions

- plain text (one sentence per line)
- CoNLL-U and other CoNLL-like formats
- SDParse (used in Stanford Dependencies & Brat)
- VISL-cg
- easy to implement other readers and writers

```
udapy write.Vislcg < x.conllu > x.vislcg
udapy read.Vislcg write.Sdparse \
      < x.vislcg > x.sdparse
```

## Querying

Udapi: (queries specified in Python)

```
cat in.conllu | udapy -T \
  util.Filter \
    mark=nonproj \
    keep_tree_if_node='node.is_nonprojective()'
```

```
cat in.conllu | udapy -TM \
  util.Mark node='node.is_nonprojective()'
```

Alternatives: (queries in special declarative languages)

- PML-TQ (Prague)
- SETS (Turku)

## Editing

Ad-hoc edits, e.g. delete the subtypes of dependency relations
(*acl:relcl* → *acl*, . . . )

```
cat in.conllu | udapy -s \
  util.Eval node='node.deprel = node.udeprel' \
  > out.conllu
```

For better reusability & maintainability use separate Python files,
e.g. udapi/block/transform/flatten.py
will be available via udapy as transform.Flatten:

```
from udapi.core.block import Block

class Flatten(Block):

    def process_node(self, node):
        node.parent = node.root
        node.deprel = 'root'
```
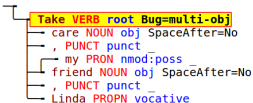
## Validation

```
udapy -HAM ud.MarkBugs skip=no-NumType \
    < en-ud-train.conllu > en-train-bugs.html
```
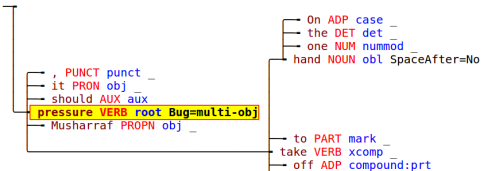
# Conversions

UDv1 to UDv2

```
udapy -s ud.Convert1to2 < in.conllu > out.conllu
```

- unsure edits marked with `ToDo` in MISC
- used for 5 UDv2 treebanks

"Google pre-UDv1" to UDv2

```
udapy -s ud.Google2ud < in.conllu > out.conllu
```

- used for 11 PUD treebanks (+ id,ko,th not released)

## Solving text vs. token mismatches

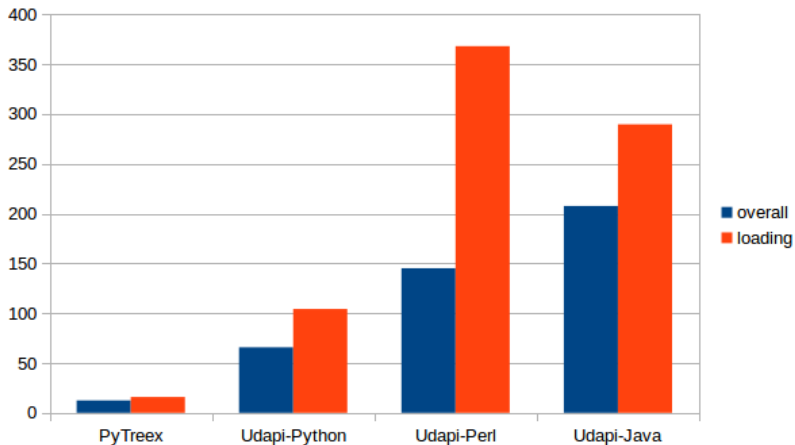Raw sentences should match the tree tokens and SpaceAfter=No.

- `ud.SetSpaceAfter` – heuristic rules for SpaceAfter=No
- `ud.SetSpaceAfterFromText` – uses the raw text
- `ud.ComplyWithText` – heuristic alignment, add MWT, add "goeswith" nodes, revert form normalization (e.g. ``TeX-like quotes'', missing thousand separators, . . . )

## Other use cases

- `util.Wc` – count of words, empty words, sents, . . .
- `util.See` – advanced statistics of nodes matching a condition
- `eval.Parsing` – UAS, LAS, LAS (udeprel only)
- `eval.F1` – Precision/Recall/F1 of various attributes
- `transform.Proj`, `transform.Deproj` – (de)projectivization
- `ud.xy.AddMwt` – split multi-word tokens into words in lang. xy
- `ud.FixPunct` – (re)attach punctuation
- `ud.FixChain`, `ud.FixRightHeaded`, . . .
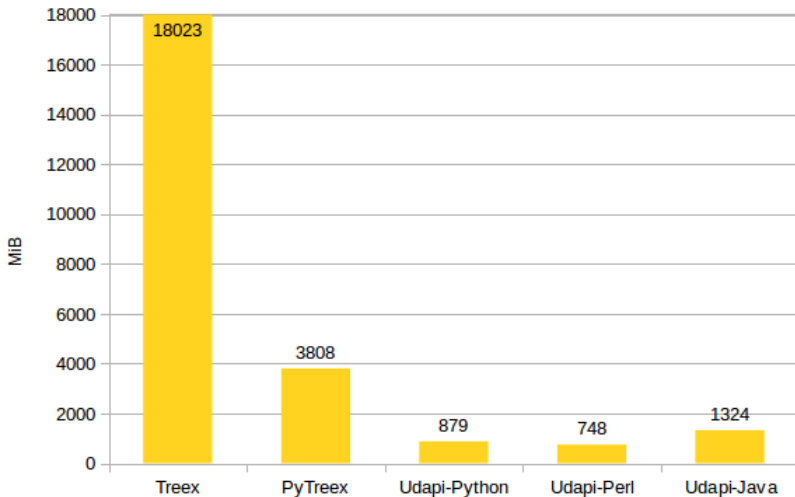- `util.MarkDiff` – diff two (CoNLL-U) files

```
udapy -HMAC \
 read.Conllu zone=old files=a.conllu \
 read.Conllu zone=new files=b.conllu \
 util.MarkDiff gold_zone=old > diff.html
```

# Benchmark: Speed-up relative to Treex



(source: https://github.com/martinpopel/newtreex)

# Benchmark: Memory (MB)



(source: https://github.com/martinpopel/newtreex)
cs-ud-train-l.conllu: 68 MiB, 41k sentences, 0.8 MWords

## Algorithmic challenges

- data structure for globally-ordered rooted trees
  `node.descendants ... ordered`
  `node.shift_before_node(another_node)`
- efficient loading&saving of CoNLL-U files
  linear-time checking of cycles
  lazy deserialization of FEATS and MISC
- `write.TextModeTrees` for non-projective trees
  minimize crossings and/or depth

## Udapi web

# http://udapi.github.io
provides links to

- Hands-on tutorial
- GitHub repo for Python, Perl, Java
- documentation
- these slides + the paper